

Materials Submitted to the National Research Council Part I: Status of Implementation of Recommendations

U.S. Environmental Protection Agency Integrated Risk Information System Program

January 30, 2013

DISCLAIMER

This document is for review purposes only. It has not been formally disseminated by EPA. It does not represent and should not be construed to represent any Agency determination or policy.

Table of Contents

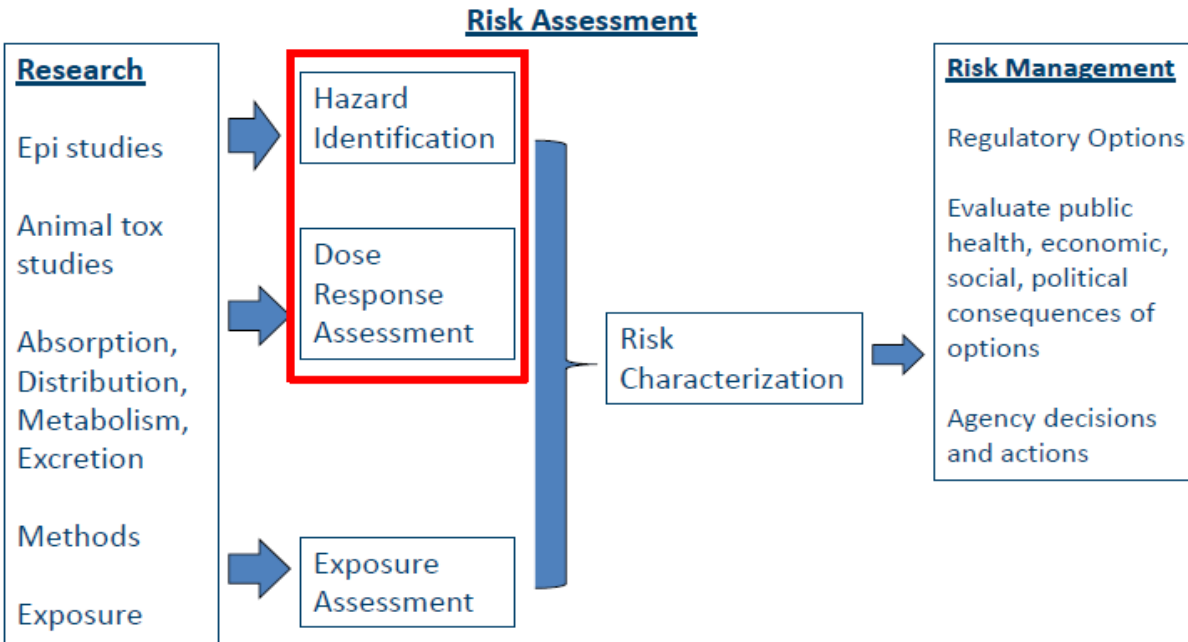
I. Introduction.....	2
II. Charge to the NRC Expert Panel.....	3
III. Overview of EPA’s Implementation of NRC’s Recommendations.....	3
IV. Additional Initiatives.....	18
V. Summary.....	19
Appendix A – IRIS Toxicological Review Template.....	A-1
Appendix B – Preamble to IRIS Toxicological Reviews.....	B-1
Appendix C – Example of IRIS Program Direction to Contractors.....	C-1
Appendix D – Information Management Tool: Comment Tracker Database	D-1
Appendix E – Scoping to Inform the Development of IRIS Assessments.....	E-1
Appendix F – Draft Handbook for IRIS Assessment Development	F-1

I. Introduction

The U.S. Environmental Protection Agency's (EPA) Integrated Risk Information System (IRIS) Program develops human health assessments that provide health effects information on environmental chemicals to which the public may be exposed, providing a critical part of the scientific foundation for EPA's decisions to protect public health. In April 2011, the National Research Council (NRC), in their report *Review of the Environmental Protection Agency's Draft IRIS Assessment of Formaldehyde*, made several recommendations to EPA for improving IRIS assessments and the IRIS Program. The NRC's recommendations were focused on Step 1 of the IRIS process, the development of draft assessments. Consistent with the advice of the NRC, the IRIS Program is implementing these recommendations using a phased approach and is making the most extensive changes to assessments that are in the earlier stages of the IRIS process.

Background on IRIS

IRIS human health assessments contain information that can be used to support the first two steps (hazard identification and dose-response analysis) of the risk assessment paradigm. IRIS assessments are scientific reports that provide information on a chemical's hazards and, when supported by available data, quantitative toxicity values for cancer and noncancer health effects. IRIS assessments are not regulations, but they provide a critical part of the scientific foundation for decisions to protect public health across EPA's programs and regions under an array of environmental laws (e.g., Clean Air Act, Safe Drinking Water Act, Comprehensive Environmental Response, Compensation, and Liability Act, etc). EPA's program and regional offices combine IRIS assessments with specific exposure information for a chemical. This information is used by EPA, together with other considerations (e.g., statutory and legal requirements, cost/benefit information, technological feasibility, and economic factors), to characterize the public health risks of environmental chemical and make risk management decisions, including regulations, to protect public health. IRIS assessments are also a resource for risk assessors and environmental and health professionals from state and local governments and other countries. Figure 1 illustrates where IRIS assessments contribute information within the risk assessment and risk management paradigms.



¹ Adapted from the National Research Council risk assessment risk management paradigm (NRC 1983).

Figure 1. Risk Assessment Risk Management Paradigm (adapted from the National Research Council’s paradigm, 1983). The red box shows the information included in IRIS assessments.

II. Charge to the NRC Expert Panel

In April 2012, EPA contracted with the NRC to conduct a comprehensive review of the IRIS assessment development process. The panel will review the IRIS process and the changes being made or planned by EPA and will recommend modifications or additional changes as appropriate to improve the process, and scientific and technical performance of the IRIS Program. The panel will focus on the development of IRIS assessments rather than the review process that follows draft development. In addition, the panel will review current methods for evidence-based reviews and recommend approaches for weighing scientific evidence for chemical hazard and dose-response assessments.

III. Overview of EPA’s Implementation of NRC’s Recommendations

EPA agrees with the NRC’s 2011 recommendations for the development of IRIS assessments and plans to fully implement the recommendations consistent with the NRC panel’s “Roadmap for Revision,” which viewed the full implementation of their recommendations by the IRIS Program as a multi-year process. In response to the NRC’s 2011 recommendations, the IRIS Program has made changes to streamline the assessment development process, improve transparency, and create efficiencies within the Program. The following sections outline the NRC’s 2011 recommendations and provide an overview of how the IRIS Program is implementing the NRC’s general and specific

recommendations. Further details regarding changes that have been made and will be made in response to the recommendations are provided in Appendices to this report.

In addition, chemical-specific examples demonstrating how the IRIS Program is currently implementing the NRC's 2011 recommendations have also been provided to the panel (see additional document provided, *Chemical-Specific Examples Demonstrating Implementation of NRC's 2011 Recommendations*). The examples cover literature search and screening, evaluation and display of individual studies, development of evidence tables, evidence integration, selecting studies for derivation of toxicity values, dose-response modeling output, and considerations for selecting organ/system-specific or overall toxicity values. The examples are not to be construed as final Agency conclusions and are provided for the sole purpose of demonstrating how the IRIS Program is implementing the NRC recommendations.

NRC's General Recommendations and Guidance

NRC Recommendations¹:

- To enhance the clarity of the document, the draft IRIS assessment needs rigorous editing to reduce the volume of text substantially and address redundancies and inconsistencies. Long descriptions of particular studies should be replaced with informative evidence tables. When study details are appropriate, they could be provided in appendices.
- Chapter 1 needs to be expanded to describe more fully the methods of the assessment, including a description of search strategies used to identify studies with the exclusion and inclusion criteria articulated and a better description of the outcomes of the searches and clear descriptions of the weight-of-evidence approaches used for the various noncancer outcomes. The committee emphasizes that it is not recommending the addition of long descriptions of EPA guidelines to the introduction, but rather clear concise statements of criteria used to exclude, include, and advance studies for derivation of the RfCs and unit risk estimates.
- Elaborate an overall, documented, and quality-controlled process for IRIS assessments.
- Ensure standardization of review and evaluation approaches among contributors and teams of contributors; for example, include standard approaches for reviews of various types of studies to ensure uniformity.
- Assess disciplinary structure of teams needed to conduct the assessments.

Implementation:

➤ New Document Structure

Implemented

In their report, the NRC recommended that the IRIS Program enhance the clarity of the document, reduce the volume of text, and address redundancies and inconsistencies. To improve the clarity of IRIS assessments, the IRIS Program has revised the assessment template to substantially reduce the volume of text and address redundancies and inconsistencies in assessments. The new template provides sections for the literature search strategy, study selection and evaluation, and methods used to develop the assessment.

¹ National Research Council, 2011. Review of the Environmental Protection Agency's Draft IRIS Assessment of Formaldehyde.

The new document structure includes an *Executive Summary* in the beginning of each assessment which provides a concise summary of the major conclusions of the assessment. Additionally, a newly developed *Preamble* describes the methods used to develop the assessment. Each assessment will include information on the literature search strategy used to identify the evidence for consideration in developing the assessment, as well as the evaluation criteria and rationale used to make decisions about including or excluding studies in the assessment.

The main body of the IRIS assessment has been reorganized into two sections, *Hazard Identification* and *Dose-Response Analysis*, to better focus on the role of IRIS assessments in the risk assessment paradigm and to further reduce the volume of text and redundancies/inconsistencies. Information on assessments by other national and international health agencies, chemical and physical properties, toxicokinetics, and individual studies has been moved to appendices (which are provided as supplemental information) to improve the flow of the document.

In the *Hazard Identification* chapter of the new document template, the IRIS Program has developed subsections based on organ/system-specific hazards to systematically integrate the available evidence for a given chemical (i.e., epidemiology, toxicological, and mechanistic data). The assessment now uses evidence tables to present the key study findings that support how toxicological hazards are identified. In addition, exposure-response arrays are being used as visual tools to inform the hazard characterization. This chapter provides for a strengthened and more integrated and transparent discussion of the weight of the available evidence supporting hazard identification. The IRIS Program is also developing standardized study summary tables, which will be included in the supplemental information, to present more detailed study characteristics and summary information.

The *Dose-Response Analysis* section of the new document structure provides a section to explain the rationale used to select and advance studies for consideration in calculating toxicity values based on conclusions regarding the potential hazards associated with chemical exposure. Key data supporting the dose-response analysis are reported and the methodology and derivation of toxicity values are described. In addition, details of the dose-response analysis—including the data, models, methods, and software—are provided as supplemental information and described in sufficient detail to allow for independent replication and verification. The *Dose-Response Analysis* section also includes tables and figures showing candidate toxicity values for comparison across studies and endpoints. Finally, this section of the new document structure includes clear documentation of the conclusions and selection of the overall toxicity values.



The IRIS assessment template demonstrating the new document structure is provided in Appendix A.

➤ **IRIS Assessment Preamble**

Implemented

In their report, the NRC recommended that the IRIS Program expand Chapter 1 of IRIS assessments to “describe more fully the methods of the assessment, including a description of search strategies used to identify studies with the exclusion and inclusion criteria clearly articulated and a better

description of the outcomes of the searches and clear descriptions of the weight-of-evidence approaches used for the various noncancer outcomes.”

In accordance with this recommendation, the IRIS Program has replaced the previous Chapter 1 of IRIS assessments with a section titled *Preamble to IRIS Toxicological Reviews* which describes the application of existing EPA guidance and the methods and criteria used in developing the assessments. The term “*Preamble*” is used to emphasize that these methods and criteria are being applied consistently across IRIS assessments. The new *Preamble* discusses the following topics:

- Scope of the IRIS Program;
- Process for developing and peer-reviewing IRIS assessments;
- Identifying and selecting pertinent studies;
- Evaluating the quality of individual studies;
- Evaluating the overall evidence of each effect;
- Selecting studies for derivation of toxicity values; and
- Deriving toxicity values.

For each of these topics, the *Preamble* summarizes and cites EPA guidance on methods used in the assessment. The *Preamble* was included in the draft IRIS assessments of ammonia and trimethylbenzenes when they were released for public comment in June 2012 and will be included in all new IRIS assessments.



The Preamble to IRIS Toxicological Reviews is included in Appendix B.

➤ **New Initiatives to Improve Overall Process, Quality Control, and Documentation**

In Progress

In their report, the NRC recommended that the IRIS Program “elaborate an overall, documented, and quality-controlled process for IRIS assessment” and “assess disciplinary structure of teams needed to conduct the assessments.” In response to these recommendations, the IRIS Program has developed several new initiatives and enhanced existing processes. These initiatives help to ensure that standardized approaches are used throughout IRIS assessments and major science decisions are rigorously vetted.

IRIS assessments are developed by interdisciplinary teams of scientists (referred to as an “Assessment Team”) internal to EPA. For each assessment, scientists with the necessary scientific backgrounds (e.g., neurotoxicology, epidemiology, developmental toxicology) are assigned to lead or assist in the development of the assessment. The expertise needed is chemical-specific and the personnel assigned to the assessment team are identified in the early stages of planning and document development.

Contractors may provide technical and analytical support to the chemical assessment teams during the development of assessments. This assistance may be provided in conducting literature searches and identifying pertinent studies; developing evidence tables and exposure-response arrays using studies identified and evaluated by the IRIS Program; and performing dose-response modeling (i.e., using EPA's benchmark dose modeling software [BMDS]). All materials provided by the contractor are evaluated in accordance with EPA policies regarding quality assurance and quality management, and specified in the contract. Contractor products are not incorporated into IRIS assessments without significant Agency scientist review. EPA is responsible for the content and conclusions within the assessments and all scientific and policy decisions are made by the Agency.



An example of instruction provided to contractors is available in Appendix C.

Additionally, discipline-specific workgroups within the IRIS Program assist the assessment teams in developing assessments. These workgroups coordinate across assessments to ensure consistency, solve cross-cutting issues, and advance scientific understanding that contributes to decision-making in IRIS assessments. The discipline-specific workgroups cover topics related to: statistics and dose-response analysis, physiologically-based pharmacokinetic modeling, and mechanistic data.

In late 2011, the IRIS Program developed a new initiative, Chemical Assessment Support Teams (CASTs), as a means of formalizing an internal process to provide continuing quality control in the development of IRIS assessments. This initiative uses a team approach to make judicious, consistent decisions during assessment development, to ensure that the necessary disciplinary expertise is available for assessment development and review, and to provide a forum for identifying and addressing key issues at each stage of the assessment. There are three CASTs and each team consists of four permanent core members: two senior scientists, a senior statistician, and a rapporteur (a staff scientist).

All on-going IRIS assessments have been distributed and assigned across the three CASTs. Each CAST meets periodically with the individual chemical assessment teams. In addition to meeting with each chemical assessment team, the CASTs convene as a group once a week to discuss issues that have surfaced in the chemical-specific CAST meetings from the previous week. Discussions at this meeting are relayed to scientists working on IRIS assessments during weekly meetings convened by the IRIS Program Director.

The CAST initiative:

- Provides a forum for problem solving;
- Ensures appropriate disciplinary structure of assessment teams;
- Pinpoints key issues early on in the assessment;
- Identifies overarching assessment issues that require Program-wide discussions;
- Increases objectivity in assessment decisions;
- Monitors progress in implementing NRC's 2011 recommendations;

- Assists in responding to Agency, interagency, external peer review, and public comments;
- Ensures consistency across assessments; and
- Serves as a mechanism for documenting and communicating decisions.

As noted above, the CASTs ensure documentation and communication of decisions. Documenting discussions and decisions from CAST meetings is the primary responsibility of the rapporteurs, who have developed a searchable database to capture comments received throughout assessment development and review as well as Agency decisions in response to these comments. This important information management tool, the *Comment Tracker Database*, allows for recording, reviewing, responding to, and analyzing comments and responses. The IRIS Program is currently testing the use of this database.

The CAST initiative is aimed at improving the quality and consistency of IRIS assessments as well as identifying overarching scientific issues to be addressed. This process facilitates communication across the organization and consistency across assessments to improve the overall efficiency of the IRIS Program.



The Comment Tracker Database is further described in Appendix D.

The IRIS Program also recognizes that it is important to understand the big picture in order to develop an assessment that is most informative and efficient for decision-makers. Having a clear understanding of the overarching environmental problems being addressed in the context of a chemical can help inform what an IRIS assessment will ultimately include. This concept was supported by the NRC in their 2009 report *Science and Decisions: Advancing Risk Assessment* when they recommended that EPA provide “greater attention on design in the formative stages of risk assessment.” While the NRC was referring to the overall risk assessment paradigm, the spirit of the recommendation supports a scoping step before developing a hazard identification and dose-response assessment (i.e., IRIS assessment). Because of the importance of considering the scope of an IRIS assessment, the IRIS Program is developing a new initiative to include a “scoping” process as an early step in developing IRIS assessments. The scoping process involves consultation with clients in EPA’s program and regional offices. This early consultation provides an opportunity to identify key questions for framing various analyses and helps ensure that the assessment meets the needs and critical timelines of Agency decision-makers.



The considerations for scoping during the development of IRIS assessments are further described in Appendix E.

The IRIS Program has recently initiated ways to improve stakeholder engagement to help ensure transparency and the use of the best available science in IRIS assessments. When IRIS toxicity values are combined with specific exposure information, government and private entities can use these values to help characterize the public health risks of chemical substances in various situations and support risk management decisions to protect public health. Environmental protection decisions can have potentially large impacts on the environment, human health, and the economy. Engaging with stakeholders can help facilitate the development of assessments and promote public

discussion of key scientific issues. Therefore, stakeholder and public scientific engagement is an important part of supporting the best decisions possible.

The IRIS Program considers a stakeholder to be any individual or group that participates in, has an impact on, or could be affected by products produced by the IRIS Program. Public and stakeholder engagement has always been an important part of the IRIS assessment development process. The May 2009 IRIS process provides multiple opportunities for engagement including: (1) public and stakeholder nomination of chemicals for assessment; (2) a public listening session for each draft assessment; (3) public review and comment of draft documents; (4) a public peer review process; and (5) two opportunities for review and comment on draft assessments by other EPA scientists, other Federal agencies, and the Executive Office of the President.

Recently, the IRIS Program convened two public meetings to engage with stakeholders. In November 2012, a public stakeholder meeting was held to discuss the IRIS Program in general. The meeting was intended to begin a series of dialogues between the IRIS Program and a broad and diverse group of stakeholders. The goals of the meeting were to: engage stakeholders in the IRIS process; listen to views and needs of IRIS users in an open and respectful environment; facilitate improvements to the IRIS process; and initiate an ongoing dialogue between the IRIS Program and stakeholders. In January 2013, a public stakeholder meeting, which focused on informing the plan for drafting a new IRIS assessment for inorganic arsenic, was convened. The meeting provided an opportunity for stakeholders to comment on their expectations for the IRIS assessment, the current state of scientific information that should be considered when developing the assessment, and the potential impacts of the completed assessment.

Another initiative involves the increased use of public peer consultation workshops to enhance the input from the scientific community as assessments are designed. Information regarding specific peer consultation workshops will be announced to the public in advance of the meetings. The goal of these workshops will vary. For example, the workshops may focus on the state-of-the-science for a particular chemical or provide a forum for discussion with experts about certain cross-cutting scientific issues that may impact the development of a scientifically complex assessment. One of the first of these peer consultation workshops will focus on mouse lung tumors as they relate to human cancer risk. This is an important issue for the IRIS assessments for naphthalene, styrene, and ethylbenzene.

The IRIS Program will also conduct public dialogue meetings to discuss the available chemical-specific data and the science issues for new IRIS assessments in the draft development stage. IRIS will share with the public the list of references and tables summarizing the key studies prior to the meeting.

NRC's Specific Recommendations and Guidance

The NRC made twenty-five specific recommendations in five broad categories:

- evidence identification,
- evidence evaluation,
- weight-of-evidence evaluation,
- selection of studies for derivation of toxicity values, and
- calculation of toxicity values

The IRIS Program has been working to improve the approaches for identifying and selecting pertinent studies; evaluating and displaying individual studies; strengthening and improving integration of evidence for hazard identification; and increasing transparency in dose-response analysis.

The IRIS Program recognized the value of providing specific information to its assessment teams and contractors in order to develop IRIS assessments that satisfy the needs of the NRC recommendations. In order to document these individual changes, the IRIS Program has compiled information into a working draft *Handbook for IRIS Assessment Development*. This document is intended to more clearly summarize the internal processes and evaluation steps used to develop IRIS assessments. The draft *Handbook* (which in its current form will be made publicly available) is a work in progress and currently does not fully discuss each step in the IRIS assessment development process. However, the draft *Handbook* contains important information that reflects the changes that have been implemented or will be implemented in response to the NRC recommendations. These changes are noted below and further described in the draft *Handbook for IRIS Assessment Development* in Appendix F.

Evidence Identification: Literature Collection and Collation Phase

NRC Recommendations:

- Select outcomes on the basis of available evidence and understanding of mode of action.
- Establish standard protocols for evidence identification.
- Develop a template for description of the search approach.
- Use a database, such as the Health and Environmental Research Online (HERO) database, to capture study information and relevant quantitative data.

Implementation:

➤ **Identifying and Selecting Pertinent Studies**

In Progress

The IRIS Program is adopting the principles of systematic review in IRIS human health assessments with regard to providing an overview of methods and points to consider in the process of developing and documenting decisions. The focus of IRIS assessments is typically on the evidence of health effects (any kind of health effects) of a particular chemical. This is, by definition, a broad

topic. The systematic review process that has been developed and applied within the clinical medicine arena (evidence-based medicine) is generally applied to narrower, more focused questions. Nonetheless, the experiences within the clinical medicine field provide a strong foundation to draw upon. The IRIS Program is planning to convene a workshop in spring 2013 on this topic in order to have a public discussion of systematic review approaches that may be applicable to IRIS assessments.

An IRIS assessment is made up of multiple systematic reviews. The initial steps of the systematic review process formulate specific strategies to identify and select studies relating to each key question, evaluate study methods based on clearly defined criteria, and transparently document the process and its outcomes. Synthesizing and integrating data also falls under the purview of systematic review. Overall, this is an iterative process that identifies relevant scientific information needed to address key, assessment-specific questions.

The IRIS Program has improved the approach to identifying and selecting studies pertinent to IRIS assessments by adopting the principles of systematic review. One of the strengths of systematic review is its ability to identify relevant studies, published and unpublished, pertaining to the question of interest (e.g., what are the health effects of a chemical?). Additionally, by transparently presenting all decision points and the rationale for each decision, bias in study selection and evaluation is eliminated.

The new IRIS assessment document structure includes a detailed description of the literature search strategy and study selection process used to develop IRIS assessments. This section describes how the scientific literature was gathered, emphasizes how studies were selected to be included in the document, and, if applicable, explains the rationale for excluding potentially relevant studies from the assessment. This section of the new document structure is specific to each chemical assessment. It is designed to provide enough information that an independent literature search would be able to replicate the results of the literature search used by the IRIS Program in developing the assessment. In this section, a link to an external database (www.epa.gov/hero) that contains the references that were cited in the document, along with those that were considered for inclusion in the assessment but not cited is provided.



For more detailed information, see the “Identifying and Selecting Pertinent Studies: Literature Search and Screening” section in the draft Handbook for IRIS Assessment Development in Appendix F.



See also Section 3 (“Identifying and selecting pertinent studies”) in the Preamble to IRIS Toxicological Reviews in Appendix B.



A chemical-specific example of the implementation of this recommendation is available as “EXAMPLE 1 – Literature Search and Screening” in the Chemical-specific Examples Demonstrating Implementation of NRC Recommendations document.

Evidence Evaluation: Hazard Identification

NRC Recommendations:

- All critical studies need to be thoroughly evaluated with standardized approaches that are clearly formulated and based on the type of research, for example, observational epidemiologic or animal bioassays. The findings of the reviews might be presented in tables to ensure transparency.
- Standardize the presentation of reviewed studies in tabular or graphic form to capture the key dimensions of study characteristics, weight of evidence, and utility as a basis for deriving reference values and unit risks.
- Standardized evidence tables for all health outcomes need to be developed. If there were appropriate tables, long text descriptions of studies could be moved to an appendix or deleted.
- Develop templates for evidence tables, forest plots, or other displays.
- Establish protocols for review of major types of studies, such as epidemiologic and bioassay.

Implementation:

➤ Evaluating and Documenting the Quality of Individual Studies

In Progress

The IRIS Program is improving the approach to evaluating and describing the strengths and weaknesses of critical studies and standardizing the documentation of this evaluation. This step in the systematic review process involves the evaluation of a variety of methodological features (e.g., study design, exposure measurement details, data analysis and presentation). The purpose of this step is generally not to eliminate studies, but rather to evaluate studies with respect to potential methodological considerations that could affect the interpretation of and relative confidence in the results. It is worth emphasizing that the systematic evaluation of the study described in this step is conducted at an early stage of assessment development (i.e., after identifying the relevant sources of primary data but before developing evidence tables and characterizing hazards associated with exposure to a chemical). The results of this systematic evaluation may inform decisions about which studies to use for hazard identification, considerations to keep in mind when interpreting the results of specific studies, and which studies to move forward for dose-response modeling for derivation of toxicity values.



For more detailed information, see “Study Quality Evaluation” and “Documentation of Study Quality Evaluations” in the Evaluation and Display of Individual Studies section in the draft Handbook for IRIS Assessment Development in Appendix F.



See also Section 4 (“Evaluating the quality of individual studies”) in the Preamble to IRIS Toxicological Reviews in Appendix B.



A chemical-specific example of the implementation of this recommendation is available as “EXAMPLE 2 – Evaluation and Display of Individual Studies” in the Chemical-specific Examples Demonstrating Implementation of NRC Recommendations document.

➤ Evidence Tables

Implemented

The IRIS Program has developed templates for evidence tables to standardize the presentation of reviewed studies in IRIS assessments. Once a literature search has been conducted and the resulting database of studies has been evaluated, evidence tables are developed to present information from the collection of studies related to a specific outcome or endpoint of toxicity. The evidence tables include studies that have been judged adequate for hazard identification and display available study results, both positive and negative results. The studies that are considered to be most informative will depend on the extent and nature of the database for a given chemical, but may encompass a range of study designs and include epidemiology, toxicology, and, other toxicity data when appropriate.



For more detailed information, see “Reporting Study Results” in the Evaluation and Display of Individual Studies section in the draft Handbook for IRIS Assessment Development in Appendix F.



A chemical-specific example of the implementation of this recommendation is available as “EXAMPLE 3 – Evidence Tables” in the Chemical-specific Examples Demonstrating Implementation of NRC Recommendations document.

Weight-of-Evidence Evaluation: Integration of Evidence for Hazard Identification

NRC Recommendations:

- Strengthened, more integrative and more transparent discussions of weight of evidence are needed. The discussions would benefit from more rigorous and systematic coverage of the various determinants of weight of evidence, such as consistency.
- Review use of existing weight-of-evidence guidelines.
- Standardize approach to using weight-of-evidence guidelines.
- Conduct agency workshops on approaches to implementing weight-of-evidence guidelines.
- Develop uniform language to describe strength of evidence on noncancer effects.
- Expand and harmonize the approach for characterizing uncertainty and variability.
- To the extent possible, unify consideration of outcomes around common modes of action rather than considering multiple outcomes separately.

Implementation:

➤ Integration of Evidence for Hazard Identification

In Progress

The IRIS Program has strengthened and increased transparency in the weight-of-evidence for identifying hazards in IRIS assessments. Hazard identification involves the integration of evidence from human, animal, and mechanistic studies in order to draw conclusions about the hazards associated with exposure to a chemical. In general, IRIS assessments integrate evidence in the context of Hill (1965), which outlines aspects — such as consistency, strength, coherence, specificity, does-response, temporality, and biological plausibility — for consideration of causality

in epidemiologic investigations that were later modified by others and extended to experimental studies (U.S. EPA, 2005a).

All results, both positive and negative, of potentially relevant studies that have been evaluated for quality are considered (U.S. EPA, 2002) to answer the fundamental question: “Does exposure to chemical X cause hazard Y?” This requires a critical weighing of the available evidence (U.S. EPA, 2005a; 1994), but is not to be interpreted as a simple tallying of the number of positive and negative studies (U.S. EPA, 2002). Hazards are identified by an informed, expert evaluation and integration of the human, animal, and mechanistic evidence streams.



For more detailed information, see “Synthesis of Observational Epidemiology Evidence”, “Synthesis of Animal Toxicology Evidence”, and “Mechanistic Considerations in Elucidating Adverse Outcome Pathways” in the Evaluating the Overall Evidence of Each Effect section in the draft Handbook for IRIS Assessment Development in Appendix F.



See also Section 5 (“Evaluating the overall evidence of each effect”) in the Preamble to IRIS Toxicological Reviews in Appendix B.



A chemical-specific example of the implementation of this recommendation is available as “EXAMPLE 4 – Evidence Integration” in the Chemical-specific Examples Demonstrating Implementation of NRC Recommendations document.

Currently, the IRIS Program is using existing guidelines that address these issues to inform assessments. In addition, the IRIS Program is taking a more systematic approach in analyzing the available human, animal, and mechanistic data is being used in IRIS assessments. In conducting this analysis and developing the synthesis, the IRIS Program evaluates the data for the:

- strength of the relationship between the exposure and response and the presence of a dose-response relationship;
- specificity of the response to chemical exposure and whether the exposure precedes the effect;
- consistency of the association between the chemical exposure and response; and
- biological plausibility of the response or effect and its relevance to humans.

The IRIS Program uses this weight of evidence approach to identify the potential hazards associated with chemical exposure.

The IRIS Program recognizes the benefit of adopting a formal weight-of-evidence framework that includes standardized classification of causality. In addition to the NRC task, in which the panel will review current methods for evidence-based reviews and recommend approaches for weighing scientific evidence for chemical hazard and dose-response assessments, the IRIS Program is planning to convene a workshop to discuss approaches to evidence integration. As part of this workshop, the various approaches that are currently in use will be acknowledged and compared for their strengths and limitations. The workshop will include scientists with expertise in the

classification of chemicals for various health effects. The workshop will be open to the public, and the details will be publicly announced.



The “Integration of Evidence Evaluation” section in the draft Handbook for IRIS Assessment Development in Appendix F is currently under development.

Selection of Studies for Derivation of Toxicity Values

NRC Recommendations:

- The rationales for the selection of the studies that are advanced for consideration in calculating the RfCs and unit risks need to be expanded. All candidate RfCs should be evaluated together with the aid of graphic displays that incorporate selected information on attributes relevant to the database.
- Establish clear guidelines for study selection.
- Balance strengths and weaknesses.
- Weigh human vs. experimental evidence.
- Determine whether combining estimates among studies is warranted.

Implementation:

➤ **Selection of Studies for Dose-Response Analysis**

Implemented

The IRIS Program has improved the process for selecting studies for derivation of toxicity values as well as increasing the transparency about this process by providing an improved discussion and rationale. Building on the individual study quality evaluations (described under *Evidence Evaluation: Hazard Identification* in this report) that identify strengths and weaknesses of individual studies, for each health effect for which there is credible evidence of hazard, a group of studies are identified and evaluated as part of the hazard identification. In evaluating these studies for selecting a subset to be considered for the derivation of toxicity values, the basic criterion is whether the quantitative exposure and response data are available to compute a point of departure (POD). The POD can be a no-observed-adverse-effect-level [NOAEL], lowest-observed-adverse-effect-level [LOAEL], or the benchmark dose/concentration lower confidence limit [BMDL/BMCL].

Additional attributes (aspects of the study, data characteristics, and relevant considerations) pertinent to derivation of toxicity values are used as criteria to evaluate the subset of studies for dose-response analysis. Thus, the most relevant, informative studies are selected to move forward. The new document structure provides for transparent discussion of the studies identified for dose-response analysis.



For more detailed information, see “Selection of Studies for Derivation of Toxicity Values” in the Dose-Response Analysis section in the draft Handbook for IRIS Assessment Development in Appendix F.



See also Section 6 (“Selecting studies for dose-response analysis”) in the Preamble to IRIS Toxicological Reviews in Appendix B.



A chemical-specific example of the implementation of this recommendation is available as “EXAMPLE 5 – Selecting Studies for Derivation of Toxicity Values” in the Chemical-specific Examples Demonstrating Implementation of NRC Recommendations document.

➤ **Considerations for Combining Data for Dose-Response Modeling**

In Progress

The IRIS Program is now routinely considering whether combining data among studies is warranted for the derivation of toxicity values. For most IRIS assessments, the POD had been derived based on data from a single study dataset. This is because in most cases, datasets are often expected to be heterogeneous for biological or study design reasons.

However, there are cases where conducting dose-response modeling after combining data from multiple studies can be considered, resulting in a single POD based on multiple datasets. For instance, this may be useful to increase precision in the POD or to quantify the impact of specific sources of heterogeneity. The IRIS Program has developed considerations for combining data for dose-response modeling to be taken into account when performing dose-response analysis for an IRIS assessment.

In addition, multiple PODs or toxicity values can be combined (considering, for example, the highest quality studies, the most sensitive outcomes, or a clustering of values) to derive a single, overall toxicity value (or “meta-value”). For example, the IRIS assessment for trichloroethylene (TCE) identified multiple candidate RfDs that fell within a narrow dose range, and selected an overall RfD that reflected the midpoint among the similar candidate RfDs. This RfD is supported by multiple effects/studies and lead to a more robust (i.e., less sensitive to limitations of individual studies) (for more information: <http://epa.gov/iris/subst/0199.htm>, U.S. EPA, 2011).



For more detailed information, see “Considerations for Combining Data for Dose-Response Modeling” in the Dose-Response Analysis section in the draft Handbook for IRIS Assessment Development in Appendix F.

Calculation of Reference Values and Unit Risks

NRC Recommendations:

- Describe and justify assumptions and models used. This step includes review of dosimetry models and the implications of the models for uncertainty factors; determination of appropriate points of departure (such as benchmark dose, no-observed-adverse-effect level, and lowest observed-adverse-effect level), and assessment of the analyses that underlie the points of departure.
- Provide explanation of the risk-estimation modeling processes (for example, a statistical or biologic model fit to the data) that are used to develop a unit risk estimate.
- Provide adequate documentation for conclusions and estimation of reference values and unit risks. As noted by the committee throughout the present report, sufficient support for conclusions in the formaldehyde draft IRIS assessment is often lacking. Given that the development of specific IRIS assessments and their conclusions are of interest to many stakeholders, it is important that they provide sufficient references and supporting documentation for their conclusions. Detailed appendixes, which might be made available only electronically, should be provided, when appropriate.
- Assess the sensitivity of derived estimates to model assumptions and end points selected. This step should include appropriate tabular and graphic displays to illustrate the range of the estimates and the effect of uncertainty factors on the estimates.

Implementation:

➤ **Conducting and Documenting Dose-Response Modeling and Deriving Toxicity Values**

Implemented

IRIS assessments, in general, include dose-response analysis to derive toxicity values. In response to NRC recommendations, the IRIS Program has improved the quality control of the overall dose-response modeling process and increased transparency by documenting the approach for conducting dose-response modeling. Part of this documentation is achieved with the addition of considerations for selecting organ/system-specific and overall toxicity values, and a streamlined dose-response modeling output (both part of the new document structure). Additionally, tools and approaches to manage data and ensure quality (e.g., Data Management and Quality Control for Dose-Response Modeling) in dose-response analyses have been developed. The objectives are to minimize errors, maintain a transparent system for data management, automate tasks where possible, and maintain an archive of data and calculations used to develop assessments.

The IRIS Program has improved the documentation of dose-response modeling. *Preamble* Section 7 provides a description of the process for dose-response analysis. In addition, the text describing the dose-response analysis will include a description of how the toxicity values were derived and will cite EPA guidelines where appropriate.



For more detailed information, see “Data Management and Quality Control for Dose-Response Modeling,” and “Considerations for Selecting Organ/System-Specific or Overall Toxicity Values” in the Dose-Response Analysis section in the draft Handbook for IRIS Assessment Development in Appendix F.



See also Section 7 (“Deriving toxicity values”) in the Preamble to IRIS Toxicological Reviews in Appendix B.



Chemical-specific examples of the implementation of this recommendation are available as “EXAMPLE 6 – Dose-Response Modeling Output” and “EXAMPLE 7 – Considerations for Selecting Organ/System-Specific or Overall Toxicity Values” in the Chemical-specific Examples Demonstrating Implementation of NRC Recommendations document.

IV. Additional Initiatives

External Peer Review Enhancements

IRIS Peer Review Basics

Implemented

Rigorous, independent peer review is a cornerstone of IRIS assessments. Every IRIS assessment is reviewed by a group of internationally recognized experts in scientific disciplines relevant for the particular assessment. The peer review process used for IRIS assessments follows EPA guidance on peer review². Most IRIS assessments are reviewed through contractor-organized or EPA’s Science Advisory Board (SAB) peer reviews. All peer reviews, regardless of the reviewing body, involve a public comment period and public meeting (usually face-to-face). Following peer review, all revised IRIS assessments include an appendix describing how peer review and public comments were addressed.

Dedicated Chemical Assessment Advisory Committee

EPA’s SAB has established a new standing committee, the Chemical Assessment Advisory Committee (CAAC), to review IRIS assessments. In the past, the SAB formed a new committee for each chemical assessment that the SAB reviewed. The new CAAC will provide the same high-level, transparent review as previous SAB reviews, but it will provide more continuous and overlapping membership for consistent advice.

The CAAC is comprised of 26 highly qualified scientists with a broad range of expertise relevant to human health assessment. The CAAC members will serve on panels reviewing individual IRIS assessments. Panels will be supplemented with added consultants who have expertise on the specific chemical substance or other areas of expertise needed to review the assessment. The CAAC review process is expected to be similar to how IRIS assessment reviews are currently conducted by the SAB and will include the following: the public will be invited to nominate peer reviewers for specific assessments; the proposed panels or pools of panelists will be posted for public comment; the proposed panelists will be screened by an Agency official for conflicts of interest; the final panel will be announced prior to the peer review phase.

² U.S. EPA (2006) *Science Policy Council Peer Review Handbook - 3rd Edition*, EPA document number EPA/100/B-06/002. (<http://www.epa.gov/peerreview/>) and the EPA National Center for Environmental Assessment Policy and Procedures for Conducting IRIS Peer Reviews (2009, http://www.epa.gov/iris/pdfs/Policy_IRIS_Peer_Reviews.pdf).

V. Summary

EPA is committed to a strong, vital, and scientifically sound IRIS Program. Over the past two years, EPA has worked to strengthen and streamline the IRIS Program, improving transparency and creating efficiencies. Significant changes have been made in response to the NRC recommendations and further efforts are underway to fully implement the recommendations.

1 **Appendix A – IRIS Toxicological Review Template**



3
4 www.epa.gov/iris

5
6
7
8 **Toxicological Review of [Chemical]**

9
10 **[CASRN X-X-X]**

11
12 **In Support of Summary Information on the**
13 **Integrated Risk Information System (IRIS)**

14
15
16
17
18 **DATE**

19
20
21 **NOTICE**

22
23 This document is an **[Agency Review, Interagency Science Consultation, Public Comment,**
24 **External Review, or Final Agency/Interagency Science Discussion] draft.** This information is
25 distributed solely for the purpose of pre-dissemination peer review under applicable information
26 quality guidelines. It has not been formally disseminated by EPA. It does not represent and should
27 not be construed to represent any Agency determination or policy. It is being circulated for review
28 of its technical accuracy and science policy implications.

29
30
31
32
33 National Center for Environmental Assessment
34 Office of Research and Development
35 U.S. Environmental Protection Agency
36 Washington, DC

This document is a draft for review purposes only and does not constitute Agency policy.

1
2
3
4
5
6
7

DISCLAIMER

This document is a preliminary draft for review purposes only. This information is distributed solely for the purpose of pre-dissemination peer review under applicable information quality guidelines. It has not been formally disseminated by EPA. It does not represent and should not be construed to represent any Agency determination or policy. Mention of trade names or commercial products does not constitute endorsement of recommendation for use.

CONTENTS

1	
2	
3	AUTHORS CONTRIBUTORS REVIEWERS.....
4	PREFACE
5	PREAMBLE TO IRIS TOXICOLOGICAL REVIEWS.....
6	EXECUTIVE SUMMARY
7	LITERATURE SEARCH STRATEGY STUDY SELECTION.....
8	1. HAZARD IDENTIFICATION
9	1.1. SYNTHESIS OF EVIDENCE
10	1.1.1. Hazard A
11	1.1.2. Hazard B
12	1.1.3. Hazard C.....
13	1.1.4. Carcinogenicity
14	1.1.5. Other Toxicological Effects
15	1.2. SUMMARY AND EVALUATION
16	1.2.1. Integration of Evidence for Effects Other Than Cancer
17	1.2.2. Integration of Evidence for Carcinogenicity
18	1.2.3. Susceptible Populations and Lifestages
19	2. DOSE-RESPONSE ANALYSIS
20	2.1. ORAL REFERENCE DOSE FOR EFFECTS OTHER THAN CANCER
21	2.1.1. Identification of Studies and Effects for Dose-Response Analysis
22	2.1.2. Methods of Analysis
23	2.1.3. Derivation of Candidate Values.....
24	2.1.4. Derivation of Organ/System-specific Reference Doses.....
25	2.1.5. Selection of the Proposed Overall Reference Dose
26	2.1.6. Uncertainties in the Derivation of Reference Dose.....
27	2.1.7. Confidence Statement.....
28	2.1.8. Previous IRIS Assessment
29	2.2. INHALATION REFERENCE CONCENTRATION FOR EFFECTS OTHER THAN CANCER.....
30	2.2.1. Identification of Studies and Effects for Dose-Response Analysis

This document is a draft for review purposes only and does not constitute Agency policy.

1	2.2.2. Methods of Analysis
2	2.2.3. Derivation of Candidate Values
3	2.2.4. Derivation of Organ/System-specific Reference Concentrations.....
4	2.2.5. Selection of the Proposed Overall Reference Concentration
5	2.2.6. Uncertainties in the Derivation of Reference Concentration
6	2.2.7. Confidence Statement.....
7	2.2.8. Previous IRIS Assessment
8	2.3. ORAL SLOPE FACTOR FOR CANCER
9	2.3.1. Analysis of Carcinogenicity Data
10	2.3.2. Dose-Response Analysis—Adjustments and Extrapolations Methods
11	2.3.3. Derivation of the Oral Slope Factor.....
12	2.3.4. Uncertainties in the Derivation of the Oral Slope Factor
13	2.3.5. Previous IRIS Assessment: Oral Slope Factor
14	2.4. INHALATION UNIT RISK FOR CANCER
15	2.4.1. Analysis of Carcinogenicity Data
16	2.4.2. Dose-Response Analysis—Adjustments and Extrapolations Methods
17	2.4.3. Inhalation Unit Risk Derivation.....
18	2.4.4. Uncertainties in the Derivation of the Inhalation Unit Risk
19	2.4.5. Previous IRIS Assessment: Inhalation Unit Risk.....
20	2.5. APPLICATION OF AGE-DEPENDENT ADJUSTMENT FACTORS.....
21	REFERENCES.....
22	

1 **TABLES**

2 Table ES-1. Summary of reference dose (RfD) derivation

3 Table ES-2. Summary of reference concentration (RfC) derivation.....

4 Table 2-1. Summary of derivation of points of departure following oral exposure

5 Table 2-2. Effects and corresponding derivation of candidate RfDs.....

6 Table 2-3. Organ/system-specific RfDs and proposed overall RfD for [chemical].....

7 Table 2-4. Summary of derivation of points of departure following inhalation exposure

8 Table 2-5. Effects and corresponding derivation of candidate RfCs.....

9 Table 2-6. Organ/system-specific RfCs and proposed overall RfC for [chemical].....

10 Table 2-11. Summary of uncertainty in the derivation of cancer risk values for [chemical].....

11 Table 2-16. Summary of uncertainty in the derivation of cancer risk values for [chemical].....

12 **FIGURES**

13 Figure 2-1. Candidate RfDs with corresponding POD and composite UF.

14 Figure 2-2. Candidate RfCs with corresponding POD and composite UF.

15

1 ABBREVIATIONS

2				
	α 2u-g	alpha2u-globulin	LOAEL	lowest-observed-adverse-effect level
	ACGIH	American Conference of Governmental Industrial Hygienists	MN	micronuclei
	AIC	Akaike's information criterion	MNPCE	micronucleated polychromatic erythrocyte
	ALD	approximate lethal dosage	MTD	maximum tolerated dose
	ALT	alanine aminotransferase	NAG	N-acetyl- β -D-glucosaminidase
	AST	aspartate aminotransferase	NCEA	National Center for Environmental Assessment
	atm	atmosphere	NCI	National Cancer Institute
	ATSDR	Agency for Toxic Substances and Disease Registry	NOAEL	no-observed-adverse-effect level
	BMD	benchmark dose	NTP	National Toxicology Program
	BMDL	benchmark dose lower confidence limit	NZW	New Zealand White (rabbit breed)
	BMDS	Benchmark Dose Software	OCT	ornithine carbamoyl transferase
	BMR	benchmark response	ORD	Office of Research and Development
	BUN	blood urea nitrogen	PBPK	physiologically based pharmacokinetic
	BW	body weight	PCNA	proliferating cell nuclear antigen
	CA	chromosomal aberration	POD	point of departure
	CASRN	Chemical Abstracts Service Registry Number	POD _[AD]	duration-adjusted POD
	CBI	covalent binding index	QSAR	quantitative structure-activity relationship
	CHO	Chinese hamster ovary (cell line cells)	RDS	replicative DNA synthesis
	CL	confidence limit	RfC	inhalation reference concentration
	CNS	central nervous system	RfD	oral reference dose
	CPN	chronic progressive nephropathy	RGDR	regional gas dose ratio
	CYP450	cytochrome P450	RNA	ribonucleic acid
	DAF	dosimetric adjustment factor	SAR	structure activity relationship
	DEN	diethylnitrosamine	SCE	sister chromatid exchange
	DMSO	dimethylsulfoxide	SD	standard deviation
	DNA	deoxyribonucleic acid	SDH	sorbitol dehydrogenase
	EPA	Environmental Protection Agency	SE	standard error
	FDA	Food and Drug Administration	SGOT	glutamic oxaloacetic transaminase, also known as AST
	FEV ₁	forced expiratory volume of 1 second	SGPT	glutamic pyruvic transaminase, also known as ALT
	GD	gestation day	SSD	systemic scleroderma
	GDH	glutamate dehydrogenase	TCA	trichloroacetic acid
	GGT	γ -glutamyl transferase	TCE	trichloroethylene
	GSH	glutathione	TWA	time-weighted average
	GST	glutathione-S-transferase	UF	uncertainty factor
	Hb/g-A	animal blood:gas partition coefficient	UF _A	animal-to-human uncertainty factor
	Hb/g-H	human blood:gas partition coefficient	UF _H	human variation uncertainty factor
	HEC	human equivalent concentration	UF _L	LOAEL-to-NOAEL uncertain factor
	HED	human equivalent dose	UF _S	subchronic-to-chronic uncertainty factor
	i.p.	intraperitoneal	UF _D	database deficiencies uncertainty factor
	IRIS	Integrated Risk Information System	U.S.	United States of America
	IVF	in vitro fertilization		
	LC ₅₀	median lethal concentration		
	LD ₅₀	median lethal dose		

This document is a draft for review purposes only and does not constitute Agency policy.

AUTHORS | CONTRIBUTORS | REVIEWERS

Assessment Team

NAME (alphabetical order) U.S. EPA/ORD/NCEA
NAME Washington, DC
NAME (Chemical Manager)
NAME

Scientific Support Team

NAME (alphabetical order) U.S. EPA/ORD/NCEA
NAME³ Washington, DC
NAME
NAME

Production Team

NAME (alphabetical order) Agency, Office, Location
NAME Agency, Office, Location
NAME Agency, Office, Location
NAME Agency, Office, Location

Contractor Support

NAME Company, Location
NAME
NAME

NAME Company, Location
NAME
NAME

Executive Direction

Kenneth Olden, Ph.D., Sc.D., L.H.D. (Center Director) U.S. EPA/ORD/NCEA
John Vandenberg, Ph.D. (National Program Director, HHRA) Washington, DC
Lynn Flowers, Ph.D., DABT (Associate Director for Health)
Vincent Cogliano, Ph.D.⁴ (IRIS Program Director—acting)
Samantha Jones, Ph.D.⁵ (IRIS Associate Director for Science)
[\[Division Director\]](#)
[\[Branch Chief\]](#)

³ Chemical Assessment Support Team (CAST) Member

⁴ Chemical Assessment Support Team (CAST) Lead

⁵ Chemical Assessment Support Team (CAST) Lead

This document is a draft for review purposes only and does not constitute Agency policy.

Internal Review Team

NAME	Agency, Office, Location
NAME	Agency, Office, Location
NAME	Agency, Office, Location
NAME	Agency, Office, Location

1

Reviewers

2 This assessment was provided for review to scientists in EPA's program and regional offices.
3 Comments were submitted by:

Office, Location
Office, Location
Office, Location
Office, Location

4 This assessment was provided for review to other federal agencies and Executive Offices of the
5 President. Comments were submitted by:

AGENCY
AGENCY
AGENCY
AGENCY

6 This assessment was released for public comment on [month] [day], [year] and comments were due
7 on [month] [day], [year]. The external peer-review comments are available on the IRIS Web site. A
8 summary and EPA's disposition of the comments received from the independent external peer
9 reviewers and from the public is included in Appendix [X] and is also available on the IRIS Web site.
10 Comments were received from the following entities:

NAME	Affiliation, Location
NAME	Affiliation, Location
NAME	Affiliation, Location
NAME	Affiliation, Location

11 This assessment was peer reviewed by independent expert scientists external to EPA and a peer-
12 review meeting was held on [month] [day], [year]. The external peer-review comments are
13 available on the IRIS Web site. A summary and EPA's disposition of the comments received from
14 the independent external peer reviewers and from the public is included in Appendix [X] and is also
15 available on the IRIS Web site.

NAME	Affiliation, Location
NAME	Affiliation, Location
NAME	Affiliation, Location
NAME	Affiliation, Location

16

PREFACE

This Toxicological Review critically reviews the publicly available studies on [chemical] in order to identify its adverse health effects and to characterize exposure-response relationships. The assessment covers [...] It was prepared under the auspices of EPA’s Integrated Risk Information System (IRIS) Program.

[Chemical] is listed as [...] [Why is EPA interested in this assessment? Is the chemical included on Agency lists (ex. HAPs, DWCL)?]

[If this is a reassessment...] This assessment updates a previous IRIS assessment of [chemical] that was developed in [year]. The previous assessment included [...]. New information has become available and this assessment reviews information on all health effects by all exposure routes. Organ/system-specific RfDs are calculated based on [applicable hazards, e.g., developmental, reproductive and immune system toxicity data]. These toxicity values may be useful for cumulative risk assessments that consider the combined effect of multiple agents acting on the same biological system.

This assessment was conducted in accordance with EPA guidance, which is cited and summarized in the *Preamble to IRIS Toxicological Reviews*. The findings of this assessment and related documents produced during its development are available on the IRIS Web site (<http://www.epa.gov/iris>). Appendices for chemical and physical properties, toxicokinetic information, and summaries of toxicology studies and other information are provided as *Supplemental Information* to this assessment.

For additional information about this assessment or for general questions regarding IRIS, please contact EPA’s IRIS Hotline at 202-566-1676 (phone), 202-566-1749 (fax), or hotline.iris@epa.gov.

Assessments by Other National and International Health Agencies

Toxicity information on [chemical] has been evaluated by [...]. The results of these assessments are presented in Appendix A of the Supplemental Information. It is important to recognize that these assessments may have been prepared for different purposes and may utilize different methods, and that newer studies may be included in the IRIS assessment.

Chemical Properties and Uses

[Appendix B...]

1

2

PREAMBLE TO IRIS TOXICOLOGICAL REVIEWS

3

This document is a draft for review purposes only and does not constitute Agency policy.

EXECUTIVE SUMMARY

Occurrence and Health Effects

[Placeholder for text]

Effects Other Than Cancer Observed Following Oral Exposure

Oral Reference Dose (RfD) for Effects Other Than Cancer

Table ES-1. Summary of reference dose (RfD) derivation

Critical effect	Point of departure*	UF	Chronic RfD
-----------------	---------------------	----	-------------

[* Conversion Factors and Assumptions—]

Confidence in the Chronic Oral RfD

Effects Other Than Cancer Observed Following Inhalation Exposure

Inhalation Reference Concentration (RfC) for Effects Other Than Cancer

Table ES-2. Summary of reference concentration (RfC) derivation

Critical effect	Point of departure*	UF	Chronic RfC
-----------------	---------------------	----	-------------

[* Conversion Factors and Assumptions—]

Confidence in the Chronic Inhalation RfC

Evidence for Human Carcinogenicity

Quantitative Estimate of Carcinogenic Risk From Oral Exposure

1 **Quantitative Estimate of Carcinogenic Risk From Inhalation Exposure**

2

3 **Susceptible Populations and Lifestages**

4

5 **Key Issues Addressed in Assessment**

6

1

2

3

LITERATURE SEARCH STRATEGY | STUDY SELECTION AND EVALUATION

4

1. HAZARD IDENTIFICATION

1 1.1. SYNTHESIS OF EVIDENCE

2

3

4 1.2. SUMMARY AND EVALUATION

5

2. DOSE-RESPONSE ANALYSIS

2.1. ORAL REFERENCE DOSE FOR EFFECTS OTHER THAN CANCER

The RfD (expressed in units of mg/kg-day) is defined as an estimate (with uncertainty spanning perhaps an order of magnitude) of a daily exposure to the human population (including sensitive subgroups) that is likely to be without an appreciable risk of deleterious effects during a lifetime. It can be derived from a no-observed-adverse-effect level (NOAEL), lowest-observed-adverse-effect level (LOAEL), or the 95% lower bound on the benchmark dose (BMDL), with uncertainty factors (UFs) generally applied to reflect limitations of the data used.

2.1.1. Identification of Studies and Effects for Dose-Response Analysis

Hazard A

Hazard B

Hazard C

2.1.2. Methods of Analysis

Table 2-1 summarizes the sequence of calculations leading to the derivation of a human-equivalent point of departure for each data set discussed above.

Table 2-1. Summary of derivation of points of departure following oral exposure

Endpoint and reference	Species/sex	Model	BMR	BMD	BMDL	POD _{ADJ}	POD _{HED}
Hazard A (ex. DEVELOPMENTAL)							
Hazard B (ex. REPRODUCTIVE)							

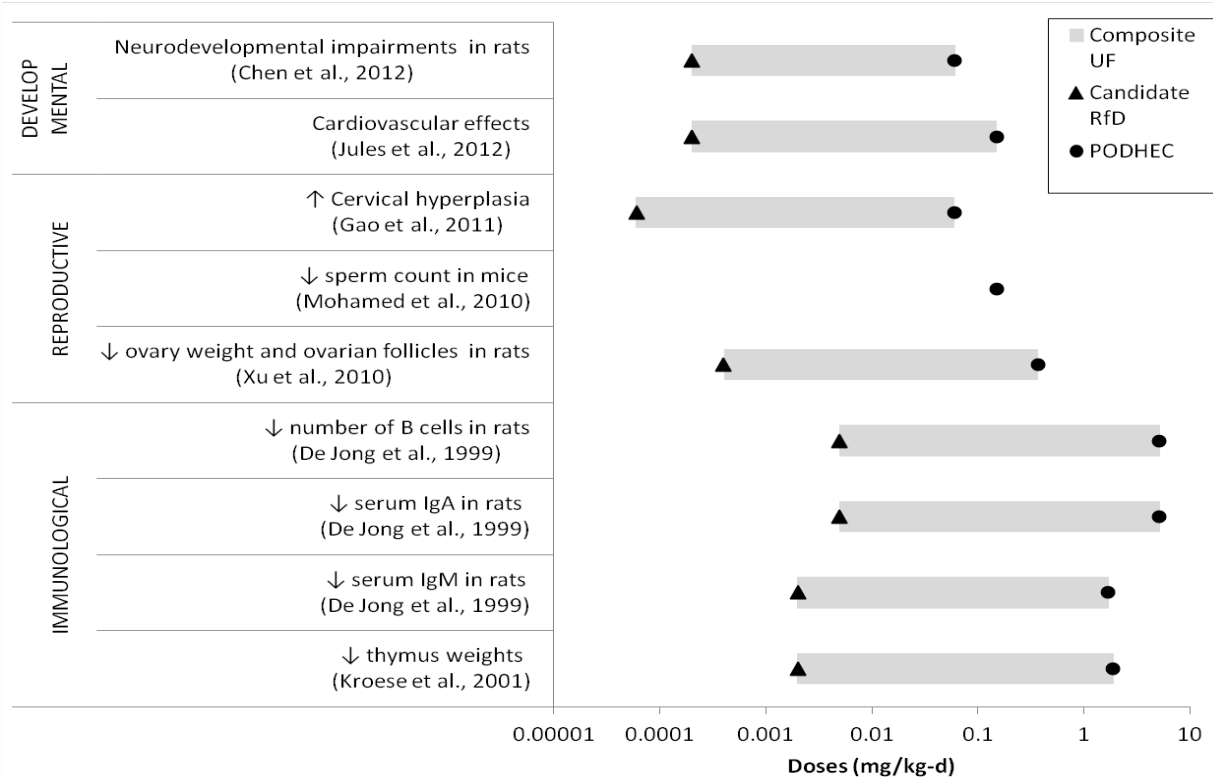
1 **2.1.3. Derivation of Candidate Values**

2 Table 2-2 is a continuation of Table 2-1 and summarizes the application of uncertainty
 3 factors to each point of departure to derive candidate values for each data set. The candidate values
 4 presented in Table 2-2 are preliminary to the derivation of reference values in subsequent sections.
 5 The selection of uncertainty factors is based on EPA’s *Review of the Reference Dose and Reference*
 6 *Concentration Processes* (U.S. EPA, 2002; Section 4.4.5) and is described in Section 7.6 of the
 7 *Preamble*. Figure 2-1 presents graphically the candidate values, uncertainty factors, and points of
 8 departure, with each bar corresponding to one data set described in Tables 2-1 and 2-2.

9 **Table 2-2. Effects and corresponding derivation of candidate RfDs**

Endpoint and reference	POD _{HED}	POD type	UF _A	UF _H	UF _L	UF _S	UF _D	Composite UF	Candidate value (mg/kg-day)
Hazard A (ex. DEVELOPMENTAL)									
Hazard B (ex. REPRODUCTIVE)									

10
 11 [Insert rationale for the application of uncertainty factors. The value (e.g., 1, 3, or 10) of the
 12 uncertainty factor will depend upon the availability of data and what is known about the
 13 chemical...]
 14



1

2

3

Figure 2-1. Candidate RfDs with corresponding POD and composite UF.
[Note: Data shown here are provided only for illustrative purposes]

4

2.1.4. Derivation of Organ/System-Specific Reference Doses

5

Table 2-3 distills the candidate values from Table 2-2 into a single value for each organ or system. These organ or system-specific reference values may be useful for subsequent cumulative risk assessments that consider the combined effect of multiple agents acting at a common site.

8

Hazard A

9

10

Hazard B

11

12

Hazard C

13

1 **Table 2-3. Organ/system-specific RfDs and proposed overall RfD for**
 2 **[chemical]**

Effect	Basis	RfD (mg/kg-day)	Exposure description	Confidence
Hazard A			Ex. chronic	
Hazard B			Ex. gestational	
Hazard C			Ex. subchronic	
Proposed overall RfD			Ex. gestational	

3
 4 **2.1.5. Selection of the Proposed Overall Reference Dose**

5
 6 **2.1.6. Uncertainties in the Derivation of Reference Dose**

7
 8 **2.1.7. Confidence Statement**

9 A confidence level of high, medium, or low is assigned to the study used to derive the RfD,
 10 the overall database, and the RfD itself, as described in Section 4.3.9.2 of EPA's *Methods for*
 11 *Derivation of Inhalation Reference Concentrations and Application of Inhalation Dosimetry* (U.S. EPA,
 12 1994).

13 **2.1.8. Previous IRIS Assessment**

15 **2.2. INHALATION REFERENCE CONCENTRATION FOR EFFECTS OTHER THAN**
 16 **CANCER**

17 The RfC (expressed in units of mg/m³) is defined as an estimate (with uncertainty spanning
 18 perhaps an order of magnitude) of a continuous inhalation exposure to the human population
 19 (including sensitive subgroups) that is likely to be without an appreciable risk of deleterious effects
 20 during a lifetime. It can be derived from a NOAEL, LOAEL, or the 95% lower bound on the
 21 benchmark concentration (BMCL), with UFs generally applied to reflect limitations of the data used.

22 **2.2.1. Identification of Studies and Effects for Dose-Response Analysis**

23 ***Hazard A***

1 **Hazard B**

2

3 **Hazard C**

4

5 **2.2.2. Methods of Analysis**

6 Table 2-4 summarizes the sequence of calculations leading to the derivation of a human-
7 equivalent point of departure for each data set discussed above.

8

9 **Table 2-4. Summary of derivation of points of departure following inhalation**
10 **exposure**

Endpoint and reference	Species/sex	Model	BMR	BMC	BMCL	POD _{ADJ}	POD _{HED}
Hazard A (ex. DEVELOPMENTAL)							
Hazard B (ex. REPRODUCTIVE)							

11

12 **2.2.3. Derivation of Candidate Values**

13 Table 2-5 is a continuation of Table 2-4 and summarizes the application of uncertainty
14 factors to each point of departure to derive a candidate values for each data set. The candidate
15 values presented in Table 2-5 are for exploratory purposes only, and are preliminary to the
16 derivation of reference values in subsequent sections. The selection of uncertainty factors was
17 based on EPA's *Review of the Reference Dose and Reference Concentration Processes* (U.S. EPA, 2002;
18 Section 4.4.5) and is described in Section 7.6 of the *Preamble*. Figure 2-2 graphically presents these
19 candidate values, uncertainty factors, and points of departure with each bar corresponding to one
20 data set described in Tables 2-4 and 2-5.

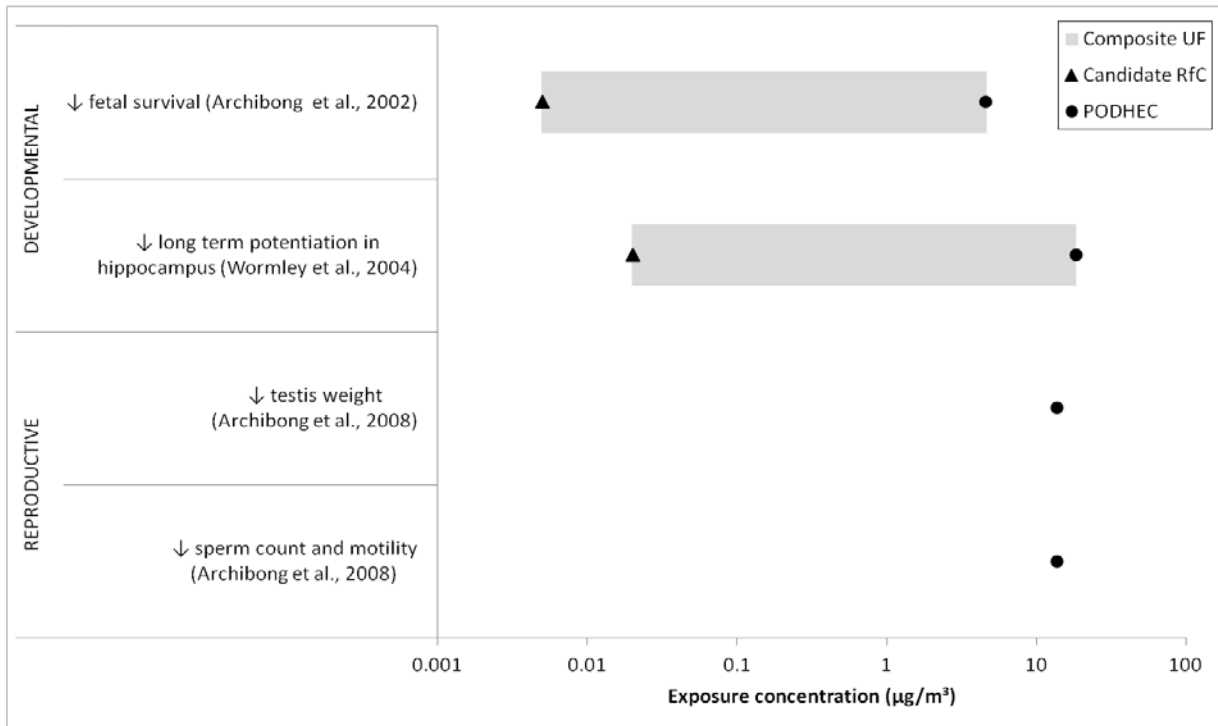
1

Table 2-5. Effects and corresponding derivation of candidate RfCs

Endpoint	POD _{HEC} (µg/m ³)	POD type	UF _A	UF _H	UF _L	UF _S	UF _D	Composite UF	Candidate value (µg/m ³)
Hazard A (ex. DEVELOPMENTAL)									
Hazard B (ex. REPRODUCTIVE)									

2
3
4
5
6
7
8
9
10

[Insert rationale for the application of uncertainty factors. The value (e.g., 1, 3, or 10) of the uncertainty factor will depend upon the availability of data and what is known about the chemical...]



11
12
13

Figure 2-2. Candidate RfCs with corresponding POD and composite UF.
 [Note: Data shown here are provided only for illustrative purposes.]

1 **2.2.4. Derivation of Organ/System-Specific Reference Concentrations**

2 Table 2-6 distills the candidate values from Table 2-5 into a single value for each organ or
3 system. These organ or system-specific reference values may be useful for subsequent cumulative
4 risk assessments that consider the combined effect of multiple agents acting at a common site.

5 **Hazard A**

6

7 **Hazard B**

8

9 **Hazard C**

10

11 **Table 2-6. Organ/system-specific RfCs and proposed overall RfC for**
12 **[chemical]**

Effect	Basis	RfC (mg/m ³)	Exposure description	Confidence
Hazard A				
Hazard B				
Hazard C				
Proposed overall RfC				

13

14 **2.2.5. Selection of the Proposed Overall Reference Concentration**

15

16 **2.2.6. Uncertainties in the Derivation of Reference Concentration**

17

18 **2.2.7. Confidence Statement**

19 A confidence level of high, medium, or low is assigned to the study used to derive the RfC,
20 the overall database, and the RfC itself, as described in Section 4.3.9.2 of EPA's *Methods for*
21 *Derivation of Inhalation Reference Concentrations and Application of Inhalation Dosimetry* (U.S. EPA,
22 1994).

1 **2.2.8. Previous IRIS Assessment**

2

3 **2.3. ORAL SLOPE FACTOR FOR CANCER**

4 The carcinogenicity assessment provides information on the carcinogenic hazard potential
5 of the substance in question, and quantitative estimates of risk from oral and inhalation exposure
6 may be derived. Quantitative risk estimates may be derived from the application of a low-dose
7 extrapolation procedure. If derived, the oral slope factor is a plausible upper bound on the estimate
8 of risk per mg/kg-day of oral exposure. [Note: Similarly, an inhalation unit risk is a plausible upper
9 bound on the estimate of risk per $\mu\text{g}/\text{m}^3$ air breathed.]

10 **2.3.1. Analysis of Carcinogenicity Data**

11

12 **2.3.2. Dose-Response Analysis—Adjustments and Extrapolations Methods**

13

14 **2.3.3. Derivation of the Oral Slope Factor**

15

16 **2.3.4. Uncertainties in the Derivation of the Oral Slope Factor**

17

18 **2.3.5. Previous IRIS Assessment: Oral Slope Factor**

19

20 **2.4. INHALATION UNIT RISK FOR CANCER**

21 The carcinogenicity assessment provides information on the carcinogenic hazard potential
22 of the substance in question and quantitative estimates of risk from oral and inhalation exposure
23 may be derived. Quantitative risk estimates may be derived from the application of a low-dose
24 extrapolation procedure. If derived, the inhalation unit risk is a plausible upper bound on the
25 estimate of risk per $\mu\text{g}/\text{m}^3$ air breathed.

26 **2.4.1. Analysis of Carcinogenicity Data**

27

28 **2.4.2. Dose-Response Analysis—Adjustments and Extrapolations Methods**

29

1 **2.4.3. Inhalation Unit Risk Derivation**

2

3 **2.4.4. Uncertainties in the Derivation of the Inhalation Unit Risk**

4

5 **2.4.5. Previous IRIS Assessment: Inhalation Unit Risk**

6

7 **2.5. APPLICATION OF AGE-DEPENDENT ADJUSTMENT FACTORS**

8

9

10

11

12

13

14

15

REFERENCES

Appendix B – Preamble to IRIS Toxicological Reviews

1. Scope of the IRIS Program

1 Soon after EPA was established in 1970, it was at
2 the forefront of developing risk assessment as a
3 science and applying it in decisions to protect
4 human health and the environment. The Clean
5 Air Act, for example, mandates that EPA provide
6 “an ample margin of safety to protect public
7 health”; the Safe Drinking Water Act, that “no
8 adverse effects on the health of persons may
9 reasonably be anticipated to occur, allowing an
10 adequate margin of safety.” Accordingly, EPA
11 uses information on the adverse effects of
12 chemicals and on exposure levels below which
13 these effects are not anticipated to occur.

15 IRIS assessments critically review the publicly
16 available studies to identify adverse health
17 effects from long-term exposure to chemicals and
18 to characterize exposure-response relationships.
19 In terms set forth by the National Research
20 Council (NRC, 1983), IRIS assessments cover the
21 hazard identification and dose-response
22 assessment steps of risk assessment, not the
23 exposure assessment or risk characterization
24 steps that are conducted by EPA’s program and
25 regional offices and by other federal, state, and
26 local health agencies that evaluate risk in specific
27 populations and exposure scenarios. IRIS
28 assessments are distinct from and do not address
29 political, economic, and technical considerations
30 that influence the design and selection of risk
31 management alternatives.

32 An IRIS assessment may cover a single chemical,
33 a group of structurally or toxicologically related
34 chemicals, or a complex mixture. Exceptions are
35 chemicals currently used exclusively as
36 pesticides, ionizing and non-ionizing radiation,
37 and criteria air pollutants listed under section
38 108 of the Clean Air Act (carbon monoxide, lead,
39 nitrogen oxides, ozone, particulate matter, and
40 sulfur oxides).

41 Periodically, the IRIS Program asks other EPA
42 programs and regions, other federal agencies,
43 state health agencies, and the general public to

44 nominate chemicals and mixtures for future
45 assessment or reassessment. These agents may
46 be found in air, water, soil, or sediment. Selection
47 is based on program and regional office priorities
48 and on availability of adequate information to
49 evaluate the potential for adverse effects. The
50 IRIS Program may assess other agents as an
51 urgent public health need arises. IRIS also
52 reassesses agents as significant new studies are
53 published.

2. Process for developing and peer-reviewing IRIS assessments

56 The process for developing IRIS assessments
57 (revised in May 2009) involves critical analysis of
58 the pertinent studies, opportunities for public
59 input, and multiple levels of scientific review.
60 EPA revises draft assessments after each review,
61 and external drafts and comments become part
62 of the public record (U.S. EPA, 2009).

63 **Step 1. Development of a draft Toxicological**
64 **Review** (generally about 11-1/2 months
65 duration). The draft assessment considers all
66 pertinent publicly available studies and
67 applies consistent criteria to evaluate study
68 quality, identify health effects, identify
69 mechanistic events and pathways, integrate
70 the evidence of causation for each effect, and
71 derive toxicity values. A public dialogue
72 meeting prior to the integration of evidence
73 and derivation of toxicity values promotes
74 public discussion of the literature search,
75 evidence, and key issues.

76 **Step 2. Internal review by scientists in EPA**
77 **programs and regions** (2 months). The
78 draft assessment is revised to address
79 comments from within EPA.

80 **Step 3. Interagency science consultation with**
81 **other federal agencies and the Executive**
82 **Offices of the President** (1-1/2 months).
83 The draft assessment is revised to address
84 the interagency comments. The science
85 consultation draft, interagency comments,
86 and EPA’s response to major comments
87 become part of the public record.

This document is a draft for review purposes only and does not constitute Agency policy.

1 **Step 4. Public review and comment, followed**
2 **by external peer review** (3-1/2 months or
3 more, depending on the review process).
4 EPA releases the draft assessment for public
5 review and comment. Another public
6 dialogue meeting provides an opportunity to
7 discuss the assessment prior to peer review.
8 EPA addresses the public comments and
9 releases a draft for external peer review. The
10 peer reviewers assess whether the evidence
11 has been assembled and evaluated according
12 to guidelines and whether the conclusions
13 are justified by the evidence. The peer
14 review meeting is open to the public and
15 includes time for oral public comments. The
16 peer review draft, peer review report, and
17 written public comments become part of the
18 public record.
19 **Step 5. Revision of draft Toxicological Review**
20 **and development of draft IRIS summary**
21 (2 months). The draft assessment is revised
22 to reflect the peer review comments, public
23 comments, and newly published studies that
24 are critical to the conclusions of the
25 assessment. The disposition of peer review
26 comments and public comments becomes
27 part of the public record.
28 **Step 6. Final EPA review and interagency**
29 **science discussion with other federal**
30 **agencies and the Executive Offices of the**
31 **President** (1-1/2 months). The draft
32 assessment and summary are revised to
33 address EPA and interagency comments. The
34 science discussion draft, written interagency
35 comments, and EPA's response to major
36 comments become part of the public record.
37 **Step 7. Completion and posting** (1 month). The
38 Toxicological Review and IRIS summary are
39 posted on the IRIS web site ([http://](http://www.epa.gov/iris/)
40 www.epa.gov/iris/).

41 The remainder of this Preamble addresses step 1,
42 the development of a draft Toxicological Review.
43 IRIS assessments follow standard practices of
44 evidence evaluation and peer review, many of
45 which are discussed in EPA guidelines (U.S. EPA,
46 1986a, 1986b, 1991, 1996, 1998, 2000, 2005a,
47 2005b) and other methods (U.S. EPA, 1994, 2002,
48 2006a, 2006b, 2011, 2012a, 2012b). A practical

49 draft *Handbook* is available for use by IRIS
50 assessment teams (U.S. EPA, 2013). Transparent
51 application of scientific judgment is of
52 paramount importance. To provide a harmonized
53 approach across IRIS assessments, this Preamble
54 summarizes concepts from these guidelines and
55 emphasizes principles of general applicability.

56 **3. Identifying and selecting pertinent** 57 **studies**

58 **3.1 Identifying studies**

59 Before beginning an assessment, EPA conducts a
60 comprehensive search of the primary scientific
61 literature. The literature search follows standard
62 practices and includes the PubMed and ToxNet
63 databases of the National Library of Medicine,
64 Web of Science, and other databases listed in
65 EPA's HERO system (Health and Environmental
66 Research Online, <http://hero.epa.gov/>). Searches
67 for information on mechanisms of toxicity are
68 inherently specialized and may include studies
69 on other agents that act through related
70 mechanisms.

71 Each assessment specifies the search strategies,
72 keywords, and cut-off dates of its literature
73 searches. EPA posts the results of the literature
74 search on the IRIS web site and requests
75 information from the public on additional studies
76 and ongoing research.

77 EPA also considers studies received through the
78 IRIS Submission Desk and studies (typically
79 unpublished) submitted under the Toxic
80 Substances Control Act or the Federal Insecticide,
81 Fungicide, and Rodenticide Act. Material
82 submitted as Confidential Business Information
83 is considered only if it includes health and safety
84 data that can be publicly released. If a study that
85 may be critical to the conclusions of the
86 assessment has not been peer-reviewed, EPA will
87 have it peer-reviewed.

88 EPA also examines the toxicokinetics of the agent
89 to identify other chemicals (for example, major
90 metabolites of the agent) to include in the
91 assessment if adequate information is available,
92 in order to more fully explain the toxicity of the
93 agent and to suggest dose metrics for subsequent
94 modeling.

1 In assessments of chemical mixtures, mixture
2 studies are preferred for their ability to reflect
3 interactions among components. The literature
4 search seeks, in decreasing order of preference
5 (U.S. EPA, 1986a, 2000):

- 6 – Studies of the mixture being assessed.
- 7 – Studies of a sufficiently similar mixture. In
8 evaluating similarity, the assessment
9 considers the alteration of mixtures in the
10 environment through partitioning and
11 transformation.
- 12 – Studies of individual chemical components of
13 the mixture, if there are not adequate studies
14 of sufficiently similar mixtures.

15 **3.2 Selecting pertinent epidemiologic**
16 **studies**

17 Study design is the key consideration for
18 selecting pertinent epidemiologic studies from
19 the results of the literature search.

- 20 – Cohort studies, case-control studies, and
21 some population-based surveys (for
22 example, NHANES) provide the strongest
23 epidemiologic evidence, especially when
24 they collect information about individual
25 exposures and effects.
- 26 – Ecological studies (geographic correlation
27 studies) relate exposures and effects by
28 geographic area. They can provide strong
29 evidence if there are large exposure
30 contrasts between geographic areas,
31 relatively little exposure variation within
32 study areas, and population migration is
33 limited.
- 34 – Case reports of high or accidental exposure
35 lack definition of the population at risk and
36 the expected number of cases. They can
37 provide information about a rare effect or
38 about the relevance of analogous results in
39 animals.

40 The assessment briefly reviews ecological studies
41 and case reports but reports details only if they
42 suggest effects not identified by other studies.

43 **3.3 Selecting pertinent experimental**
44 **studies**

45 Exposure route is a key design consideration for
46 selecting pertinent experimental animal studies
47 or human clinical studies.

- 48 – Studies of oral, inhalation, or dermal
49 exposure involve passage through an
50 absorption barrier and are considered most
51 pertinent to human environmental exposure.
- 52 – Injection or implantation studies are often
53 considered less pertinent but may provide
54 valuable toxicokinetic or mechanistic
55 information. They also may be useful for
56 identifying effects in animals if deposition or
57 absorption is problematic (for example, for
58 particles and fibers).

59 Exposure duration is also a key design
60 consideration for selecting pertinent
61 experimental animal studies.

- 62 – Studies of effects from chronic exposure are
63 most pertinent to lifetime human exposure.
- 64 – Studies of effects from less-than-chronic
65 exposure are pertinent but less preferred for
66 identifying effects from lifetime human
67 exposure. Such studies may be indicative of
68 effects from less-than-lifetime human
69 exposure.

70 Short-duration studies involving animals or
71 humans may provide toxicokinetic or
72 mechanistic information.

73 For developmental toxicity and reproductive
74 toxicity, irreversible effects may result from a
75 brief exposure during a critical period of
76 development. Accordingly, specialized study
77 designs are used for these effects (U.S. EPA, 1991,
78 1996, 1998, 2006b).

79 **4. Evaluating the quality of individual**
80 **studies**

81 After the subsets of pertinent epidemiologic and
82 experimental studies have been selected from the
83 literature searches, the assessment evaluates the
84 quality of each individual study. This evaluation
85 considers the design, methods, conduct, and
86 documentation of each study, but not whether
87 the results are positive, negative, or null. The

1 objective is to identify the stronger, more
2 informative studies based on a uniform
3 evaluation of quality characteristics across
4 studies of similar design.

5 **4.1 Evaluating the quality of** 6 **epidemiologic studies**

7 The assessment evaluates design and
8 methodological aspects that can increase or
9 decrease the weight given to each epidemiologic
10 study in the overall evaluation (U.S. EPA, 1991,
11 1994, 1996, 1998, 2005a):

- 12 – Documentation of study design, methods,
13 population characteristics, and results.
- 14 – Definition and selection of the study group
15 and comparison group.
- 16 – Ascertainment of exposure to the chemical
17 or mixture.
- 18 – Ascertainment of disease or health effect.
- 19 – Duration of exposure and follow-up and
20 adequacy for assessing the occurrence of
21 effects.
- 22 – Characterization of exposure during critical
23 periods.
- 24 – Sample size and statistical power to detect
25 anticipated effects.
- 26 – Participation rates and potential for selection
27 bias as a result of the achieved participation
28 rates.
- 29 – Measurement error (can lead to
30 misclassification of exposure, health
31 outcomes, and other factors) and other types
32 of information bias.
- 33 – Potential confounding and other sources of
34 bias addressed in the study design or in the
35 analysis of results. The basis for
36 consideration of confounding is a reasonable
37 expectation that the confounder is related to
38 both exposure and outcome and is
39 sufficiently prevalent to result in bias.

40 For developmental toxicity, reproductive toxicity,
41 neurotoxicity, and cancer there is further
42 guidance on the nuances of evaluating
43 epidemiologic studies of these effects (U.S. EPA,
44 1991, 1996, 1998, 2005a).

45 **4.2 Evaluating the quality of** 46 **experimental studies**

47 The assessment evaluates design and
48 methodological aspects that can increase or
49 decrease the weight given to each experimental
50 animal study, in-vitro study, or human clinical
51 study (U.S. EPA, 1991, 1994, 1996, 1998, 2005a).
52 Research involving human subjects is considered
53 only if conducted according to ethical principles.

- 54 – Documentation of study design, animals or
55 study population, methods, basic data, and
56 results.
- 57 – Nature of the assay and validity for its
58 intended purpose.
- 59 – Characterization of the nature and extent of
60 impurities and contaminants of the
61 administered chemical or mixture.
- 62 – Characterization of dose and dosing regimen
63 (including age at exposure) and their
64 adequacy to elicit adverse effects, including
65 latent effects.
- 66 – Sample sizes and statistical power to detect
67 dose-related differences or trends.
- 68 – Ascertainment of survival, vital signs, disease
69 or effects, and cause of death.
- 70 – Control of other variables that could
71 influence the occurrence of effects.

72 The assessment uses statistical tests to evaluate
73 whether the observations may be due to chance.

74 The standard for determining statistical
75 significance of a response is a trend test or
76 comparison of outcomes in the exposed groups
77 against those of concurrent controls. In some
78 situations, examination of historical control data
79 from the same laboratory within a few years of
80 the study may improve the analysis. For an
81 uncommon effect that is not statistically
82 significant compared with concurrent controls,
83 historical controls may show that the effect is
84 unlikely to be due to chance. For a response that
85 appears significant against a concurrent control
86 response that is unusual, historical controls may
87 offer a different interpretation (U.S. EPA, 2005a).

88 For developmental toxicity, reproductive toxicity,
89 neurotoxicity, and cancer there is further
90 guidance on the nuances of evaluating
91 experimental studies of these effects (U.S. EPA,

1 1991, 1996, 1998, 2005a). In multi-generation
2 studies, agents that produce developmental
3 effects at doses that are not toxic to the maternal
4 animal are of special concern. Effects that occur
5 at doses associated with mild maternal toxicity
6 are not assumed to result only from maternal
7 toxicity. Moreover, maternal effects may be
8 reversible, while effects on the offspring may be
9 permanent (U.S. EPA, 1991, 1998).

10 **4.3 Reporting study results**

11 The assessment uses evidence tables to present
12 the design and key results of pertinent studies.
13 There may be separate tables for each site of
14 toxicity or type of study.

15 If a large number of studies observe the same
16 effect, the assessment considers the study quality
17 characteristics in this section to identify the
18 strongest studies or types of study. The tables
19 present details from these studies, and the
20 assessment explains the reasons for not
21 reporting details of other studies or groups of
22 studies that do not add new information.
23 Supplemental information provides references to
24 all studies considered, including those not
25 summarized in the tables.

26 The assessment discusses strengths and
27 limitations that affect the interpretation of each
28 study. If the interpretation of a study in the
29 assessment differs from that of the study authors,
30 the assessment discusses the basis for the
31 difference.

32 As a check on the selection and evaluation of
33 pertinent studies, EPA asks peer reviewers to
34 identify studies that were not adequately
35 considered.

36 **5. Evaluating the overall evidence of** 37 **each effect**

38 **5.1 Concepts of causal inference**

39 For each health effect, the assessment evaluates
40 the evidence as a whole to determine whether it
41 is reasonable to infer a causal association
42 between exposure to the agent and the
43 occurrence of the effect. This inference is based
44 on information from pertinent human studies,
45 animal studies, and mechanistic studies of

46 adequate quality. Positive, negative, and null
47 results are given weight according to study
48 quality.

49 Causal inference involves scientific judgment,
50 and the considerations are nuanced and complex.
51 Several health agencies have developed
52 frameworks for causal inference, among them the
53 U.S. Surgeon General (DHEW, 1964; DHHS,
54 2004), the International Agency for Research on
55 Cancer (2006), the Institute of Medicine (2008),
56 and the U.S. Environmental Protection Agency
57 (2005a, 2010). Although developed for different
58 purposes, the frameworks are similar in nature
59 and provide an established structure and
60 language for causal inference. Each considers
61 aspects of an association that suggest causation,
62 discussed by Hill (1965) and elaborated by
63 Rothman and Greenland (1998) (U.S. EPA, 1994,
64 2002, 2005a).

65 **Strength of association:** The finding of a large
66 relative risk with narrow confidence
67 intervals strongly suggests that an
68 association is not due to chance, bias, or
69 other factors. Modest relative risks, however,
70 may reflect a small range of exposures, an
71 agent of low potency, an increase in an effect
72 that is common, exposure misclassification,
73 or other sources of bias.

74 **Consistency of association:** An inference of
75 causation is strengthened if elevated risks
76 are observed in independent studies of
77 different populations and exposure
78 scenarios. Reproducibility of findings
79 constitutes one of the strongest arguments
80 for causation. Discordant results sometimes
81 reflect differences in study design, exposure,
82 or confounding factors.

83 **Specificity of association:** As originally
84 intended, this refers to one cause associated
85 with one effect. Current understanding that
86 many agents cause multiple effects and many
87 effects have multiple causes make this a less
88 informative aspect of causation, unless the
89 effect is rare or unlikely to have multiple
90 causes.

91 **Temporal relationship:** A causal interpretation
92 requires that exposure precede development
93 of the effect.

1 **Biologic gradient (exposure-response**
 2 **relationship):** Exposure-response
 3 relationships strongly suggest causation. A
 4 monotonic increase is not the only pattern
 5 consistent with causation. The presence of an
 6 exposure-response gradient also weighs
 7 against bias and confounding as the source of
 8 an association.

9 **Biologic plausibility:** An inference of causation
 10 is strengthened by data demonstrating
 11 plausible biologic mechanisms, if available.
 12 Plausibility may reflect subjective prior
 13 beliefs if there is insufficient understanding
 14 of the biologic process involved.

15 **Coherence:** An inference of causation is
 16 strengthened by supportive results from
 17 animal experiments, toxicokinetic studies,
 18 and short-term tests. Coherence may also be
 19 found in other lines of evidence, such as
 20 changing disease patterns in the population.

21 **“Natural experiments”:** A change in exposure
 22 that brings about a change in disease
 23 frequency provides strong evidence, as it
 24 tests the hypothesis of causation. An example
 25 would be an intervention to reduce exposure
 26 in the workplace or environment that is
 27 followed by a reduction of an adverse effect.

28 **Analogy:** Information on structural analogues or
 29 on chemicals that induce similar mechanistic
 30 events can provide insight into causation.

31 These considerations are consistent with
 32 guidelines for systematic reviews that evaluate
 33 the quality and weight of evidence. Confidence is
 34 increased if the magnitude of effect is large, if
 35 there is evidence of an exposure-response
 36 relationship, or if an association was observed
 37 and the plausible biases would tend to decrease
 38 the magnitude of the reported effect. Confidence
 39 is decreased for study limitations, inconsistency
 40 of results, indirectness of evidence, imprecision,
 41 or reporting bias (Guyatt et al., 2008a,b).

42 **5.2 Evaluating evidence in humans**

43 For each effect, the assessment evaluates the
 44 evidence from the epidemiologic studies as a
 45 whole. The objective is to determine whether a
 46 credible association has been observed and, if so,
 47 whether that association is consistent with
 48 causation. In doing this, the assessment explores

49 alternative explanations (such as chance, bias,
 50 and confounding) and draws a conclusion about
 51 whether these alternatives can satisfactorily
 52 explain any observed association.

53 To make clear how much the epidemiologic
 54 evidence contributes to the overall weight of the
 55 evidence, the assessment may select a standard
 56 descriptor to characterize the epidemiologic
 57 evidence of association between exposure to the
 58 agent and occurrence of a health effect.

59 **Sufficient epidemiologic evidence of an**
 60 **association consistent with causation:** The
 61 evidence establishes a causal association for
 62 which alternative explanations such as
 63 chance, bias, and confounding can be ruled
 64 out with reasonable confidence.

65 **Suggestive epidemiologic evidence of an**
 66 **association consistent with causation:** The
 67 evidence suggests a causal association but
 68 chance, bias, or confounding cannot be ruled
 69 out as explaining the association.

70 **Inadequate epidemiologic evidence to infer a**
 71 **causal association:** The available studies do
 72 not permit a conclusion regarding the
 73 presence or absence of an association.

74 **Epidemiologic evidence consistent with no**
 75 **causal association:** Several adequate studies
 76 covering the full range of human exposures
 77 and considering susceptible populations, and
 78 for which alternative explanations such as
 79 bias and confounding can be ruled out, are
 80 mutually consistent in not finding an
 81 association.

82 **5.3 Evaluating evidence in animals**

83 For each effect, the assessment evaluates the
 84 evidence from the animal experiments as a whole
 85 to determine the extent to which they indicate a
 86 potential for effects in humans. Consistent results
 87 across various species and strains increase
 88 confidence that similar results would occur in
 89 humans. Several concepts discussed by Hill
 90 (1965) are pertinent to the weight of
 91 experimental results: consistency of response,
 92 dose-response relationships, strength of
 93 response, biologic plausibility, and coherence
 94 (U.S. EPA, 1994, 2002, 2005a).

1 In weighing evidence from multiple experiments,
2 U.S. EPA (2005a) distinguishes
3 **Conflicting evidence** (that is, mixed positive and
4 negative results in the same sex and strain
5 using a similar study protocol) from
6 **Differing results** (that is, positive results and
7 negative results are in different sexes or
8 strains or use different study protocols).

9 Negative or null results do not invalidate positive
10 results in a different experimental system. EPA
11 regards all as valid observations and looks to
12 explain differing results using mechanistic
13 information (for example, physiologic or
14 metabolic differences across test systems) or
15 methodological differences (for example, relative
16 sensitivity of the tests, differences in dose levels,
17 insufficient sample size, or timing of dosing or
18 data collection).

19 It is well established that there are critical
20 periods for some developmental and
21 reproductive effects. Accordingly, the assessment
22 determines whether critical periods have been
23 adequately investigated (U.S. EPA, 1991, 1996,
24 1998, 2005a, 2005b, 2006b). Similarly, the
25 assessment determines whether the database is
26 adequate to evaluate other critical sites and
27 effects.

28 In evaluating evidence of genetic toxicity:

- 29 – Demonstration of gene mutations,
30 chromosome aberrations, or aneuploidy in
31 humans or experimental mammals (*in vivo*)
32 provides the strongest evidence.
- 33 – This is followed by positive results in lower
34 organisms or in cultured cells (*in vitro*) or for
35 other genetic events.
- 36 – Negative results carry less weight, partly
37 because they cannot exclude the possibility
38 of effects in other tissues (IARC, 2006).

39 For germ-cell mutagenicity, EPA has defined
40 categories of evidence, ranging from positive
41 results of human germ-cell mutagenicity to
42 negative results for all effects of concern (U.S.
43 EPA, 1986b).

44 **5.4 Evaluating mechanistic data to**
45 **identify adverse outcome pathways**
46 **and modes of action**

47 Mechanistic data can be useful in answering
48 several questions.

- 49 – The biologic plausibility of a causal
50 interpretation of human studies.
- 51 – The generalizability of animal studies to
52 humans.
- 53 – The susceptibility of particular populations
54 or lifestages.

55 The focus of the analysis is to describe, if
56 possible, *adverse outcome pathways* that lead to a
57 health effect. An adverse outcome pathway
58 encompasses:

- 59 – *Toxicokinetic processes* of absorption,
60 distribution, metabolism, and elimination
61 that lead to the formation of an active agent
62 and its presence at the site of initial biologic
63 interaction.
- 64 – *Toxicodynamic processes* that lead to a health
65 effect at this or another site (also known as a
66 *mode of action*).

67 For each effect, the assessment discusses the
68 available information on its *modes of action* and
69 associated *key events* (*key events* being
70 empirically observable, necessary precursor
71 steps or biologic markers of such steps; *mode of*
72 *action* being a series of key events involving
73 interaction with cells, operational and anatomic
74 changes, and resulting in disease). Pertinent
75 information may also come from studies of
76 metabolites or of compounds that are
77 structurally similar or that act through similar
78 mechanisms. Information on mode of action is
79 not required for a conclusion that the agent is
80 causally related to an effect (U.S. EPA, 2005a).

81 The assessment addresses several questions
82 about each hypothesized mode of action (U.S.
83 EPA, 2005a).

84 (1) **Is the hypothesized mode of action**
85 **sufficiently supported in test animals?**
86 Strong support for a key event being
87 necessary to a mode of action can come from
88 experimental challenge to the hypothesized
89 mode of action, in which studies that

1 suppress a key event observe suppression of
 2 the effect. Support for a mode of action is
 3 meaningfully strengthened by consistent
 4 results in different experimental models,
 5 much more so than by replicate experiments
 6 in the same model. The assessment may
 7 consider various aspects of causation in
 8 addressing this question.

9 (2) **Is the hypothesized mode of action**
 10 **relevant to humans?** The assessment
 11 reviews the key events to identify critical
 12 similarities and differences between the test
 13 animals and humans. Site concordance is not
 14 assumed between animals and humans,
 15 though it may hold for certain effects or
 16 modes of action. Information suggesting
 17 quantitative differences in doses where
 18 effects would occur in animals or humans is
 19 considered in the dose-response analysis.
 20 Current levels of human exposure are not
 21 used to rule out human relevance, as IRIS
 22 assessments may be used in evaluating new
 23 or unforeseen circumstances that may entail
 24 higher exposures.

25 (3) **Which populations or lifestages can be**
 26 **particularly susceptible to the**
 27 **hypothesized mode of action?** The
 28 assessment reviews the key events to
 29 identify populations and lifestages that might
 30 be susceptible to their occurrence.
 31 Quantitative differences may result in
 32 separate toxicity values for susceptible
 33 populations or lifestages.

34 The assessment discusses the likelihood that an
 35 agent operates through multiple modes of action.
 36 An uneven level of support for different modes of
 37 action can reflect disproportionate resources
 38 spent investigating them (U.S. EPA, 2005a). It
 39 should be noted that in clinical reviews, the
 40 credibility of a series of studies is reduced if
 41 evidence is limited to studies funded by one
 42 interested sector (Guyatt et al., 2008b).

43 For cancer, the assessment evaluates evidence of
 44 a mutagenic mode of action to guide
 45 extrapolation to lower doses and consideration
 46 of susceptible lifestages. Key data include the
 47 ability of the agent or a metabolite to react with
 48 or bind to DNA, positive results in multiple test

49 systems, or similar properties and structure-
 50 activity relationships to mutagenic carcinogens
 51 (U.S. EPA, 2005a),

52 **5.5 Characterizing the overall weight of** 53 **the evidence**

54 After evaluating the human, animal, and
 55 mechanistic evidence pertinent to an effect, the
 56 assessment answers the question: Does the agent
 57 cause the adverse effect? (NRC, 1983, 2009). In
 58 doing this, the assessment develops a narrative
 59 that integrates the evidence pertinent to
 60 causation. To provide clarity and consistency, the
 61 narrative includes a standard hazard descriptor.
 62 For example, the following standard descriptors
 63 combine epidemiologic, experimental, and
 64 mechanistic evidence of carcinogenicity (U.S.
 65 EPA, 2005a).

66 ***Carcinogenic to humans:*** There is convincing
 67 epidemiologic evidence of a causal
 68 association (that is, there is reasonable
 69 confidence that the association cannot be
 70 fully explained by chance, bias, or
 71 confounding); or there is strong human
 72 evidence of cancer or its precursors,
 73 extensive animal evidence, identification of
 74 key precursor events in animals, and strong
 75 evidence that they are anticipated to occur in
 76 humans.

77 ***Likely to be carcinogenic to humans:*** The
 78 evidence demonstrates a potential hazard to
 79 humans but does not meet the criteria for
 80 *carcinogenic*. There may be a plausible
 81 association in humans, multiple positive
 82 results in animals, or a combination of
 83 human, animal, or other experimental
 84 evidence.

85 ***Suggestive evidence of carcinogenic potential:***
 86 The evidence raises concern for effects in
 87 humans but is not sufficient for a stronger
 88 conclusion. This descriptor covers a range of
 89 evidence, from a positive result in the only
 90 available study to a single positive result in
 91 an extensive database that includes negative
 92 results in other species.

93 ***Inadequate information to assess carcinogenic***
 94 ***potential:*** No other descriptors apply.
 95 *Conflicting evidence* can be classified as
 96 *inadequate information* if all positive results

1 are opposed by negative studies of equal
 2 quality in the same sex and strain. *Differing*
 3 *results*, however, can be classified as
 4 *suggestive evidence* or as *likely to be*
 5 *carcinogenic*.

6 ***Not likely to be carcinogenic to humans:*** There
 7 is robust evidence for concluding that there
 8 is no basis for concern. There may be no
 9 effects in both sexes of at least two
 10 appropriate animal species; positive animal
 11 results and strong, consistent evidence that
 12 each mode of action in animals does not
 13 operate in humans; or convincing evidence
 14 that effects are not likely by a particular
 15 exposure route or below a defined dose.

16 Multiple descriptors may be used if there is
 17 evidence that carcinogenic effects differ by dose
 18 range or exposure route (U.S. EPA, 2005a).

19 Another example of standard descriptors comes
 20 from EPA’s Integrated Science Assessments,
 21 which evaluate causation for the effects of the
 22 criteria pollutants in ambient air (U.S. EPA,
 23 2010).

24 ***Causal relationship:*** Sufficient evidence to
 25 conclude that there is a causal relationship.
 26 Observational studies cannot be explained by
 27 plausible alternatives, or they are supported
 28 by other lines of evidence, for example,
 29 animal studies or mechanistic information.

30 ***Likely to be a causal relationship:*** Sufficient
 31 evidence that a causal relationship is likely,
 32 but important uncertainties remain. For
 33 example, observational studies show an
 34 association but co-exposures are difficult to
 35 address or other lines of evidence are limited
 36 or inconsistent; or multiple animal studies
 37 from different laboratories demonstrate
 38 effects and there are limited or no human
 39 data.

40 ***Suggestive of a causal relationship:*** At least one
 41 high-quality epidemiologic study shows an
 42 association but other studies are
 43 inconsistent.

44 ***Inadequate to infer a causal relationship:*** The
 45 studies do not permit a conclusion regarding
 46 the presence or absence of an association.

47 ***Not likely to be a causal relationship:*** Several
 48 adequate studies, covering the full range of

49 human exposure and considering susceptible
 50 populations, are mutually consistent in not
 51 showing an effect at any level of exposure.

52 EPA is investigating and may on a trial basis use
 53 these or other standard descriptors to
 54 characterize the overall weight of the evidence
 55 for effects other than cancer.

56 **6. Selecting studies for derivation of**
 57 **toxicity values**

58 For each effect where there is credible evidence
 59 of an association with the agent, the assessment
 60 derives toxicity values if there are suitable
 61 epidemiologic or experimental data. The decision
 62 to derive toxicity values may be linked to the
 63 hazard descriptor.

64 Dose-response analysis requires quantitative
 65 measures of dose and response. Then, other
 66 factors being equal (U.S. EPA, 1994, 2005a):

- 67 – Epidemiologic studies are preferred over
 68 animal studies, if quantitative measures of
 69 exposure are available and effects can be
 70 attributed to the agent.
- 71 – Among experimental animal models, those
 72 that respond most like humans are
 73 preferred, if the comparability of response
 74 can be determined.
- 75 – Studies by a route of human environmental
 76 exposure are preferred, although a validated
 77 toxicokinetic model can be used to
 78 extrapolate across exposure routes.
- 79 – Studies of longer exposure duration and
 80 follow-up are preferred, to minimize
 81 uncertainty about whether effects are
 82 representative of lifetime exposure.
- 83 – Studies with multiple exposure levels are
 84 preferred for their ability to provide
 85 information about the shape of the exposure-
 86 response curve.
- 87 – Studies with adequate power to detect
 88 effects at lower exposure levels are
 89 preferred, to minimize the extent of
 90 extrapolation to levels found in the
 91 environment.

92 Studies with non-monotonic exposure-response
 93 relationships are not necessarily excluded from
 94 the analysis. A diminished effect at higher

1 exposure levels may be satisfactorily explained
 2 by factors such as competing toxicity, saturation
 3 of absorption or metabolism, exposure
 4 misclassification, or selection bias.
 5 If a large number of studies are suitable for dose-
 6 response analysis, the assessment considers the
 7 study characteristics in this section to focus on
 8 the most informative data. The assessment
 9 explains the reasons for not analyzing other
 10 groups of studies. As a check on the selection of
 11 studies for dose-response analysis, EPA asks peer
 12 reviewers to identify studies that were not
 13 adequately considered.

14 **7. Deriving toxicity values**

15 **7.1 General framework for dose-response**
 16 **analysis**

17 EPA uses a two-step approach that distinguishes
 18 analysis of the observed dose-response data from
 19 inferences about lower doses (U.S. EPA, 2005a).

20 Within the observed range, the preferred
 21 approach is to use modeling to incorporate a
 22 wide range of data into the analysis. The
 23 modeling yields a *point of departure* (an exposure
 24 level near the lower end of the observed range,
 25 without significant extrapolation to lower doses)
 26 (sections 7.2-7.3).

27 Extrapolation to lower doses considers what is
 28 known about the modes of action for each effect
 29 (sections 7.4-7.5). When response estimates at
 30 lower doses are not required, an alternative is to
 31 derive *reference values*, which are calculated by
 32 applying factors to the point of departure in
 33 order to account for sources of uncertainty and
 34 variability (section 7.6).

35 For a group of agents that induce an effect
 36 through a common mode of action, the dose-
 37 response analysis may derive a *relative potency*
 38 *factor* for each agent. A full dose-response
 39 analysis is conducted for one well-studied *index*
 40 *chemical* in the group, then the potencies of other
 41 members are expressed in relative terms based
 42 on relative toxic effects, relative absorption or
 43 metabolic rates, quantitative structure-activity
 44 relationships, or receptor binding characteristics
 45 (U.S. EPA, 2000, 2005a).

46 Increasingly, EPA is basing toxicity values on
 47 combined analyses of multiple data sets or
 48 multiple responses. EPA also considers multiple
 49 dose-response approaches when they can be
 50 supported by robust data.

51 **7.2 Modeling dose to sites of biologic**
 52 **effects**

53 The preferred approach for analysis of dose is
 54 toxicokinetic modeling because of its ability to
 55 incorporate a wide range of data. The preferred
 56 dose metric would refer to the active agent at the
 57 site of its biologic effect or to a close, reliable
 58 surrogate measure. The active agent may be the
 59 administered chemical or a metabolite.
 60 Confidence in the use of a toxicokinetic model
 61 depends on the robustness of its validation
 62 process and on the results of sensitivity analyses
 63 (U.S. EPA, 1994, 2005a, 2006a).

64 Because toxicokinetic modeling can require
 65 many parameters and more data than are
 66 typically available, EPA has developed standard
 67 approaches that can be applied to typical data
 68 sets. These standard approaches also facilitate
 69 comparison across exposure patterns and
 70 species.

- 71 – Intermittent study exposures are
 72 standardized to a daily average over the
 73 duration of exposure. For chronic effects,
 74 daily exposures are averaged over the
 75 lifespan. Exposures during a critical period,
 76 however, are not averaged over a longer
 77 duration (U.S. EPA, 1991, 1996, 1998,
 78 2005a).
- 79 – Doses are standardized to equivalent human
 80 terms to facilitate comparison of results from
 81 different species.
 - 82 – Oral doses are scaled allometrically
 83 using $\text{mg}/\text{kg}^{3/4}\text{-d}$ as the equivalent dose
 84 metric across species. Allometric scaling
 85 pertains to equivalence across species,
 86 not across lifestages, and is not used to
 87 scale doses from adult humans or
 88 mature animals to infants or children
 89 (U.S. EPA, 2005a, 2011).
 - 90 – Inhalation exposures are scaled using
 91 dosimetry models that apply species-
 92 specific physiologic and anatomic factors

1 and consider whether the effect occurs
 2 at the site of first contact or after
 3 systemic circulation (U.S. EPA, 1994,
 4 2012b).

5 It can be informative to convert doses across
 6 exposure routes. If this is done, the assessment
 7 describes the underlying data, algorithms, and
 8 assumptions (U.S. EPA, 2005a).

9 In the absence of study-specific data on, for
 10 example, intake rates or body weight, EPA has
 11 developed recommended values for use in dose-
 12 response analysis (U.S. EPA, 1988).

13 **7.3 Modeling response in the range of**
 14 **observation**

15 Toxicodynamic (“biologically based”) modeling
 16 can incorporate data on biologic processes
 17 leading to an effect. Such models require
 18 sufficient data to ascertain a mode of action and
 19 to quantitatively support model parameters
 20 associated with its key events. Because different
 21 models may provide equivalent fits to the
 22 observed data but diverge substantially at lower
 23 doses, critical biologic parameters should be
 24 measured from laboratory studies, not by model
 25 fitting. Confidence in the use of a toxicodynamic
 26 model depends on the robustness of its
 27 validation process and on the results of
 28 sensitivity analyses. Peer review of the scientific
 29 basis and performance of a model is essential
 30 (U.S. EPA, 2005a).

31 Because toxicodynamic modeling can require
 32 many parameters and more knowledge and data
 33 than are typically available, EPA has developed a
 34 standard set of empirical (“curve-fitting”) models
 35 (<http://www.epa.gov/ncea/bmds/>) that can be
 36 applied to typical data sets, including those that
 37 are nonlinear. EPA has also developed guidance
 38 on modeling dose-response data, assessing
 39 model fit, selecting suitable models, and
 40 reporting modeling results (U.S. EPA, 2012a).
 41 Additional judgment or alternative analyses are
 42 used when the procedure fails to yield reliable
 43 results, for example, if the fit is poor, modeling
 44 may be restricted to the lower doses, especially if
 45 there is competing toxicity at higher doses (U.S.
 46 EPA, 2005a).

47 Modeling is used to derive a point of departure
 48 (U.S. EPA, 2005a, 2012a). (See section 7.6 for
 49 alternatives if a point of departure cannot be
 50 derived by modeling.)

- 51 – When linear extrapolation is used, selection
 52 of a response level corresponding to the
 53 point of departure is not highly influential, so
 54 standard values near the low end of the
 55 observable range are generally used (for
 56 example, 10% extra risk for cancer bioassay
 57 data, 1% for epidemiologic data, lower for
 58 rare cancers).
- 59 – For nonlinear approaches, both statistical
 60 and biologic considerations are taken into
 61 account.
 - 62 – For dichotomous data, a response level
 63 of 10% extra risk is generally used for
 64 minimally adverse effects, 5% or lower
 65 for more severe effects.
 - 66 – For continuous data, a response level is
 67 ideally based on an established
 68 definition of biologic significance. In the
 69 absence of such definition, one control
 70 standard deviation from the control
 71 mean is often used for minimally
 72 adverse effects, one-half standard
 73 deviation for more severe effects.

74 The point of departure is the 95% lower bound
 75 on the dose associated with the selected
 76 response level.

77 **7.4 Extrapolating to lower doses and**
 78 **response levels**

79 The purpose of extrapolating to lower doses is to
 80 estimate responses at exposures below the
 81 observed data. Low-dose extrapolation is
 82 typically used for cancer data. Low-dose
 83 extrapolation considers what is known about
 84 modes of action (U.S. EPA, 2005a).

- 85 (1) If a biologically based model has been
 86 developed and validated for the agent,
 87 extrapolation may use the fitted model below
 88 the observed range if significant model
 89 uncertainty can be ruled out with reasonable
 90 confidence.
- 91 (2) Linear extrapolation is used if the dose-
 92 response curve is expected to have a linear

1 component below the point of departure.
 2 This includes:
 3 – Agents or their metabolites that are
 4 DNA-reactive and have direct mutagenic
 5 activity.
 6 – Agents or their metabolites for which
 7 human exposures or body burdens are
 8 near doses associated with key events
 9 leading to an effect.
 10 Linear extrapolation is also used if there is
 11 an absence of sufficient information on
 12 modes of action.
 13 The result of linear extrapolation is
 14 described by an *oral slope factor* or an
 15 *inhalation unit risk*, which is the slope of the
 16 dose-response curve at lower doses or
 17 concentrations, respectively.
 18 (3) Nonlinear models are used for extrapolation
 19 if there are sufficient data to ascertain the
 20 mode of action and to conclude that it is not
 21 linear at lower doses, and the agent does not
 22 demonstrate mutagenic or other activity
 23 consistent with linearity at lower doses. If
 24 nonlinear extrapolation is appropriate but no
 25 model is developed, an alternative is to
 26 calculate reference values.

27 If linear extrapolation is used, the assessment
 28 develops a candidate slope factor or unit risk for
 29 each suitable data set. These results are arrayed,
 30 using common dose metrics, to show the
 31 distribution of relative potency across various
 32 effects and experimental systems. The
 33 assessment then derives or selects an overall
 34 slope factor and an overall unit risk for the agent,
 35 considering the various dose-response analyses,
 36 the study preferences discussed in section 6, and
 37 the possibility of basing a more robust result on
 38 multiple data sets.

39 **7.5 Considering susceptible populations**
 40 **and lifestages**

41 The assessment analyzes the available
 42 information on populations and lifestages that
 43 may be particularly susceptible to each effect. A
 44 tiered approach is used (U.S. EPA, 2005a).

45 (1) If an epidemiologic or experimental study
 46 reports quantitative results for a susceptible
 47 population or lifestage, these data are

48 analyzed to derive separate toxicity values
 49 for susceptible individuals.
 50 (2) If data on risk-related parameters allow
 51 comparison of the general population and
 52 susceptible individuals, these data are used
 53 to adjust the general-population toxicity
 54 values for application to susceptible
 55 individuals.
 56 (3) In the absence of chemical-specific data, EPA
 57 has developed *age-dependent adjustment*
 58 *factors* for early-life exposure to potential
 59 carcinogens that have a mutagenic mode of
 60 action. There is evidence of early-life
 61 susceptibility to various carcinogenic agents,
 62 but most epidemiologic studies and cancer
 63 bioassays do not include early-life exposure.
 64 To address the potential for early-life
 65 susceptibility, EPA recommends (U.S. EPA,
 66 2005b):
 67 – 10-fold adjustment for exposures before
 68 age 2 years.
 69 – 3-fold adjustment for exposures
 70 between ages 2 and 16 years.

71 **7.6 Reference values and uncertainty**
 72 **factors**

73 *An oral reference dose* or an *inhalation reference*
 74 *concentration* is an estimate of an exposure
 75 (including in susceptible subgroups) that is likely
 76 to be without an appreciable risk of adverse
 77 health effects over a lifetime (U.S. EPA, 2002).
 78 Reference values are typically calculated for
 79 effects other than cancer and for suspected
 80 carcinogens if a well characterized mode of
 81 action indicates that a necessary key event does
 82 not occur below a specific dose. Reference values
 83 provide no information about risks at higher
 84 exposure levels.

85 The assessment characterizes effects that form
 86 the basis for reference values as adverse,
 87 considered to be adverse, or a precursor to an
 88 adverse effect. For developmental toxicity,
 89 reproductive toxicity, and neurotoxicity there is
 90 guidance on adverse effects and their biologic
 91 markers (U.S. EPA, 1991, 1996, 1998).

92 To account for uncertainty and variability in the
 93 derivation of a lifetime human exposure where
 94 adverse effects are not anticipated to occur,

1 reference values are calculated by applying a
 2 series of *uncertainty factors* to the point of
 3 departure. If a point of departure cannot be
 4 derived by modeling, a no-observed-adverse-
 5 effect level or a lowest-observed-adverse-effect
 6 level is used instead. The assessment discusses
 7 scientific considerations involving several areas
 8 of variability or uncertainty.

9 **Human variation.** The assessment accounts for
 10 variation in susceptibility across the human
 11 population and the possibility that the
 12 available data may not be representative of
 13 individuals who are most susceptible to the
 14 effect. A factor of 10 is generally used to
 15 account for this variation. This factor is
 16 reduced only if the point of departure is
 17 derived or adjusted specifically for
 18 susceptible individuals (not for a general
 19 population that includes both susceptible
 20 and non-susceptible individuals) (U.S. EPA,
 21 1991, 1994, 1996, 1998, 2002).

22 **Animal-to-human extrapolation.** If animal
 23 results are used to make inferences about
 24 humans, the assessment adjusts for cross-
 25 species differences. These may arise from
 26 differences in toxicokinetics or
 27 toxicodynamics. Accordingly, if the point of
 28 departure is standardized to equivalent
 29 human terms or is based on toxicokinetic or
 30 dosimetry modeling, a factor of $10^{1/2}$
 31 (rounded to 3) is applied to account for the
 32 remaining uncertainty involving
 33 toxicodynamic differences. If a biologically
 34 based model adjusts fully for toxicokinetic
 35 and toxicodynamic differences across
 36 species, this factor is not used. In most other
 37 cases, a factor of 10 is applied (U.S. EPA,
 38 1991, 1994, 1996, 1998, 2002, 2011).

39 **Adverse-effect level to no-observed-adverse-
 40 effect level.** If a point of departure is based
 41 on a lowest-observed-adverse-effect level,
 42 the assessment must infer a dose where such
 43 effects are not expected. This can be a matter
 44 of great uncertainty, especially if there is no
 45 evidence available at lower doses. A factor of
 46 10 is applied to account for the uncertainty
 47 in making this inference. A factor other than
 48 10 may be used, depending on the magnitude
 49 and nature of the response and the shape of

50 the dose-response curve (U.S. EPA, 1991,
 51 1994, 1996, 1998, 2002).

52 **Subchronic-to-chronic exposure.** If a point of
 53 departure is based on subchronic studies, the
 54 assessment considers whether lifetime
 55 exposure could have effects at lower levels of
 56 exposure. A factor of 10 is applied to account
 57 for the uncertainty in using subchronic
 58 studies to make inferences about lifetime
 59 exposure. This factor may also be applied for
 60 developmental or reproductive effects if
 61 exposure covered less than the full critical
 62 period. A factor other than 10 may be used,
 63 depending on the duration of the studies and
 64 the nature of the response (U.S. EPA, 1994,
 65 1998, 2002).

66 **Incomplete database.** If an incomplete database
 67 raises concern that further studies might
 68 identify a more sensitive effect, organ
 69 system, or lifestage, the assessment may
 70 apply a database uncertainty factor (U.S.
 71 EPA, 1991, 1994, 1996, 1998, 2002). The size
 72 of the factor depends on the nature of the
 73 database deficiency. For example, EPA
 74 typically follows the suggestion that a factor
 75 of 10 be applied if both a prenatal toxicity
 76 study and a two-generation reproduction
 77 study are missing and a factor of $10^{1/2}$ if
 78 either is missing (U.S. EPA, 2002).

79 In this way, the assessment derives candidate
 80 values for each suitable data set and effect that is
 81 credibly associated with the agent. These results
 82 are arrayed, using common dose metrics, to show
 83 where effects occur across a range of exposures
 84 (U.S. EPA, 1994).

85 The assessment derives or selects an *organ- or
 86 system-specific reference value* for each organ or
 87 system affected by the agent. The assessment
 88 explains the rationale for each organ/system-
 89 specific reference value (based on, for example,
 90 the highest quality studies, the most sensitive
 91 outcome, or a clustering of values). By providing
 92 these organ/system-specific reference values,
 93 IRIS assessments facilitate subsequent
 94 cumulative risk assessments that consider the
 95 combined effect of multiple agents acting at a
 96 common site or through common mechanisms
 97 (U.S. EPA, 2002; NRC, 2009).

1 The assessment then selects an overall reference
 2 dose and an overall reference concentration for
 3 the agent to represent lifetime human exposure
 4 levels where effects are not anticipated to occur.
 5 This is generally the most sensitive
 6 organ/system-specific reference value, though
 7 consideration of study quality and confidence in
 8 each value may lead to a different selection.

9 **7.7 Confidence and uncertainty in the**
 10 **reference values**

11 The assessment selects a standard descriptor to
 12 characterize the level of confidence in each
 13 reference value, based on the likelihood that the
 14 value would change with further testing.
 15 Confidence in reference values is based on
 16 quality of the studies used and completeness of
 17 the database, with more weight given to the
 18 latter. The level of confidence is increased for
 19 reference values based on human data supported
 20 by animal data (U.S. EPA, 1994).

21 **High confidence:** The reference value is not
 22 likely to change with further testing, except
 23 for mechanistic studies that might affect the
 24 interpretation of prior test results.

25 **Medium confidence:** This is a matter of
 26 judgment, between high and low confidence.

27 **Low confidence:** The reference value is
 28 especially vulnerable to change with further
 29 testing.

30 These criteria are consistent with guidelines for
 31 systematic reviews that evaluate the quality of
 32 evidence. These also focus on whether further
 33 research would be likely to change confidence in
 34 the estimate of effect (Guyatt et al., 2008a).

35 All assessments discuss the significant
 36 uncertainties encountered in the analysis. EPA
 37 provides guidance on characterization of
 38 uncertainty (U.S. EPA, 2005a). For example, the
 39 discussion distinguishes model uncertainty (lack
 40 of knowledge about the most appropriate
 41 experimental or analytic model) and parameter
 42 uncertainty (lack of knowledge about the
 43 parameters of a model). Assessments also
 44 discuss human variation (interpersonal
 45 differences in biologic susceptibility or in
 46 exposures that modify the effects of the agent).

References

48 DHEW (1964) *Smoking and Health: Report of the*
 49 *Advisory Committee to the Surgeon General of*
 50 *the Public Health Service*. Public Health
 51 Service Pub. No. 1103.
 52 <http://profiles.nlm.nih.gov/ps/access/NNBB>
 53 [MQ.pdf](#).
 54 DHHS (2004) *The Health Consequences of*
 55 *Smoking: A Report of the Surgeon General*.
 56 http://www.cdc.gov/tobacco/data_statistics
 57 [/sgr/2004/index.htm](#).
 58 Guyatt GH, Oxman AD, Vist GE, Kunz R, Falck-
 59 Ytter Y, Alonso-Coello P, Schünemann HJ
 60 (2008a) GRADE: an emerging consensus on
 61 rating quality of evidence and strength of
 62 recommendations. *British Medical Journal*
 63 336: 924-926, <http://www.bmj.com/>
 64 [content/336/7650/924.full](#).
 65 Guyatt GH, Oxman AD, Kunz R, Vist GE, Falck-
 66 Ytter Y, Schünemann HJ (2008b) GRADE:
 67 what is “quality of evidence” and why is it
 68 important to clinicians? *British Medical*
 69 *Journal* 336: 995-998, <http://www.bmj.com/>
 70 [content/336/7651/995.full](#).
 71 Hill AB (1965) The environment and disease:
 72 association or causation? *Proceedings of the*
 73 *Royal Society of Medicine* 58: 295-300.
 74 <http://www.ncbi.nlm.nih.gov/pmc/articles/>
 75 [PMC1898525/](#).
 76 IARC (2006) Preamble to the IARC Monographs.
 77 <http://monographs.iarc.fr/>.
 78 IOM (2008) *Improving the Presumptive Disability*
 79 *Decision-Making Process for Veterans*.
 80 Washington: National Academy Press.
 81 http://www.nap.edu/catalog.php?record_id
 82 [=11908](#).
 83 NRC (1983) *Risk Assessment in the Federal*
 84 *Government: Managing the Process*.
 85 Washington: National Academy Press.
 86 http://www.nap.edu/catalog.php?record_id
 87 [=366](#).
 88 NRC (2009) *Science and Decisions: Advancing Risk*
 89 *Assessment*. Washington: National Academy
 90 Press.
 91 http://www.nap.edu/catalog.php?record_id
 92 [=12209](#).
 93 Rothman KJ, Greenland S (1998) *Modern*
 94 *Epidemiology*. Philadelphia: Lippincott
 95 Williams and Wilkins.

DRAFT MATERIALS FOR REVIEW ONLY – DO NOT CITE OR QUOTE

- 1 U.S. EPA (1986a) Guidelines for the Health Risk
2 Assessment of Chemical Mixtures.
3 EPA/630/R-98/002.
- 4 U.S. EPA (1986b) Guidelines for Mutagenicity
5 Risk Assessment. EPA/630/R-98/003.
- 6 U.S. EPA (1988) Recommendations for and
7 Documentation of Biological Values for Use
8 in Risk Assessment. EPA/600/6-87/008.
- 9 U.S. EPA (1991) Guidelines for Developmental
10 Toxicity Risk Assessment. EPA/600/FR-
11 91/001.
- 12 U.S. EPA (1994) Methods for Derivation of
13 Inhalation Reference Concentrations and
14 Application of Inhalation Dosimetry.
15 EPA/600/8-90/066F.
- 16 U.S. EPA (1996) Guidelines for Reproductive
17 Toxicity Risk Assessment. EPA/630/R-
18 96/009.
- 19 U.S. EPA (1998) Guidelines for Neurotoxicity Risk
20 Assessment. EPA/630/R-95/001F.
- 21 U.S. EPA (2000) Supplementary Guidance for
22 Conducting Health Risk Assessment of
23 Chemical Mixtures. EPA/630/R-00/002.
- 24 U.S. EPA (2002) A Review of the Reference Dose
25 and Reference Concentration Processes.
26 EPA/630/P-02/002F.
- 27 U.S. EPA (2005a) Guidelines for Carcinogen Risk
28 Assessment. EPA/630/P-03/001F.
- 29 U.S. EPA (2005b) Supplemental Guidance for
30 Assessing Susceptibility from Early-Life
31 Exposure to Carcinogens. EPA/630/R-
32 03/003F.
- 33 U.S. EPA (2006a) Approaches for the Application
34 of Physiologically Based Pharmacokinetic
35 (PBPK) Models and Supporting Data in Risk
36 Assessment. EPA/600/R-05/043F.
- 37 U.S. EPA (2006b) A Framework for Assessing
38 Health Risks of Environmental Exposures to
39 Children. EPA/600/R-05/093F.
- 40 U.S. EPA (2009) IRIS Process.
41 <http://www.epa.gov/iris/process.htm>.
- 42 U.S. EPA (2010) Integrated Science Assessment
43 for Carbon Monoxide. EPA/600/R-09/019F.
- 44 U.S. EPA (2011) Recommended Use of Body
45 Weight^{3/4} as the Default Method in
46 Derivation of the Oral Reference Dose.
47 EPA/100/R11/0001.
- 48 U.S. EPA (2012a) Benchmark Dose Technical
49 Guidance. EPA/100/R-12/001.
- 50 U.S. EPA (2012b) Advances in Inhalation Gas
51 Dosimetry for Derivation of a Reference
52 Concentration (RfC) and Use in Risk
53 Assessment. EPA/600/R-12/044.
- 54 U.S. EPA (2013) draft Handbook for IRIS
55 Assessment Development.
56 <http://www.epa.gov/iris/xxx>.
- 57 January 2013

58

1 Appendix C – Example of IRIS Program Direction to 2 Contractors

3 CONDUCTING DOSE-RESPONSE MODELING

4 *NOTE: This section addresses only dose-response modeling of animal bioassay data from standard*
5 *experimental designs. Analysis of animal data from complex experimental designs (for example,*
6 *repeated measurements on the same test subjects) and analysis of epidemiological studies require*
7 *specialized methods that are documented on a case by case basis.*
8

9 The IRIS Program’s Statistics Workgroup (SWG) has provided instructions for conducting
10 the majority of dose-response modeling of animal bioassay data using BMDS software. This was
11 written to provide unambiguous instructions to contractors regarding what the IRIS Program
12 requires as an initial analysis before reviewing results. Analysis is often iterative, because review of
13 results may suggest corrections or additional analyses, or prompt a second QA review of data. This
14 draft cannot provide completely detailed advice; modelers should consult a statistician familiar
15 with the relevant methods for detailed advice and for a review of work.

16 This section summarizes the process of assembling data for dose-response modeling from
17 animal bioassay toxicology studies and then conducting and reporting those analyses. It is intended
18 for use by the IRIS Program. It can also be included with a Statement of Work (Performance Work
19 Statement) to provide detailed directions to a contractor and can also serve as a review guide. .
20 The objective is to reduce errors, streamline operations, and ensure that analyses and model
21 selection are done consistently for each assessment.

22 The section is intended to describe a common core of best practices, but is not necessarily
23 exhaustive and may be modified in the future. In addition, Chemical Managers and data analysts
24 within EPA may modify these practices to suit the particular needs of an assessment or the features
25 of a particular data set. This section assumes that the user has a good working knowledge of EPA’s
26 Benchmark Dose Software (BMDS).

27 The selection and justification of a final point of departure (POD) is the responsibility of
28 EPA and not to be delegated to a contractor, although contractors may be asked to make
29 recommendations on a wide range of topics and to draft related justifications. This section is
30 intended to be applied in conjunction with detailed guidance, including EPA’s *Benchmark Dose*
31 *Technical Guidance* (EPA/100/R-12/001, June 2012).

32 Any potential problems or questions that the analyst or modeler identifies should be
33 brought to the attention of the chemical manager as soon as possible. For example, there may be
34 questions about which data sets should be used for dose-response modeling (i.e., suitability, study
35 quality, biological significance, and significance of the response to dose), problems with data quality
36 or missing data, or difficulties encountered while modeling the data. Decisions to include or
37 exclude studies for DRA are significant and should be reviewed by the chemical manager or by EPA
38 staff designated by the chemical manager. The IRIS Program’s Statistics Workgroup (SWG) can
39 provide advice on implementing this section and on complicated or non-routine statistical analysis
40 issues.
41

1 **A. Preparing Data for Dose-Response Modeling**

2 **The steps in this section, especially steps 1 and 2, will be most efficiently accomplished**
3 **during preparation of data for hazard evaluation. When a number of studies need to be**
4 **compared, use of common units and allowance for different dosing regimens as**
5 **described below will facilitate the exercise.**

6 **1. For continuous response data, verify quantities described as either standard deviations**
7 **or standard errors.** BMDS requires standard deviation (S), which characterizes variability in
8 responses among individual animals. Some investigators report the standard error of the group
9 mean, which is S/\sqrt{n} . Some studies may report confidence limits. To calculate S from these
10 data, you will need to know the confidence level and method used. Have a statistician review
11 the report and your calculations. Document your assumptions and calculations.⁶

12 **2. Convert doses to standard units.**

13 a. Standard units for oral and inhalation exposures are mg/kg-day and mg/m³, respectively. If
14 the original study does not provide these data, the data analyst will need to apply the
15 necessary assumptions (especially those regarding body weight and food or water
16 consumption) and must document the necessary calculations.

17
18 **Best Practices for Data Management to Support Modeling**

19 **All data should be proofed against the original cited source before being used**
20 **in any evaluations, comparisons, or analyses. This should be done using**
21 **double inspection after data entry or by double entry followed by machine**
22 **comparison.**

23 **Permanently document all calculations and conversions, using a database or**
24 **spreadsheet to record the individual terms, used to do calculations and to**
25 **produce the adjusted dose or concentration.**

26

27 b. For oral exposures in food or water, exposure is treated as if continuous over 24 hours, 7
28 days per week. The source publication will usually report the oral intake of the chemical in
29 mg/kg-day. If not, it will be necessary to calculate this from study-specific data on intake
30 and body weights or to make inferences specific to the species, strain, and sex (e.g., U.S. EPA,
31 1988). For oral exposures by gavage, the entire dose is assumed to be distributed across 24
32 hours (in effect, assuming 24-hour dosing), so the most common adjustment is for
33 days/week of dosing (typically 5 days per week).

34 c. One may need to obtain animal body weights applicable to the dose-response data. Use data
35 from the study, if available. If body-weight data are presented only in figures, numerical
36 data can be recovered by digitizing the figures. Average body weight should be calculated
37 across the period of exposure, for the purpose of calculating exposure per unit body weight.

⁶ BMDS provides a "Transformation" option to convert SE to SD.

1 If there are no study-specific data, there are sources of default body weights for adults of
2 each sex for various strains of rodents (e.g., U.S. EPA, 1988).

- 3 d. For the purpose of animal to human extrapolation of oral exposures, the IRIS Program uses
4 a standard weight for each species, regardless of the availability of study-specific
5 information, as recommended by Agency guidance (U.S. EPA, 2011).
- 6 e. For inhalation exposures, RfCs are typically expressed in mg/m³ by multiplying the POD in
7 ppm by [(Molecular Weight)/24.45]. This conversion can be made before or after dose-
8 response modeling. RfCs may be extrapolated to humans using RfC methodology (U.S. EPA,
9 1994). This conversion can also be made before dose-response modeling, or after if there
10 were no substantive differences in breathing rate with exposure level.
- 11 f. The foregoing may not apply if you are using internal dose metrics from a PBPK model.
12 Instructions for such analyses are beyond the scope of this draft.

13 **3. Make standard dose adjustments to account for intermittent and less-than-lifetime**
14 **exposures.** Report the individual terms and final multiplier employed in the adjustments, in
15 the data summary table in the modeling appendix of the assessment.

- 16 a. This section may not apply if you are using internal dose metrics from a PBPK model. The
17 PBPK dose metrics may partially or fully account for intermittent exposure. The dose-
18 response analyst must communicate with the PBPK analyst regarding adjustments to dose
19 metrics.

20 b. Cancer bioassays:

21
$$\text{dose} * (\text{hours/day})/24 * (\text{days/week})/7 * (\text{weeks exposed} / \text{weeks on study})$$

22 If sacrifice time is less than the standard 'lifetime' (104 weeks for rats and mice), also
23 multiply by [(weeks on study before sacrifice / 104)³] (see, e.g., Portier et al., 1986).

24 Example: 8-week-old mice were exposed by inhalation for 6 hours per day, 5 days per
25 week, for 78 weeks, and were sacrificed at 91 weeks after starting exposure; the adjusted
26 exposure concentration is (nominal exposure concentration) (6/24) * (5/7) * (78/91) *
27 [(91 wk / 104 wk)³].

28 c. Noncancer endpoints (in cancer bioassays or other studies):

29
$$\text{dose} * (\text{hours/day})/24 * (\text{days/week})/7 * (\text{weeks exposed} / \text{weeks on study})$$

30 If exposure lasted until the final sacrifice, the last term has no effect on the conversion.
31 Unlike with cancer bioassays, exposures are not extrapolated to longer durations without
32 chemical-specific data.

- 33 d. Time-varying exposures: These dosing regimens require calculation of a time-weighted
34 average dose or concentration, before applying the factors above.

35 Example: a study applied one dose, D1, for weeks 1-12 and another, slightly different dose,
36 D2, for weeks 13-78 (and nothing thereafter to week 91). The time-weighted average dose
37 is (12 weeks/91 weeks)*D1 + (66 weeks/91 weeks)*D2 + (13 weeks/91 weeks)*0.

1 **4. If survival rates differ substantially among dose groups for cancer bioassays**

- 2 a. **If data are available for individual animals on time of death:** Assemble data on
3 individual times of death and tumor incidence for use in time-to-tumor modeling (described
4 below).
- 5 b. **If only grouped data are available:** Estimate the number of animals at risk in each dose
6 group using the number alive at 52 weeks (if the first tumor was observed later than 52
7 weeks) or the number alive at the week when the first tumor was observed. Note the use of
8 an adjusted number at risk when reporting the DRA. This adjustment is not as effective as
9 using individual animal data for survival adjustment, so some bias in estimates is to be
10 expected. The results of DRA must be qualified by noting the possible inaccuracy (bias)
11 caused by incomplete survival adjustment.
- 12 c. If no survival adjustment is possible, results of DRA must be qualified by noting the possible
13 inaccuracy (bias) caused by lack of a survival adjustment.

14 **5. If survival rates differ substantially among dose groups for noncancer effects:** The
15 foregoing methods have not been used in IRIS assessments for noncancer effects, but they could
16 be. If there was severe early mortality (either different or similar among dose groups), this
17 should be called to the attention of the chemical manager (and noted in the assessment, if the
18 data are used).

19 **6. Sort/organize endpoints by type of health effect (i.e., different target systems). Seek the**
20 **Chemical Manager's advice on this. An example list:**

- 21 • kidney
- 22 • liver
- 23 • other organs
- 24 • body weight
- 25 • neurological
- 26 • immunological (includes thymus, adrenal)
- 27 • respiratory tract
- 28 • reproductive
- 29 • developmental (when summarized by dam; excluding nested data)

30 **Physiologically-Based Pharmacokinetic (PBPK) models**

31 PBPK models relate external exposures, also referred to as applied exposure, to internal
32 measures of concentration or exposure. Preferably a PBPK model will describe the concentration of
33 the active toxicant (often a metabolite of the material whose risk is being assessed) in the target
34 tissue. PBPK models may also describe early (precursor) molecular interactions, such as binding to
35 a receptor, inhibition of an enzyme, or formation of DNA adducts or cross-links. Even if a PBPK
36 model only describes the blood concentration of a toxicant, or rate of metabolic formation of the
37 toxicant (in the liver), these measures of exposure or dose are closer to and hence presumed to be

1 more predictive of a toxicant response than the applied dose. There are a number of potential
2 advantages to using a PBPK model:

- 3 • Since the internal dose metric is mechanistically closer to the toxic response,
4 subsequent dose-response modeling (BMD modeling in particular) should more
5 accurately interpolate among the dose-response data and better characterize
6 uncertainty in that relationship.
- 7 • If a PBPK model is calibrated for a route of exposure for which toxicity data are not
8 available (e.g., toxicity data are only available for oral exposure but the PBPK model is
9 calibrated for both oral and inhalation exposure), then the model can also be used for
10 route-to-route extrapolation.
- 11 • Given PBPK models which are calibrated for both a test animal species and humans
12 should allow for a more accurate prediction of the human equivalent (exposure)
13 concentration (HEC) or dose (HED), since they use chemical- and species-specific data
14 to match the exposure associated with a toxicological point-of-departure in the animal
15 with the corresponding exposure in humans.
- 16 • Also, in conjunction with Monte Carlo sampling, if the distributions of PBPK parameters
17 are defined for a human population, then the model can be used to obtain a data-derived
18 uncertainty factor for human variability in pharmacokinetics (UF_{H,PK}).

19 However, it must be first noted that these benefits are dependent first on the PBPK model's
20 ability to predict a metric which is in fact mechanistically closer to the effect of concern. If toxicity
21 is caused by a metabolite for which the model has not been calibrated, or occurs in a portal-of-entry
22 tissue for which PK data are lacking, then other internal metrics that the model does predict may
23 not be more predictive of toxicity than applied dose. And because PBPK models are used to predict
24 internal doses under exposure conditions (e.g., over chronic periods) and at concentrations for
25 which direct data may not be available, there is higher uncertainty in the prediction than in directly
26 measured or known exposure data. To assure that these issues of applicability and certainty are
27 adequately addressed, PBPK models should be subject to careful scientific and quality reviews
28 before they are used. It is not possible to be as specific and proscriptive about when and how a
29 PBPK model should be used as is the case for the statistical dose-response modeling, because PBPK
30 modeling has not reached the level of maturity that exists for statistics, but also because the models
31 are much more compound-specific. Each model tends to contain chemical-specific aspects, and
32 each PK data set presents unique challenges for modeling. Because PBPK modeling therefore
33 involves a significant use of scientific judgment, it is more important for models to be reviewed by a
34 team of experts, a primary role of the Pharmacokinetics Work Group (PKWG). Having multiple
35 modelers consider a given model (or set of models) provides input from multiple experts with a
36 variety of backgrounds and perspectives.

37 The following items should be considered in reviewing a PBPK model and considering how
38 to apply it.

- 39 **(i)** Both the science and the model code must be evaluated. The science -- model structure,
40 equation forms, and hypotheses that these represent -- will be described or at least implied
41 in the corresponding publication or report. But a QA of the model code is also essential, to

1 assure that it accurately represents the science as described, correctly converts units, and
2 uses the correct parameters (again as listed in a paper or report). If a model cannot
3 describe all the data with a consistent set of parameters, or ones which vary in a predictable
4 or measurable way (e.g., respiration rate can be varied to match measured values), then that
5 calls into question its ability to predict kinetics for exposure situations (bioassays) for
6 which calibration data are not available.

7 **(ii)** The degree of mechanistic complexity should be evaluated against the available data.
8 This is a matter of professional judgment, specific rules can't be laid out, but a more
9 complex model may represent hypotheses that have not been adequately tested and hence
10 is not necessarily considered the best for use in the assessment. Likewise a metric which is
11 closer to a toxic endpoint (e.g., binding to a receptor in the brain) may have a higher degree
12 of uncertainty due to limited calibration/validation data than one which is intermediately
13 close (e.g., blood concentration).

14 **(iii)** Use of PBPK models for animal test species to replace applied dose with internal dose
15 metrics can improve the quality of the statistical modeling if it takes out (really explains)
16 some of the exposure-response nonlinearity. Therefore this approach is suggested, though
17 it is not necessary. Using PBPK-predicted internal doses in the dose-response modeling
18 makes it harder to update an assessment if the PBPK model is changed or there is an update
19 in the analyzed dose-response data, so conducting the BMD modeling using applied doses is
20 an acceptable alternative. Using PBPK-derived internal doses should be most considered,
21 though, when the statistical BMD models have trouble fitting dose-response data.

22 **(iv)** When toxicity data are available for multiple exposure routes (i.e., inhalation, oral, or
23 dermal), ideally a PBPK model can explain apparent route sensitivity differences.
24 Specifically, the dose-response may appear discrepant when water ingestion/inhalation
25 rates are used, but become aligned when an internal metric is used. This is mentioned in
26 the current draft as a possible means of combining data sets and it is suggested that this
27 possibility be evaluated where possible. However the metric will depend on how well the
28 model captures any portal-of-entry/first-pass effects. Predictions for oral exposure of some
29 compounds have been found to depend strongly on the assumed drinking water pattern.
30 Oral ingestion is typically not continuous, and a bolus exposure can saturate metabolic
31 processes when a continuous exposure to the same total dose would not. So some thought
32 should be put into providing realistic oral ingestion patterns. The PKWG can provide
33 guidance on this.

34 **(v)** When there is a background or endogenous level for a material being considered, some
35 adjustment should be made to account for the fact that animals or people with no
36 exogenous exposure will still have some internal dose. A simple approach that may be
37 considered is to simply subtract the background levels from PK data and then calibrate a
38 PBPK model by fitting the resulting background-subtracted data. However this approach
39 implicitly assumes that the background is additive and that it is constant. Since PBPK
40 models are particularly useful when the PK are nonlinear, the first assumption in particular
41 may be inconsistent with the model and background subtraction can also distort the
42 apparent linearity or non-linearity in dose-response data. Incorporating a background term

1 into a PBPK model may make the model somewhat more complex, but it will allow for a
2 more accurate and transparent description and analysis of the dose-response relationship.

3 **B. Conducting Dose-Response Modeling**

4 In general, follow the advice on dose-response modeling in EPA's *Benchmark Dose Technical*
5 *Guidance* (EPA/100/R-12/001, June 2012), aka BMD-TG. The instructions that follow assume data
6 from 'chronic' studies, and rely mainly on use of EPA's Benchmark Dose Software
7 (<http://www.epa.gov/ncea/bmds/>). After the data and endpoints are selected (previous section),
8 these general principles apply:

- 9 ○ identify important or unusual statistical issues
- 10 ○ model all orders of multistage and polynomial models, up to and including the number
11 of dose groups minus one.
- 12 ○ select a best-fitting model using model selection procedures at BMD-TG §2.3.9
- 13 ○ the decision flow chart at BMD-TG §2.5 is a useful guide
- 14 ○ IRIS assessment modeling usually aims at getting a lower confidence limit (BMDL) for
15 the benchmark dose (BMD). For that end, a profile-likelihood interval or a Bayesian
16 interval is preferred to a Wald interval because the latter is less accurate for BMD.
17 BMDs uses profile-likelihood (BMD-TG 2.3.8). Most commercial statistical software (e.g.,
18 SAS) will report Wald intervals. Profile intervals can be obtained using custom
19 programs in SAS, R, and other software.
- 20 ○ The POD may be below the lowest non-zero dose. Confidence in the inference will
21 depend upon the degree of extrapolation, BMD/BMDL ratio, the type of model, the
22 model fit, and what is known about the chemical and endpoint, including the probable
23 MOA.

24 Advice to consider when initial modeling attempts are unsuccessful or the results are highly
25 uncertain (e.g., poor model fit, large confidence intervals, large differences between models).

- 26 ○ Non-standard approaches (e.g., adding parameter constraints) or additional models
27 (not included in BMDs) might then be considered (consult a statistician).
- 28 ○ Lack of fit might be evident for high doses (especially if the response decreases at the
29 high dose, which can be owed to mortality or to a change in the response pattern or
30 mechanisms at high doses). In some cases, the high dose group(s) may be omitted
31 (BMD-TG §2.3.6) after an unsuccessful attempt at modeling.
- 32 ○ Model fit might be improved by using internal dose measures based on toxicokinetic
33 modeling, if available.
- 34 ○ If adequately-fitting models differ greatly in BMDLs, and an adequately-fitting model
35 has a very large ratio of BMD/BMDL, the data do not permit accurate estimation of a
36 BMD and the data may be judged "not amenable to modeling." These facts should be
37 reported. An alternative is to use LOAEL or NOAEL as a basis for the POD. The response
38 and its confidence interval should always be reported with the LOAEL or NOAEL.
- 39 ○ other advice on improving model fit is given at BMD-TG §2.3.6

1 General advice concerning hypothesis tests:

- 2 ○ Calculating a LOAEL or NOAEL from bioassay data (in most cases this is available from the
3 source publication): seek advice from a statistician. Various methods are available,
4 depending on the properties of the data. Individual observations will be needed unless the
5 distributional properties of the data are well known from other studies. In general, the
6 method should account for multiple comparisons.
- 7 ○ Trend tests: seek advice from a statistician. The SWG is developing more detailed advice.
8 There are various methods, having different strengths and limitations.

9 **C. Modeling Cancer Endpoints (Single Tumor Sites)**

10 Three approaches may be available: (1) a biologically based dose-response model, (2)
11 time-to-tumor modeling, using survival times for individual animals, (3) modeling
12 cumulative incidence of cancer in dose groups.

13 A biologically based model could be used if one is available and deemed appropriate. This
14 approach requires specialized knowledge and custom software, and is not discussed
15 further.

16 If dose groups differed substantially in survival, and if data for individual animals are
17 available, time-to-tumor modeling or survival-adjusted quantal modeling is appropriate
18 (see below).

19 Incidence of cancers, grouped by dose, may be modeled using BMDS, as described further
20 below. This approach currently accounts for the majority of cancer dose-response
21 modeling.

22 Modeling cancer incidence:

- 23 ○ If dose groups differ substantially in survival, adjustments to the number of animals at
24 risk (below) will be needed
- 25 ○ Adenomas and carcinomas are combined (i.e., counting the number of animals with
26 either adenomas or carcinomas) when they arise from the same cell type and the
27 adenomas are believed to progress to carcinomas (U.S. EPA, 2005a)⁷
- 28 ○ When low-dose linearity is expected, current IRIS practice is to use the “cancer model”
29 in BMDS (Gehlhaus et al. 2011)⁸. The decision to use this approach must be made by
30 EPA.
- 31 ○ Apply cancer (multistage) models from the highest order model and all lower order
32 models down to first order. E.g., given 4 dose groups, fit models of 3rd, 2nd and 1st order.
33 Selection among these models is usually based upon minimum AIC (Akaike Information
34 Criterion)

⁷ U.S. EPA Guidelines for Carcinogen Risk Assessment (2005a), Section 2.2.2.1.2. *Statistical considerations*, p.2-19, states: “Statistical analysis of a long-term study should be performed for each tumor type separately. The incidence of benign and malignant lesions of the same cell type, usually within a single tissue or organ, are considered separately but may be combined when scientifically defensible (McConnell et al., 1986).”

⁸ The BMDS “cancer model” is a multistage model with non-negative coefficients, and it reports the “cancer slope factor” or potency, i.e., the ratio BMR/BMDL.

- 1 ○ If the “cancer model” does not fit adequately—i.e., if p is not greater than 0.05 (BMD-TG
2 Sections 2.3.5 and 2.3.9)— fit other BMDS models and select the best-fit model, as for
3 noncancer.
- 4 ○ If low-dose non-linearity is expected based upon mode-of-action information (U.S. EPA,
5 2005a), fit the full suite of BMDS dichotomous models to the relevant precursor effect
6 data and then select the best-fit model as for noncancer. The decision to use this
7 approach must be made by EPA.

8 **D. Combined Cancer Risk (multiple tumor types in a single animal study)**

9 When there is increased risk from multiple tumor types (sites), when the tumor types can
10 be assumed to be approximately independent, and when no single type substantially
11 dominates risk (e.g., by over 10-fold), then composite (“total”) cancer risk estimates are
12 derived by combining risk across tumor types (NRC 1994). Evaluating composite cancer
13 risk by modeling data for the total incidence for all cancers (incidence of “tumor-bearing
14 animals”) is not preferred (NRC 1994).

15 BMDS provides a ‘multitumor’ option for modeling composite risk. This does not refer to
16 multiple tumors in one animal. It refers to multiple types of tumors. For this model,
17 incidence is measured by the number of animals exhibiting a type of tumor, not by
18 counting numbers of a certain tumor in each animal.

19 The BMDS ‘multitumor’ option requires the use of a single, common dose metric.
20 Occasionally there is a need to estimate composite risk using different dose metrics for
21 different tumors (e.g., best-fitting cancer models are based on different dose metrics; e.g.,
22 one cancer is modeled using a BMDS quantal model and another is modeled using a time to
23 tumor model). For these and other special circumstances, consult the IRIS Program’s SWG
24 for alternative methods for estimating composite risk.

25 **E. Time-to-Tumor Analysis and Survival Adjustment**

26 Studies should be reviewed to identify those that may benefit from a survival adjusted
27 analysis. At present, we have no set criterion for the survival difference among dose groups
28 that would trigger such an analysis. However, if the incidence appears to plateau (at less
29 than 100% of animals) or decreases at a high dose, and survival is also lower in higher dose
30 groups, these methods should be used (either to augment or to replace standard BMDS
31 modeling of cumulative incidence).

32 Two approaches are available for conducting a survival-adjusted analysis, (a) using a time-
33 to-tumor model, and (b) calculating a survival-adjusted number at risk and then using a
34 standard BMDS quantal model. Both require data for each animal on time of death, and
35 both assume that tumors are observed ‘incidental’ to the cause of death. In most cases,
36 cause of death is not available for most animals. If there is a need to analyze data in which
37 death of each animal can be attributed to the tumor of interest (vs. some other cause),
38 consult a statistician.

1 **1. Time-to-Tumor Modeling**

- 2 ○ Use the “MSW” (multistage Weibull) program (at www.epa.gov/ncea/bmds) for time-
- 3 to-tumor analysis. A guide to using the program is available at the BMDS web site, as is
- 4 an external review document. Report any failures of the MSW program to solve the
- 5 BMDL.
- 6 ○ Evaluate estimates for all model orders between 1 and the number of dose groups
- 7 minus one. Select a model based on minimum AIC (Akaike Information Criterion).
- 8 ○ These programs do not report a goodness of fit statistic. Evaluating goodness of fit for
- 9 models based on interval-censored or current-status data is still an open issue (Lawless,
- 10 2002), but recent literature may point to solutions.⁹

11 **2. Survival-Adjusted Number at Risk (Poly-3 Method)**

- 12 ○ Calculate a survival-adjusted number at risk using the poly-3 method¹⁰
- 13 ○ Use the observed incidence and adjusted N in BMDS quantal models
- 14 ○ This method gave BMDs and BMDLs similar to those from the MSW model in a
- 15 limited number of comparisons.¹¹

16 **F. Modeling Noncancer Endpoints**

17 This section describes analysis of data from animal bioassays in which animals are

18 randomized to treatments and a single response is of interest. The randomization may

19 involve restricted or stratified schemes that ensure similarity among dose groups in

20 animal weights or other attributes. However, no repeated measures or cluster sampling is

21 involved. An average or other summary statistic of repeated measures might be used, but

22 its standard deviation must be correctly estimated.

23 **1. Fit all ‘standard’ BMDS models**

- 24 **a.** For continuous responses, apply the power, polynomial, Hill, and exponential
- 25 models in BMDS. All orders of polynomial model (up to order equal to the number
- 26 of dose groups-1) should be evaluated
- 27 **b.** For dichotomous responses, apply¹² the gamma, logistic, log-logistic, probit, log-
- 28 probit, multistage, and Weibull models¹³. All orders of multistage model (up to order
- 29 equal to the number of dose groups-1) should be evaluated

⁹ References are available on request.

¹⁰ See Piegorsch, W.W., and A.J. Bailer, 1997, *Statistics for Environmental Biology and Toxicology*, London: Chapman & Hall. Section 6.3.2, pp. 235-236. SWG can provide a template spreadsheet for poly-3 calculation or an R program that makes the calculation along with a survival-adjusted trend test. Note that use of a weighting method to produce “adjusted” numbers at risk (i.e., poly-3, causes the dose-response modeling assumption of binomially distributed observations no longer to be exact. Software is also provided at <http://www.jstatsoft.org/v16/i07>, "A Computational Tool for Testing Dose-related Trend Using an Age-adjusted Bootstrap-based Poly-k Test", by H. Moon et al. (2006) *J. Statist. Software* 16(7), 14 pages.

¹¹ It is not clear that either method should always be preferred. There may be a variance-bias trade-off (greater accuracy and lower precision for MSW) such that neither is uniformly ‘best’

¹² The “quantal-linear” and “quantal-quadratic” models are special cases of other models in BMDS; there is no need to use them routinely. Those models should only be applied in the unusual case that a mechanistic hypothesis (based on independent evidence and reasoning) supports one of these models.

1 **2. Identify a best-fitting model.** Apply model selection procedures and identify a best-fitting
2 model following EPA’s Benchmark Dose Technical Guidance §2.3.9

3 **G. Modeling Clustered Developmental Data**

4 [Specific advice for this type of data is being developed by SWG]

5 **H. Other Types of Data**

6 Other forms of data are not handled in BMDS and will require collaboration with a
7 statistician to ensure correct analysis and interpretation. Some of these are listed below and
8 briefly described.

9 **1. Multivariate Response Data**

10 Data collected as multivariate measurements on animals are usually reported and analysed
11 singly as univariate measures. There can be advantages to joint analysis, including a more
12 accurate estimation of risk and BMD (eg, Catalano et al., 1993, 1997; Dunson, 2000; Krewski
13 and Zhu, 1995; Najita et al., 2009; Ryan, 1992).

14 **2. Repeated Measures Data**

15 **a. Toxic-Diffusion Model for Time-Dependent Neurobehavioral Data**

16 The BMDS web site provides a toxico-diffusion model implemented in R (e.g., Zhu et
17 al., 2005a, b). The model and program are based on published work and are
18 intended for modeling dose-response for time-dependent (i.e. repeated measures)
19 neurobehavioral data from neurotoxicity experiments.

20 **b. Growth Models**

21 There is a large literature on growth models with repeated measures. Mixed effects
22 models (for example using NLMIXED procedure in SAS) are suitable for such
23 analyses. Those who need to apply such models should consult a statistician and
24 pertinent textbooks and monographs. Applicability of the toxico-diffusion model
25 (above) to such data has not been evaluated.

26 **3. Categorical Response Data**

27 The BMDS web site provides the “CatReg” model. This software or related methods
28 have been applied to risk estimation and meta-analysis in a variety of cases (eg, Brown
29 and Strickland, 2003; Krewski et al., 2010a, b; Simpson et al., 1996; Guth et al., 1997).

30 Modeling of categorical data is a well developed topic encompassing many statistical
31 models and methods described in a number of textbooks. It is unwise to proceed with
32 modeling categorical data without the assistance of a toxicologist and a statistician
33 familiar with these methods. Users may consult examples of categorical modeling
34 applied to risk estimation (Teuschler et al, 1999).

35 **4. Concentration-Time Data**

¹³ As of December, 2012, there remain some occasional convergence difficulties with the dichotomous Hill model. If it is plausible that response might plateau at < 100%, this model should be applied, but some trial and error (initializing parameter values) may be needed occasionally to achieve convergence.

1 Typically these data involve responses after short-duration exposures (i.e., acute
2 responses). The BMDS web site provides software implementing the Ten Berge CT
3 model. It is unwise to proceed with modeling such data without the assistance of a
4 statistician familiar with these methods. Users may consult examples from the
5 literature as a starting point (Brown and Foureman, 2005).

6
7 **References**

- 8 Brown, K.G. and G.L. Foureman, 2005, Regul. Toxicol. Pharmacol. 43: 45-54.
9 Brown, K.G. and J.A. Strickland. 2003, Regul. Toxicol. Pharmacol. 37:305-317.
10 Catalano, P.J., 1997, Statistics in Medicine 16:883-900.
11 Catalano, P.J. et al, 1993, Teratology 47: 281-290.
12 Dunson, D.B., 2000, Risk analysis 20:429-437.
13 Guth et al., 1997, Risk Analysis 17: 321-332.
14 Krewski, D. and Y. Zhu, 1995, Risk Analysis 15:29-39.
15 Krewski et al., 2010, J. Toxicol. Environ. Health Part A 73:187-207.
16 Krewski et al., 2010, J. Toxicol. Environ. Health Part A 73:208-216.
17 Lawless, J.F. 2002. Statistical Models and Methods for Lifetime Data. 2nd ed. Wiley-Interscience.
18 NRC (National Research Council). 1994. Science and Judgment in Risk Assessment. Washington,
19 DC: National Academy Press [Chapter 11, Appendix I-1, Appendix I-2]
20 Najita, J.S., et al., 2009, Applied Statistics 58:555-573
21 Ryan, L., 1992, Biometrics 48:163-174.
22 Simpson, D.G. et al., 1996, J. Agric. Biol. Environ. Stat. 1:354-376.
23 Teuschler, L.K. et al., 1999, Regul. Toxicol. Pharmacol. 30 (Supplement):S19-S26.
24 U.S. EPA. (1988) Recommendations for and Documentation of Biological Values for Use in Risk
25 Assessment. U.S. Environmental Protection Agency, Washington, D.C., EPA/600/6-87/008 (NTIS
26 PB88179874). <http://cfpub.epa.gov/ncea/cfm/recordisplay.cfm?deid=34855>.
27 U.S. EPA Guidelines for Carcinogen Risk Assessment (2005). U.S. Environmental Protection Agency,
28 Washington, DC, EPA/630/P-03/001F, 2005.
29 Zhu, Y., 2005, Environmetrics 16:603-617.
30 Zhu, Y., et al., 2005, Regulatory Toxicology and Pharmacology 41: 190-201 and 240-255;

31

1 Appendix D – Information Management Tool: Comment 2 Tracker Database

3 During 2012, the IRIS Program as part of the CAST initiative identified a need for better
4 documentation and communication of decisions. To address this need, an information management
5 tool that allowed for recording, reviewing, responding to, and analyzing comments and responses
6 was determined to be of value to the IRIS Program. To this end, two databases were developed to
7 (1) facilitate the analysis of, and response to, comments received during the course of developing
8 an assessment and review, and (2) allow comparison of comments and recommendations as well as
9 Agency responses made in multiple assessments.

10 The first database consists of a MS Access 2007 database with a form designed to streamline data
11 entry (Figure C-1). The form collects the following information fields on a given comment:

Database ID #	Overarching Issues*
Charge Question ID (<i>if relevant</i>)	Reviewer Agreement with EPA*
Verbatim Charge Question (<i>if relevant</i>)	Assessment Team Response/Level of Effort*
Reviewer	Revisions to Toxicological Review
Topic*	Response to Comment Appendix Location (Pg # and Charge Question)
Stage at which Comment was Received*	Official Response to Comment
Verbatim Reviewer Comment	Individual Addressing Comment
Summary of Reviewer Points/Recommendations	Completion Date
Major Comment*	Type of Review*

12 *Fields contain a limited number of options to facilitate comparison across chemicals.
13

14 Some fields contain a limited number of options for the user to choose from; for example, the topic
15 field allows the user to link comments with various sections of the assessment, or with broader
16 topics that may not be limited to a specific place in the document. The form also includes a means
17 of navigating the list of records as well as adding records for individuals with less experience with
18 MS Access. Additionally, pre-defined templates have been created that generate reports for
19 different purposes (e.g., project management for the chemical manager, CAST review of comments
20 received in peer review). Further templates will be developed as needed. Finally, the database
21 contains a query function (see Figures C-2 and C-3), allowing the user to look for specific comments,
22 or patterns of comments using different selection criteria. For example, a user can determine if
23 certain comments are repeated at multiple points during assessment development, or restrict a
24 search on a given term to a specific stage or type of review (e.g., limit the search to peer review
25 comments only).

26 The Comment Tracker Database provides a number of benefits to the IRIS Program, including:

- 27 • **The database serves as a quality-control tool.** With the flexibility to track public
28 comments received during the life of assessment development as well as formal comments
29 received during peer review, the database ensures that comments received from the public or
30 external peer reviewers are adequately considered.

- 1 • **The database facilitates project management.** Chemical managers will now have a tool for
2 efficiently assigning initial responsibility for addressing comments to assessment team
3 members and clearly defining roles and responsibilities of team members on an assessment.
- 4 • **The database simplifies management review.** For example, as part of the post-peer review
5 CAST process, assessment teams provide reports to their CAST teams detailing their response
6 to comments, including estimations of feasibility and level of effort required to address a
7 comment. Comments that would require a significant level of effort to address can quickly be
8 identified and decisions made in consultation with management on resource allocation to
9 resolve a scientific issue raised by a reviewer/commenter.
- 10 • **The database promotes a deeper review of comments.** The ability to sort, limit, and query
11 the full text in the database encourages team members to look for broader issues raised by
12 commenters. For example, the recurrence of a comment at multiple stages of review, even if
13 previously addressed by the assessment team, may indicate areas in a document where
14 further attention is warranted, or where the clarity of a section of the document could be
15 improved. Alternatively, comments that are positive towards a particular analysis or issue
16 presented in an assessment may be instructive to other assessment teams or for team
17 members in their work on other assessments.

18 The IRIS Program does not anticipate that using the database will significantly alter the length of
19 time it takes to complete an assessment. Entering information into the database is unlikely to take
20 longer than current methods used to report comments and Agency responses. After an assessment
21 is finalized and posted, further modifications to that database will be restricted. A copy of the
22 database will remain with the files specific to that assessment, while another copy will be reserved
23 for use in the cross-chemical database.

Comment Data Entry Form

Database ID: (New)

Charge Question ID:

Verbatim Charge Question:

Reviewer:

Topic:

Stage at which comment was received:

Verbatim Reviewer Comment:

Summary of Reviewer Points/Recommendations:

Major Comment:

Overarching Issues:

Reviewer Agreement With EPA:

Assessment Team Response / Level of Effort:

Revisions to Tox Review:

Rtc Appendix Location - Pg # and Charge Question:

Official Response to Comment:

Individual addressing comment:

Completion Date:

Type of Review

- Public Comment
- Agency Review
- Interagency Review
- External Peer Review
- SAB Review
- NAS Review

Record Management

Add Record Save Record

Print Record Delete Record

Clean Records

Record Navigation

Next Record Previous Record

First Record Last Record

Return to Welcome Screen

Record: 32 of 32 No Filter Search

1
2 **Figure C-1. Screen capture of Comment Data Entry Form. Drop-down menus contain pre-defined lists of content**
3 **to facilitate management review and searchability of the database.**
4

DRAFT MATERIALS FOR REVIEW ONLY – DO NOT CITE OR QUOTE

Database Query Form

Charge Question ID:

Verbatim Charge Question:

Reviewer:

Topic:

Assessment Team Response / Level of Effort:

Revisions to Tox Review:

RtC Appendix Location - Pg # and Charge Question:

Official Response to Comment:

Individual addressing comment:

Stage at which comment was received:

Verbatim Reviewer Comment:

Summary of Reviewer Points/Recommendations:

Major Comment:

Overarching Issues:

Reviewer Agreement With EPA:

Type of Review:

Select Fields to Show in the Report:

Charge Question ID

Verbatim Charge Question

Reviewer

Topic

Stage at which comment was received

Verbatim Reviewer Comment

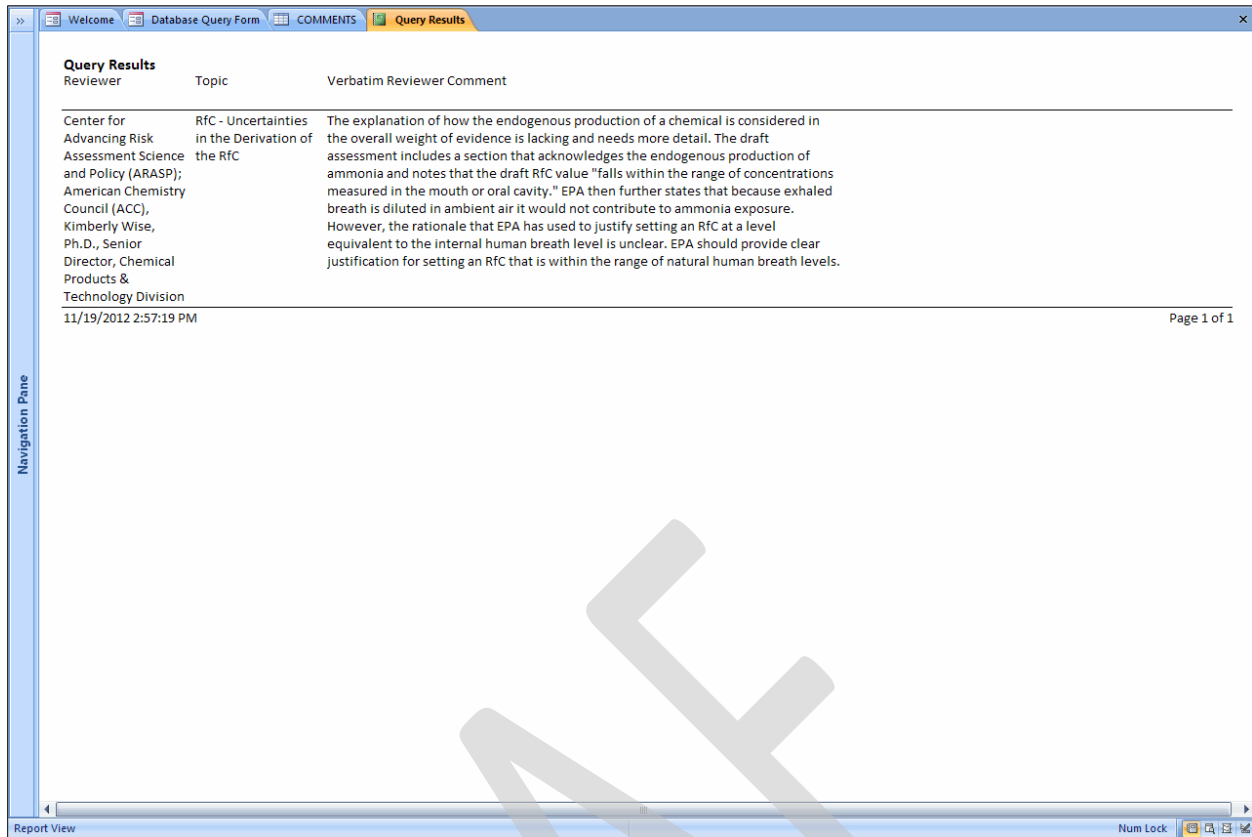
Search Return to Welcome Screen

Record: 1 of 1 No Filter Search

1

2 Figure C-2. Screen capture of the Database Query Form. Fields in the query form are largely the same as in the
3 Comment Data Entry Form to allow the user to define any search parameters they may find useful.

4



1
2 **Figure C-3. Screen capture of a database query result, where the user searched the public comments on the**
3 **draft ammonia assessment for use of the term “endogenous.” The search returned one result, with the user**
4 **generating a report with the Reviewer, Topic (i.e., section of the toxicological review associated with the**
5 **comment), and Verbatim Reviewer Comment fields.**
6

7 **Cross-Chemical Comparison Database**

8 Along with the chemical-specific comment tracker databases, the IRIS Program is developing an
9 additional database tool to query multiple chemicals. Also developed in MS Access, the query tool is
10 similar to the query tool used in the Comment Tracker Database; however this query form includes
11 an additional field for the user to select which chemical(s) they wish to query (Figure C-4).

12 The Cross-Chemical Comparison Database allows the user to search for comments on a specific
13 topic across chemical assessments. For example, if a chemical manager wanted to find all of the
14 peer review comments related to a specific mode of action, a simple search would identify all of the
15 relevant comments across all of the assessments in the database. The search results are based on
16 information entered at the time the comment was added to a database. Users can conduct a full-
17 text search within multiple fields, including the verbatim reviewer comment, Agency response to
18 comment, and other fields where the responses can be variable dependent on user input. Figure C-
19 5 shows an example output selecting from public comments received on the draft IRIS assessments
20 for ammonia and trimethylbenzenes.

21 There are several benefits to the cross-chemical database, some of which overlap with the
22 chemical-specific database. Similar to the chemical-specific databases, the cross-chemical database

1 is expected to aid staff during assessment development as well as management review. For
2 example, the cross-chemical database will:

3 • **Help staff identify how scientific issues have been addressed in different assessments.**

4 Early in draft development, chemical managers and assessment teams identify major
5 scientific issues that will need to be addressed in the assessment. The cross-chemical
6 database will allow staff to identify how these issues were addressed in other chemical
7 assessments, as well as how the public and/or peer reviewers responded to EPA's analysis of
8 that scientific issue. Understanding how a scientific issue was considered during peer review
9 or public comment can highlight important points for chemical managers and assessment
10 teams to consider early in draft development. Recognizing and considering these issues
11 earlier in draft development help reduce the time needed to complete an assessment.

12 • **Simplify post-hoc analysis of comments across chemicals.** A searchable database will
13 allow management and staff to review comments received across assessments and identify
14 when issues are being raised consistently by single or multiple reviewers. These comments
15 may point to cross-cutting issues that warrant further attention.

16 • **Allow for real-time searches during the review and public comment process.** During
17 any of the review or public comment steps (Agency/Interagency, public comment, peer
18 review), the database will allow scientists to identify, in real time, comments received on a
19 topic at different steps in the IRIS process and look for consistency or inconsistencies. For
20 example, if a public comment contradicts advice received during a peer review meeting, that
21 inconsistency can be brought up during the meeting to ensure that peer reviewers are aware
22 of previous arguments, and if needed, have further discussion on the issue.

DRAFT MATERIALS FOR REVIEW ONLY – DO NOT CITE OR QUOTE

1

2 **Figure C-4. Database Query Form for multiple chemical comparisons. The query form is identical to the query**
3 **form in the Comment Tracker Database except for the inclusion of a field labeled “Chemical,” which allows the**
4 **user to select specific chemicals for comparison.**

Reviewer	Topic	Verbatim Reviewer Comment	Chemical
ACC	Literature Search Strategy & Study Selection	EPA's draft assessment includes an explanation of the literature search strategy and study selection criteria the Agency used to identify studies for inclusion in the draft assessment. Table LS-1 provides the search parameters and terms used to identify studies; Figure LS-1 provides a schematic for how the Agency narrowed the available scientific literature. Table LS-1 provides useful and sufficient detail and should be maintained, in its current form, in future toxicological reviews. We have included below several areas where the transparency of the literature search could be greatly improved: - Figure LS-1 needs to be further expanded to include more detailed information regarding the criteria EPA used to include or exclude studies from consideration in the ammonia assessment. For example, Figure LS-1 indicates that 220 human studies, 203 animal studies and 599 supporting studies were found for a total of 1022 studies which were considered for inclusion in the draft assessment. 781 of these studies were excluded for various reasons (e.g. inadequate exposure characterization) but no breakdown has been included regarding the number of studies that were excluded for each of the exclusion categories provided.	Ammonia
American Chemistry Council: Hydrocarbon Solvents Panel	Literature Search Strategy & Study Selection	II. EPA's utilization of a consistent and transparent procedure for identifying, selecting and evaluating appropriate studies for inclusion in the Draft IRIS Assessment is critical to ensure and maximize the document's quality. The Draft IRIS Assessment's classification as "influential" information requires its content to meet rigorous standards of consistency and transparency. In failing to employ a consistent and transparent procedure for identifying, evaluating and selecting appropriate studies for the Draft IRIS Assessment, EPA has not complied with its IQ guidelines and severely undermined the legitimacy and utility of the document. What follows is a number of suggestions that, if adopted, can bring the Draft IRIS Assessment into accordance with both OMB and EPA IQ guidelines and significantly improve the accuracy and value of the document.	TMBs
American Chemistry Council: Hydrocarbon Solvents Panel	Literature Search Strategy & Study Selection	INTRODUCTION The literature search strategy and study selection as presented is significantly flawed. The process by which some studies were either not considered or not used was not transparent, or consistently and reliably applied. The inclusion of the extensive body of published data available on C9 aromatic hydrocarbon solvents tested by inhalation under the TSCA Section 4(a) test rule (FR50 20662, 1985) would greatly enhance the database available on TMB isomers individually and address many of the uncertainties raised in the Draft IRIS Assessment.	TMBs

1

2

Figure C-5. Results of a multi-chemical query on the ammonia and trimethylbenzenes databases (public comments only). The user searched for the term “transparent*” in the verbatim text field and limited results to comments associated with the Literature Search Strategy and Study Selection topic/section of the toxicological review, and the query returned information on the Reviewer, Topic, Verbatim Comment, and Chemical. Use of the wildcard “*” insured that derivatives including “transparency” and “transparent” were both captured in the search. The search returned three comments, two for the trimethylbenzenes draft assessment and one for the draft ammonia assessment; all the comments were made by the American Chemistry Council.

3

4 **Summary**

5

The databases described above are currently being tested with assessments of varying size and complexity. Wider implementation of the databases is planned in January 2013. Once in place, the databases will serve as an important information management tool for assessment development as well as documentation and quality control of IRIS assessments.

6

1 **Appendix E – Scoping to Inform the Development of IRIS** 2 **Assessments**

3 The following provides considerations for scoping to inform assessment development. These
4 considerations are not to be construed as a standard process for scoping an IRIS assessment as the
5 process will likely evolve as the IRIS Program gains experience in this area.

6 **Purpose**

7 The primary purpose of the scoping process is to understand the needs of clients in EPA's program
8 and regional offices with regards to a chemical or group of chemicals addressed by an IRIS
9 assessment. The scoping process builds upon information developed during the process of
10 identifying chemicals for the IRIS agenda which helps to outline clear objectives for the IRIS
11 assessment by defining and clarifying hazard identification and dose-response needs. During the
12 scoping phase, more detailed questions are asked to seek greater understanding of the specific
13 needs of the client offices.

14 **Process**

15 Scoping involves gathering specific information from EPA's program and regional offices (through
16 either a meeting or other communications) before beginning an IRIS assessment. This allows client
17 offices to identify their needs by explaining the environmental issues they need to address and
18 what type of information is needed in a hazard identification and dose-response assessment to
19 inform the decisions they will need to make. This exchange helps the IRIS Program understand the
20 types of information needed and the level of detail necessary to address client needs.
21 Understanding the clients' timelines is also important and may be factored into decisions about the
22 scope of the IRIS assessment. The following provides examples of the types of questions that might
23 be asked during the scoping process. It is important to note that these questions focus on the
24 "what" rather than the "how" of developing an IRIS assessment.

- 25 • What are the environmental issues and the types of decisions that will have to be made? If
26 there is more than one client interested in the IRIS assessment, do their decisions have different
27 scopes?
- 28 • What risk assessment activities, if any, have been carried out for this chemical by EPA's
29 program and regional offices, and other State and Federal stakeholders? Are there experts
30 within these organizations with whom the IRIS Program can consult?
- 31 • What routes of exposure are of most *a priori* concern (e.g., oral, dermal, inhalation)?
- 32 • What form(s) of the compound are most relevant for the clients' needs (e.g., elemental forms or
33 certain compounds for metals)? What ionic forms are of concern? In addition, should the IRIS
34 assessment consider the effect of counter ions on the toxicity of the ion in question (e.g., NO₃⁻ in
35 NH₄NO₃ or Br₂⁻ in PbBr₂)?

- 1 • Are there particular concerns regarding dose-response issues related to bioavailability of the
2 compound? How stable is the chemical in
3 physiological media? (If it readily
4 decomposes, decomposition products may
5 need to be considered in the IRIS
6 assessment.)
- 7 • Do other chemicals routinely co-occur or are
8 co-released with the compound(s) under
9 consideration? If so, what are they? Where
10 there is typically co-exposure, will the
11 potential regulatory decision address only
12 one chemical in the mixture, or the mixture as
13 a whole? Would it be useful to develop a
14 single IRIS assessment for the group of
15 chemicals? For example: a single chemical
16 (TCDD) versus all dioxin-like compounds.
- 17 • Does the decision-making or potential
18 regulatory action pertain to a group of
19 chemicals (such as substitutes for each
20 other)?
- 21 • Are specific lifestages and windows of
22 exposure of particular concern (e.g., children,
23 geriatric, in utero, perinatal, lactation)?
- 24 • What are the typical durations of exposures
25 of concern (e.g., acute, short-term,
26 subchronic, and chronic)? Do exposure levels
27 from scenarios of concern fluctuate
28 significantly over time (which might impact
29 the importance of short-term values to
30 evaluate peak levels)?
- 31 • Are there urgent or time-sensitive decision-
32 making needs faced by EPA's program or
33 regional offices? (This question will help the
34 IRIS Program address assessment deadlines
35 by weighing priorities among IRIS assessments and by possibly re-aligning, if necessary, the
36 scope of the IRIS assessment.)
- 37 • Would a conclusion on hazard identification without, or prior to, dose-response analysis be
38 useful?
- 39 • Is it critical to have dose-response information that enables cost-benefit analysis, with some
40 estimates of changes in health impacts between decision-making options?

Example of scoping: IRIS Assessment of Inorganic Arsenic

IRIS implemented a scoping process to inform the development of the inorganic arsenic assessment. Scoping meetings were held with EPA programs and regions as well as with interested stakeholders from the public and other federal offices. These meetings will inform the final planning and scoping statement for the inorganic arsenic assessment.

The following factors were discussed and identified by EPA clients at the arsenic planning and scoping meetings:

Hazard Identification Needs:

- Consideration of cancer and noncancer endpoints due to inorganic arsenic exposure.
- Inclusion of inhalation route in addition to oral.
- Consideration of exposure through occupational uses.
- Consideration of metabolites and oxidation state.
- Consideration of sensitive populations and lifestages: in particular children, and in utero and perinatal exposure. Evaluation of genetic and epigenetic factors affecting susceptibility, and impact of non-chemical stressors.

Dose-Response Needs:

- Impact of measures of exposure, bioavailability, and arsenic speciation.
- Risk at exposure to naturally occurring levels of inorganic arsenic.
- Estimation of risk beyond a reference concentration.
- Dose-response analyses that facilitate cost-benefit analyses.
- Impact of uncertainties in the dose-response analysis.
- Transparent presentation of choices made in the assessment and the supporting rationale.

- 1 • Is there a particular subpopulation already known to be of greater concern so that it is
2 especially important to understand that particular susceptibility?
- 3 • Do the EPA decisions involve occupational risks (which may be at ranges above what are typical
4 environmental exposures)?
- 5 • Is this a decision for which uncertainty and variability assessments might be particularly
6 important? What kind of uncertainty-variability information will best inform the decision-
7 making (taking into consideration the resource and time-intensive aspect of certain extensive
8 uncertainty and variability analyses)?

9 **Outcome**

10 The outcome of the scoping process is a statement that outlines the focus of the assessment, the
11 nature of the hazard characterization needed, and a clear indication of issues that are beyond the
12 scope of the IRIS assessment.

13 **Conclusion**

14 The scoping process is an evolving tool. Although some level of scoping takes place for every IRIS
15 assessment, the IRIS Program is just beginning to implement it as an early step in developing an
16 assessment. As the Program gains more experience in this area, standard procedures may be
17 developed. The IRIS Program needs to develop institutional experience and knowledge with the
18 planning and scoping process for several assessments before formulating standard procedures.
19 While face-to-face meetings may be necessary for some chemicals, email or other virtual
20 consultation may be sufficient in other cases.

21

1 **Appendix F – Draft Handbook for IRIS Assessment** 2 **Development**

3 The draft *Handbook for IRIS Assessment Development* provides information to IRIS
4 assessment teams regarding internal processes and evaluation steps used in the development of
5 IRIS assessments; however, the draft *Handbook* is a work in progress and currently does not fully
6 discuss each step in the IRIS assessment development process. The draft *Handbook* is designed to
7 provide the chemical assessment team with instructions and considerations involved in conducting
8 a literature search; screening for relevance to identify and select pertinent studies; evaluating and
9 documenting the quality of individual studies; reporting individual study results; synthesizing and
10 integrating evidence for epidemiological, toxicological, and mechanistic data; selecting studies for
11 derivation of toxicity values; considering combining data for dose-response modeling; managing
12 data for dose-response modeling; and selecting an organ/system-specific or overall toxicity value.

13 Components that remain missing from the working draft are integrating across evidence
14 (epidemiological, toxicological, and mechanistic data) to identify hazards and transition to dose-
15 response analysis; conducting dose-response modeling; extrapolating to lower doses and response
16 levels; considering susceptible populations and lifestages; developing candidate toxicity values;
17 characterizing confidence and uncertainty in toxicity values; and selecting final toxicity values.

1 Identifying and Selecting Pertinent Studies:

2 Literature Search and Screening

3 The focus of IRIS assessments is typically on the evidence of all types of health effects of a
4 particular chemical, and their exposure-response relationships. This is, by definition, a broad topic.
5 The systematic review process that has been developed and applied within the clinical medicine
6 arena (evidence-based medicine) is generally applied to narrower, more focused questions.
7 Nonetheless, the experiences within the clinical medicine arena provide a strong foundation for a
8 similar endeavor in IRIS assessments. (Some useful references describing systematic review within
9 clinical medicine are described at the end of this section.)

10 Systematic review, as applied in IRIS health assessments, is an iterative process that
11 identifies relevant scientific information needed to address key, assessment-specific questions. The
12 initial steps of the systematic review process formulate specific strategies to identify and select
13 studies relating to each key question, evaluate study methods based on clearly defined criteria, and
14 transparently document the systematic review process and its outcomes. The systematic review
15 process must be conducted in a way that protects from bias in study selection and evaluation by
16 transparently presenting all decision points and the rationale for each decision.

17 This section of the draft *Handbook* provides a discussion of the principles, overview of
18 methods, and points to consider as you go through the process of developing and documenting a
19 systematic review within the context of IRIS health assessments. It is not meant to be a “cookbook”
20 or a checklist of procedures. The topics covered are the first two steps in the systematic review
21 process: literature search and screening for relevance. Other steps involving evaluating the quality
22 of individual studies and evaluating and synthesizing data across multiple studies will be covered in
23 subsequent sections of the draft *Handbook*.

24 Throughout this draft *Handbook*, the examples provided are meant to be generalizable and
25 are presented in a simplified form, such as would be expected for health assessments with a small
26 literature database. Chemical-specific examples that have been drafted or completed can illustrate
27 how these approaches may be envisioned for more complex datasets.
28

29 **STEP 1: LITERATURE SEARCH**

30 The strength of a systematic review of research rests on its ability to identify relevant
31 studies, both published and unpublished, pertaining to the question of interest (e.g., health effects
32 of a chemical). All search strategies balance competing needs for “sensitivity” (i.e., the ability to
33 identify all potentially relevant studies) and specificity (i.e., the ability to avoid identification of
34 non-relevant studies), using a process that is both manageable and reproducible. The efficiency of
35 this process depends on optimizing the approaches used in initial searching and screening steps.

36 The goal of the search strategy is to identify full reports of **primary studies** (i.e., original
37 sources of data) pertaining to the key question(s). These studies can be published papers or
38 unpublished reports, but need to provide sufficient detail to allow evaluation of the study methods.

39 The initial search strategy “casts a wide net”; subsequent steps in the process are used to screen
40 and exclude articles that are not relevant, and to sort the relevant studies into categories (e.g.,
41 experimental studies in animals, observational studies in humans) for further evaluation.

1 In addition to the process described below, the IRIS Program takes other steps to identify
2 relevant studies that may have been missed by the formal search strategy. The IRIS Program
3 invites public review of the results of the literature search as one of the early steps in the
4 development of the assessment. In addition, more targeted requests for review, for example by
5 investigators active in the field of research, may also be a useful way to find studies that were not
6 otherwise identified. Studies that are identified through this process can be included in the
7 evaluation phase described in Study Quality Evaluation. It is also important to try to identify why
8 the studies were missed in the initial literature search (e.g., the chemical was not part of the
9 indexing terms used for that particular publication; publication not in any of the databases
10 searched) and consider if any modifications to the search strategy are warranted.

11 The following sections discuss the key steps in the literature search process: selecting
12 databases, selecting search terms, augmentation of a database search, documenting the search
13 strategy, and updating the literature search.

14 ***1A. Selecting Databases***

15
16
17 **Systematic Reviews Conducted For IRIS Assessments**
18 **Should Include Several (Types of) Databases,**
19 **Including Sources of Unpublished Studies**
20
21

22 The IRIS assessment team is responsible for devising and executing the literature search
23 and screening strategy, with assistance from EPA resources (e.g., HERO staff) and from contractors
24 as needed. If contractors are used, it is essential that EPA expertise guides the search process.

25 A search of PubMed is one way to start the process, to gain familiarity with the subject
26 matter, but by itself is not sufficient. Table F-1 describes the databases that serve as the foundation
27 for IRIS assessments. **PubMed**, **Web of Science**, and **Toxline** are overlapping databases of journals
28 focusing on medical and life science, science and social science, and toxicology literature,
29 respectively. These three databases are the core sources the IRIS Program uses for published
30 studies. Other databases may be useful for specific chemicals or questions, so IRIS is not limited
31 only to these core sources.

32 Another source of primary studies is the bioassays conducted by the National Toxicology
33 Program (NTP). Although some of these reports may be found through the searches of the
34 databases described above, it is also useful to directly search the NTP web site. This will assure you
35 that you have found reports that have not yet been published. (All of the NTP reports have
36 undergone external peer review, regardless of publication status.)

37 The category of unpublished studies can be quite broad, and could vary from anecdotal
38 reports to standardized animal bioassays conducted under established protocols by reputable
39 laboratories. These unpublished data are sometimes referred to as “gray literature.” Under the
40 federal Toxic Substances Control Act (TSCA), companies that manufacture, process or commercially
41 distribute a chemical are required to submit to EPA results of chemical testing and health and safety
42 studies. The **Toxic Substances Control Act Test Submissions (TSCATS)** database is a repository
43 of unpublished studies submitted to EPA under TSCA. Submissions from 1985 to 2004 can be
44 found through TOXLINE; subsequent submissions can be found through an EPA web site (see

1 TSCATS2, Table F-1). There is no requirement that these studies also be submitted for publication,
2 so this database may be the only source of the data contained in these studies. Reports in this
3 database should be included in the identification and evaluation process.

4 Another type of “gray literature” is conference proceedings and abstracts. IRIS assessments
5 generally do not include these types of publications as primary research literature because the level
6 of detail is insufficient to evaluate the methods. This group of references is kept as a separate
7 category to facilitate its evaluation, but in general a study that is only available in abstract form
8 would not be included in an IRIS assessment.

9 Chemicals used as pesticides must be registered for use in the U.S. by EPA’s Office of
10 Pesticide Programs (OPP). To support registration determinations, EPA requires more than 100
11 different scientific studies and tests from applicants, including toxicology studies. Searches of OPP
12 databases should be performed for chemicals that are also used as pesticides. OPP’s two main data
13 bases are PRISM Documentum and the Integrated Hazard Assessment Database (IHAD). PRISM
14 Documentum contains toxicology studies for all pesticides, nearly all of which are Good Laboratory
15 Practices (GLP) guideline studies. IHAD contains summaries of the studies and reviews (Data
16 Evaluation Records, DERs) which describe the studies and rank them for study quality and
17 guideline compliance. Access to PRISM Documentum and IHAD is limited to EPA employees with
18 FIFRA confidential business information access authorization, which requires one hour of online
19 training.

20 The primary options for conducting searches are 1) using the HERO interface to selected
21 databases, 2) directly searching databases, downloading citations into EndNote for review (and
22 eventual import into HERO), and 3) supervising the search process conducted by contractors. It is
23 possible that a combination of approaches will need to be used.

24 When using HERO for the search process, you will need access to the LitSearch, LitCiter, and
25 LitTagger functions. It is important to test the search string in each of the selected databases, select
26 the “no pdfs” option for this initial search, and include “tags” for each database as part of initial
27 project page set-up. These tags can be used to track the source(s) for each citation identified in the
28 search.

29 When doing direct searching of databases, you will need to take additional steps to
30 eliminate duplicate references after combining the results of more than one database. (HERO does
31 this automatically through the PMID number, but EndNote uses an algorithm that is more prone to
32 errors based on differences in punctuation or capitalization). Also, there is no easy method to keep
33 track of the source(s) of each individual citation. However, this option may provide greater
34 flexibility in searching than the HERO interface provides.

35 When working with contractors, you will need to review each of the key decisions (i.e.,
36 selection of databases, search strings, additional sources). You do not want to simply receive a
37 product; rather you want to be part of the process that creates the product by providing oversight
38 and technical direction in accordance with the procedures specified in the contract and task order.

Table F-1. Description of Core Databases For Primary Literature

Database	Description	Notes
PubMed^a	Approximately 5,600 medical, biology, and other life sciences journals (through MEDLINE), most back to 1966. www.pubmed.com	Uses Medical Subjects Headings (MeSH) terms
Web of Science^a	12,000 science and social science journals, back to 1970. Also includes conference abstracts. Maintained by Thompson Reuters. http://apps.webofknowledge.com	Can also do “forward” searching, i.e., searching for publications that cite a specified reference.
TOXLINE^a	Toxicology journals, including developmental and reproductive toxicology (DART), technical reports and research projects, and archival collections; back to 1965 (a few citations dating back to the 1940's); run by NLM. http://toxnet.nlm.nih.gov/cgi-bin/sis/htmlgen?TOXLINE	CASRN and synonyms
TSCATS	Toxic Substances Control Act Test Submissions Unpublished, studies submitted to EPA under TSCA Section 4 (chemical testing results), Section 8(d) (health and safety studies), Section 8(e) (substantial risk of injury to health or the environment notices) and FYI (voluntary documents). TSCATS is included in the TOXLINE database via HERO from 1985 through 2004; for submissions after 2004, use TSCATS2 at: http://yosemite.epa.gov/oppts/epatscat8.nsf/ReportSearch?OpenForm Or for recent 8E and FYI submissions, search: http://www.epa.gov/oppt/tsca8e/pubs/8eandfyisubmissions.html	Chemical name; CASRN Section 8(e) submissions most relevant
Office of Pesticide Programs (EPA)		
PRISM Documentum	Contains GLP guideline toxicology study reports for all pesticides from 1996 to present. Study reports older than 1996 can be acquired within a few days. Accessible to any EPA employee with FIFRA confidential business information access authorization. Go to: <ul style="list-style-type: none"> • OPP@Work - http://intranet.epa.gov/opp00002/ <ul style="list-style-type: none"> ○ OPP Applications (under popular sites in green box on left) ○ e-Registration Workflow (Documentum Login) 	
Integrated Hazard Assessment Database (IHAD)	Contains Data Evaluation Records (DERs; reviews of toxicology study reports), memoranda, cancer reports, metabolism reports, etc. for all of OPP. Accessible through Lotus NOTES database to any EPA employee with FIFRA confidential business information access authorization	

^a Accessible through HERO

1
2

1B. Selecting Search Terms

**Consult a Reference Librarian Early and Often
When Developing and Refining Your Search Strategy**

IRIS assessments are not limited to a single, narrowly defined study question. The focus of IRIS assessments is typically the evidence of toxicity or health effect (of any kind) of a particular chemical. The search strategy thus would generally begin by selecting the appropriate forms of the chemical name, CAS number, and if relevant, major metabolite(s).

The process of selecting the search terms should be done in close collaboration between the EPA assessment team and a reference librarian, either with HERO or with a contractor working on the assessment. Both of these resources offer extensive experience with database searching and information management. Correctly using limits in the search strategy, and correctly constructing complex structure using AND, OR, and NOT terms requires a high level of training and experience.

For some chemicals, the initial search using the chemical terms will yield an easily manageable reference base (e.g., < 1000 citations). In other situations, it may be necessary to refine the search strategy. For example, if the chemical is used as a positive control in certain assays, or used as an extraction solvent, you may find yourself in the position of looking for the few relevant studies among thousands of citations. This situation presents challenges from the standpoint of efficiency, as well as accuracy, because even with the various computer-based systems to facilitate screening titles and abstracts, reviewer fatigue (and subsequent error) is possible.

For very large databases, the searching and screening processes may also be improved by developing a series of secondary searches, each focused on a particular question (i.e., reproductive toxicity, cancer, pulmonary function). These more focused searches would result in smaller collections of citations that can more easily be reviewed by people with the appropriate scientific background.

There are situations in which expansion of general search terms within a category of effects is warranted. Review articles and other key documents should be consulted for information about specific types of effects that are of particular concern for the chemical under study. For example, for some chemicals, focusing the immune-related effects on allergic sensitization may make sense. For others, autoimmune effects may be most relevant. Male infertility may be a primary endpoint of interest within the category of reproductive effects for one chemical, and ovulatory disorders may be of more interest for another chemical. In each of these cases, more targeted development of search strings may be warranted.

Examples of secondary search terms used in conjunction with the chemical name terms to focus a search are shown in Table F-2. The first set is an example of terms that could be used to focus a search on pulmonary effects. The second set is an example of terms that could be used to exclude studies that use a chemical in certain types of assays (in this case, formaldehyde), but which are not studies of the effects of that chemical. NOTE: Be very careful when using "NOT" as a Boolean operator. An abstract that contains the term will be removed, even if other parts of the study are relevant to the primary focus of the search.

It is important to evaluate the set of studies that are excluded to see if the exclusion process was overly broad. For example, if this secondary filter eliminated 10,000 out of 12,000 initial references, a random sample of the excluded articles (small enough to be manageable, e.g., but large

1 enough to be representative, 100-200) should be selected and manually reviewed (based on titles
2 and abstracts) to determine if the “error rate” is acceptable; further refinements or additional
3 manual review may be necessary.

4
5
6 **The Literature Search Is Often an Iterative Process:
7 What You Learn From Reviewing the Results Should
8 Feed Back Into Ways to Expand or Refine the Strategy**
9

Table F-2. Illustrative Example of Strategies to Improve Literature Search Results

Chemical Terms (include synonyms and relevant metabolites)	Purpose of Additional Search	Example Secondary Terms For Refinement of Search
Beryllium and beryllium compounds	Focus on pulmonary toxicity	“Chemical” AND (alveolar OR BAL OR brochoalveolar OR “carbon monoxide” OR “chest pain” OR “chest tightness” OR cough OR crackles OR DL _{CO} OR dyspnea OR FVC OR “pulmonary edema” OR FEV OR fibrosis OR granuloma* OR hypoxemia OR pneumon* OR pulmon* OR spirometry OR “radiographic X-ray”)
Formaldehyde	Exclude use of formaldehyde as a fixative	“Chemical” NOT (“formaldehyde fixation” OR “formalin fixation” OR “formalin fixed” OR “formaldehyde fixed”)

10 * indicates wild-card search term; search will include all permutations of the word with the specified backbone
11 (e.g., *toxic* will include neurotoxic, toxicity, immunotoxicant, etc.)
12

13 **1C. Augmentation of a Database Search**

14
15 **Do Not Rely Solely On the Initial Computer-Based
16 Search of Databases to Identify Relevant Studies**
17
18

19 Some publications will be missed with even the best-designed algorithm-based search strategy.
20 Publications can be missed because they are not indexed correctly, or because the relevant data in
21 the paper are not mentioned in the abstract. Articles published before 1965 are likely to be missed
22 because of the coverage of the primary databases used in the search process. In addition, many
23 older papers (e.g., published before 1970) do not include an abstract and so are difficult to evaluate
24 in the initial screening process.

25 The team responsible for the IRIS assessment should identify key studies (review articles,
26 other comprehensive documents, and articles with primary (i.e., original) data) that will serve as
27 the basis for additional searches. It is useful to include reviews from different time periods,
28 because earlier review papers may have more descriptive information about earlier studies.

29 Additional search strategies that can be employed through a database (e.g., Web of Science)
30 include “forward searching”, and “backward searching” based on articles identified as key studies.
31 Forward searching identifies articles that cite the key study, and backward searching identifies
32 articles that the key study cites. Using the forward and backward search options through Web of

1 Science does not eliminate the need for additional review, particularly of earlier time periods, given
2 the coverage limitations of the database.

3 Regulatory resources and other web sites for information pertaining to a chemical of
4 interest should be checked for additional resources.

5 As an additional check on the completeness of your search, you can send the list of
6 identified relevant studies (i.e., studies that pass the screening step and are moved to Study Quality
7 Evaluation) to researchers who are currently or were previously active in the specified area of
8 research. They may know of other studies, including “file drawer” unpublished studies.

9 When using the “forward” or “backward” searching strategies through a database, the
10 resulting set of references will need to undergo the screening for relevance step described in the
11 next section. Other strategies will result in a more refined set of additional references that can
12 bypass that step. For example, when reviewing the references in a discussion section of a paper,
13 you may find six citations for similar studies, one of which is not already included in your search
14 results. That one reference would be added to your literature search, but would not need to be
15 further screened for relevance.

16

17 ***1D. Documenting the Search Strategy***

18

**Tell the Story of Your Strategy, Preferably in
Such a Way that Someone Else Can Reproduce It**

19

20

21

22 Accurate documentation of the search strategy is an essential component of the systematic
23 review process. You want to be able to provide the reader with the information needed to
24 understand what you did. Documentation of database searches should include, at a minimum, the
25 database(s) and date range covered by the search, search terms used and the fields (e.g., title,
26 abstract) to which they were applied, and dates the searches were performed (Table F-3). It is also
27 useful to keep a log of the sources and approaches you used to augment the initial database search
28 (Table F-4). In addition to these details, information pertaining to the context of your search
29 strategy (not what you did but why you did it), such as why you focused the search in particular
30 ways and other ways in which the search strategy evolved, should be included in the text describing
31 your literature search.

32

Table F-3. Example Worksheet Summarizing the Database Search Process (Note: this is a research aid; this is not expected to be included in the finalized assessment)

Database	Set #	Terms	Hits
PubMed Date range Search date	1A	CHEMICAL TERMS; ADDITIONAL TERMS	
Web of Science Date range Search date	1B		
ToxNet Date range Search date	1C		
Other Database Date range Search date	1D		
Merged Reference Set	1	(duplicates eliminated through electronic screen)	

Table F-4. Summary of Additional Search Strategies

System Used	Selected Key Reference(s)	Date	Additional References Identified
Web of Science, forward search	Review study: Yuko et al., 2000		N, citation(s)
	Review article: Smith et al., 2010		N, citation(s)
	Primary study: Kim et al., 2006		N, citation(s)
Manual search of cited references	Primary study: Kim et al., 2006		N, citation(s)
	Review article: Drew et al., 1966		N, citation(s)

1E. Updating the Literature Search

Establish a System to Regularly Update the Literature Search for IRIS Assessments

IRIS assessments can take over 2 years to complete. The length of this process necessitates the establishment of a system to regularly update the literature search. The team responsible for the IRIS assessment is responsible for this updating process. An “alert” system can be set up through the core literature databases for automatic notification of new citations based on a designated search string. The frequency of the updates depends on personal preferences and the relative amount of research activity for the chemical under review, ranging from weekly for very large and active research areas, to bi-monthly for chemicals with relatively little active research. A cut-off date can be established for various steps. For example, although the updating process would continue throughout the development of the assessment, only studies identified up to a certain point (e.g., 45 days) before the literature search results are posted for public review would be

1 included. An additional cut-off date would be used for subsequent steps, such as finalization of the
2 external peer review draft. Notably, after a certain step of the process (e.g., external peer review),
3 additional studies will only be added if it is expected that they will substantially change the
4 conclusions of the assessment; thus, a full literature search update will no longer be necessary.

Review of General Principles: Literature Search

- Use your initial search strategy to “cast a wide net”
- Include databases of published papers (PubMed, Web of Science, Toxline) and unpublished studies (TSCATS, TSCATS2; OPP documents for pesticides)
- Consult a reference librarian (through EPA resources or contractor resources) to develop search terms
- If using HERO, include “tags” for each database in initial project page set-up
- Augment database search with additional sources (e.g., lists of references from reviews and primary studies); solicited review of the identified studies by knowledgeable investigators may yield additional references (published and unpublished)
- Update literature search regularly (at least bi-monthly), with specified cut-off dates for study inclusion before key steps
- Keep a record of database(s), dates, search terms, results (n citations) and augmentation strategies
- It may be necessary to update your search criteria based on discoveries made at later stages of the systematic review process, such as evidence synthesis and integration. Be flexible to change, but be sure these criteria are documented and applied systematically

STEP 2: SCREENING FOR RELEVANCE

This Step Addresses the Following Question: Is This Study Relevant to the Question of Interest?

33 It is highly likely that many, or even most, of the studies identified using your search
34 strategy are not useful because they do not address the question of interest (i.e., the health effects
35 of a chemical, or for more focused searches, a particular type of health effect or specified mode of
36 action). This result should not be viewed as a deficiency of the search strategy process, but rather
37 is expected given your goal of “casting a wide net.” Your next step is to undergo a systematic review
38 of each of the citations to determine its relevance; neither the quality of the study, nor the results
39 are considered in this step. Depending on just how wide a net you ended up casting, this process
40 could be somewhat akin to finding a needle in a haystack. In an effort to efficiently identify the non-
41 relevant studies, this screening step is broken down into two sequential stages, title and abstract
42 screening followed by full text screening. In some situations, a three-stage process may be more
43 efficient, with an initial screen based on title, followed by screening based on abstract, followed by

1 full text screening. There is not a “right” or ‘wrong” choice; however, whichever you choose, be
 2 sure to document the process you use.

3 The articles identified as relevant are then organized into topic-specific bins to aid you in
 4 performing the next step of evaluating the quality of individual studies.

5

6 **2A. Review Process**

7 This step in the review process is based on review of the title and abstract, and in some
 8 cases, the full text of the article, and should be conducted by two reviewers. If a contractor is used
 9 for this step, one of the reviewers should be an EPA staff member. It is meant to be a limited review
 10 of each citation as a way to relatively quickly exclude the large portion of citations that are not
 11 relevant.

12 There are numerous reasons a study may not be relevant to the subject matter of interest.
 13 Some reasons are common to all chemicals, and some may be tailored based on the specifics of the
 14 chemical. In some cases it is not possible to determine if the paper contains relevant data based on
 15 the information contained in the title and abstract; these citations should be set aside for additional
 16 perusal. Examples of the decisions that may be made about a citation, and reasons for these
 17 decisions, are shown in Table F-5, and discussed below.

18

Table F-5. Examples of Decisions Made Regarding Relevancy of Citation to Research Question (e.g., health effects of Chemical X)

Decision	Reasons
1. Exclude from consideration	<ul style="list-style-type: none"> • Duplicate • Abstract only (full report not available) • Not relevant - define categories as appropriate, for example: <ul style="list-style-type: none"> - study that uses chemical in sample preparation or assay - study that uses chemical as a positive control - study of effects on ecosystems
2. Not a primary data source of health effects data, but keep as additional resource	<ul style="list-style-type: none"> • Review articles, meta-analyses, editorials, risk assessments (use as source of additional references, discussion of key issues) • Articles describing development of measurement methods or exposure levels • Absorption, distribution, metabolism, and excretion studies • other (to be specified)
3. Further review needed	<ul style="list-style-type: none"> • No abstract • Language other than English • Case reports • Not enough information in title and abstract to determine relevancy • other (to be specified)
4. Move to full text screening	<ul style="list-style-type: none"> • Seems to be relevant to question of health effects of Chemical X

19

1 **Decision Category 1. Exclude from consideration**
2
3

4 **Excluded Studies**

- 5 • This category should contain a set of studies that two reviewers agree
6 can be eliminated from further review; these studies are not included in
7 the assessment.
- 8 • The goal of this step is to eliminate studies that do not contain original
9 research that addresses the relevant question; neither the quality of the
10 methods nor the details of the results should be considered in assessing
11 the question of relevance at this stage.
12

13 This category will be relatively large. It will contain duplicates that were not caught
14 electronically after the merging of the citations from the different databases (i.e., differences
15 in punctuation or capitalization style may result in two citations being counted as two
16 separate entries rather than as identical entries; when encountered, these types of
17 duplicates need to be manually eliminated from the database). It also includes studies
18 available only in abstract form (i.e., presented as a poster or presentation at a conference,
19 but never published). As mentioned previously, this abstract-only group is treated as a
20 separate category. In general, abstracts do not present enough information to allow
21 evaluation of the study details. The assessment team can review this category and decide
22 the appropriate level of effort to be used to pursue these studies; it may be possible to
23 contact the study author(s), particularly if it is a relatively recent study, to obtain additional
24 information. The decision to seek additional information should be consistently applied
25 across all studies within the database that are similarly lacking information and should not
26 be based on criteria that rely on the direction or magnitude of the study results.

27 This category will also include the studies that are “not relevant” – that is, studies
28 that do not address the question of the health effects of the chemical of interest. Some of
29 these types of “not relevant” studies are shown in Table F-5. You will find many other
30 reasons that a study does not pass the “relevancy” test. The process of sorting through the
31 database is facilitated by development of a list of common types of “not relevant” studies
32 you are likely to encounter. The IRIS assessment team is responsible for developing this
33 list, drawing upon the experiences of the reference librarians at EPA or a contractor you
34 may be working with. You can do this by reviewing the results of other chemicals that have
35 undergone this type of review (i.e., what are the categories that have been used
36 previously?). Since each database may present unique issues, it may also be useful,
37 particularly for chemicals with large databases, to systematically review a sample of the
38 citations to develop a set of relevancy-exclusion categories – i.e., sort by year, take 5-10
39 citations per year, review titles/abstracts to get a sense of what this database includes, and
40 use this information to develop the categories or questions that can be used to more easily
41 sort through the entire database. You may initially have a large “miscellaneous reasons”
42 category; this category can be examined and organized further as part of the review
43 process.

1 There may be multiple reasons that a study can be considered “not relevant”; the
2 reviewers should agree that the reason is correctly applied to each of the studies included in
3 that category, but it is not necessary to count a study within every category that applies to
4 it. It may be possible to create a hierarchy of categories, with those that are likely to be
5 most easily determined (e.g., duplicates, reviews) placed first, to facilitate the review
6 process. The two reviewers need to assure they have the same interpretation of the
7 meaning of each category. For large databases especially, this may involve working through
8 selected batches of 50-100 citations as “training” exercises. New categories should be
9 documented with a written description of its definition, with enough detail that someone
10 else could read it and determine that it was correctly applied. If a hierarchy of categories is
11 used for the review process, this, too, should be documented.

12 One strategy for accomplishing this task is to have one member do the initial
13 screening and sorting of the database, with the second member responsible for checking the
14 accuracy of each of the resulting group (i.e., assuring that the reason for exclusion applies to
15 each study in this group). A final step is resolution of the differences or discrepancies that
16 are found. This approach allows for each study to be reviewed using two different
17 frameworks: one asking “Does this study belong in the “not relevant” category, and if so
18 why?” and the other asking “Is it true that [the specific lack of relevance category] applies to
19 this study?” As with all of the other steps in the systematic review process, be sure to
20 document the procedure(s) you use.

21 What you end up with in this category is a set of studies that two reviewers agree
22 does not need to be considered further. If there are cases where the reviewers do not reach
23 resolution and you are unsure of a study’s relevance, set it aside for further review.

24 The main options for conducting the literature screening step are 1) “tagging”
25 through HERO, 2) sorting sets of citations in EndNote (with eventual importation into
26 HERO), and 3) supervising the screening process conducted by contractors. As discussed in
27 the literature searching step, when working with contractors, you want to take an active
28 role in decision making and quality assurance.

29 The most common approach for “tagging” in HERO is through use of the “LitTagger”
30 function and EndNote. It is possible to directly “tag” citations in HERO, but that option does
31 not work well for more than a minimal number of citations at a time. The “tags” used to
32 represent the different exclusion categories should be specified when setting up the initial
33 project page; modifications can be made but will need to be requested through the HERO
34 librarians. When working with the downloaded HERO citations through EndNote, you will
35 need to save the file to your desktop or file system if you want to complete the tagging
36 process in multiple sessions. The HERO team is working on enhancements to the HERO
37 database that will, in the future, allow you to complete the tagging over multiple sessions
38 directly in HERO.

39 If you conducted your initial searches through the individual databases, rather than
40 through HERO, you can use the EndNote grouping function for the screening process. After
41 the database is uploaded into HERO, you can use the EndNote groupings to generate lists of
42 HERO IDs for each of the exclusion categories. This process can be somewhat cumbersome
43 for long lists, so you may need to ask for help from a HERO librarian. In brief, the list is
44 generated by changing style to “HERO ID”, selecting the group of references, right clicking

1 the mouse, and copying into the project page; a preface of ‘hero.’ then needs to be added to
2 each number.

3
4 **Decision Category 2. Not a primary data source of health effects data: keep as**
5 **additional resource**

6
7 **Additional Resources**

- 8 • This category includes reviews and types of studies that can serve as a
9 useful additional source of potentially relevant primary articles, and
10 studies that provide background information that could be useful in
11 evaluating the health effects literature.
12

13 Review articles may address the question of the health effects of a chemical, but they
14 are not considered relevant in that they are not a source of original data (i.e., a “primary”
15 article). These types of studies, including meta-analyses, should be eliminated from
16 consideration of primary data, but should be kept as an additional resource. For example,
17 earlier reviews may contain information about studies that were not obtained in your
18 search strategy because of limitations of the online databases and changes in indexing
19 terms over time. In addition, as noted previously, reference lists of review articles also
20 serve as a good source to augment your algorithm-based search strategy (“backward
21 searching”). Reviews can also provide background information on issues you will need to
22 consider as you evaluate the literature. Finally, some review articles do contain primary
23 data (i.e., updates of previously reported data by the review authors), so additional review
24 of the paper to specifically look for new primary data should also be conducted by one of
25 the members of the screening review team.

26 Depending on the database, there may be other sets of studies that do not contain
27 primary data pertaining to the toxicity of the chemical, but do contain background
28 information that may be useful. For example, you may find studies describing the
29 sensitivity or specificity of a particular type of effect measure, or studies of exposure levels
30 in the general population or in different types of occupational settings. These studies can
31 also be set aside as additional resources to draw upon in your evaluation of the primary
32 studies. Absorption, distribution, metabolism, and excretion (ADME) studies are other
33 examples of studies that you want to retain for use in the assessment.

34 The options for documenting this category are the same as those discussed for the
35 first group of exclusion categories, and will most likely involve working through HERO and
36 EndNote, or through a contractor.
37

1 **Decision Category 3. Possible further review**
2

3 **Possible Further Review**

- 4 • This category includes sets of studies that may be relatively easy to
5 define, but difficult to process.
6

7 Some of the groups of studies in this category may be the easiest sets of studies to
8 create. That is, it is relatively easy to select studies with no abstract and case reports
9 through standard database search and sorting capabilities (e.g., through EndNote). The
10 reference librarian(s) working with you on your search strategy can help with this sorting
11 process.

12 Once organized, however, another question that must be tackled is what is to be
13 done with them. For example, review of the set of case reports can give you an idea of the
14 types of effects that have been seen, but only a limited number of these citations would need
15 to be included if more extensive epidemiological studies examining the types of effects
16 described in the case reports are available.

17 The category of citations with no abstract can be a relatively large group, and can
18 include letters to the editor, commentaries, older studies (e.g., published before the 1970's),
19 and some "brief reports" or "brief communications" found in some journals. It is often very
20 difficult to determine the type of article, or the topic of the article, solely from the title. In
21 some cases, the title may be sufficient to allow you to move it to another category (e.g.,
22 review article). Neither decision at either end of the spectrum (i.e., exclude them all or
23 spend the time and resources to obtain, translate if necessary, and review them all) is likely
24 to be an optimal decision.

25 The category of non-English language studies can also be relatively large, and it can
26 be difficult to determine relevancy based on the limited information available. Often a
27 relatively easily-obtained translation (such as through Google Translate) of the title and
28 abstract will be enough to determine if an article belongs in one of the "not relevant"
29 categories. Another approach is to use a "forward search" for papers which cite the foreign
30 language studies; this can provide a better sense of the information that each contains.
31 Decisions to translate the full text of articles that appear to be primary data sources of
32 health effects data (and which are not also published as an English-language report) need to
33 consider characteristics of the database, and available resources.

34 As the IRIS Program gains more experience with this process, more definitive advice
35 may be developed as to how to proceed (i.e., whether attempts are made to obtain the
36 complete publication, and to translate it if necessary). At this time, however, it is up to the
37 assessment team to review each of these batches of citations to develop a decision making
38 process that works given the scope of what is found for a given chemical. As has been noted
39 previously, information that is available about the magnitude or direction of effects seen in
40 a given study (e.g., from the phrasing of the title) should not be used in the decision
41 regarding how to proceed. That is, you do not want to decide to translate a study because it
42 looks like it has "positive" or "negative" results; rather, your decision should be based on
43 your perception of the likelihood that the citation contains primary data pertaining to the
44 research question (i.e., the health effects of a specific chemical). The options for

1 documenting this category are the same as those discussed previously, and will most likely
2 involve working through HERO and EndNote, or through a contractor.
3

4 **Decision Category 4. Move to full text screening**

5
6 **Move to Full Text Screening**

- 7 • This category consists of the set of articles that appears to contain
8 primary data pertaining to chemical toxicity; there is enough
9 information available in the title and abstract to warrant further review.
10

11 The initial screening process should leave you with a much smaller set of studies
12 than you started with. This smaller set of studies that will be subject of additional review
13 through examination of the full text. During this process, it is likely that for some, the
14 additional perusal of the full article will result in the realization that the study does not, in
15 fact, belong in the group of “relevant” studies either because of one or more of the reasons
16 used to define Decision Category 1 (Excluded Studies) or because of some other reason. In
17 these situations, the citation should be “tagged” into the appropriate exclusion category.

18 Steps 1 and 2 of the literature search process, literature search and screening for
19 relevance, are summarized in Figure F-1.
20

21 **Review of General Principles: Screening for Relevance**

- 22 • Focus is on this question: Does the study provide primary data relevant to the
23 question of health effects of exposure to the chemical?
- 24 • Quality of the study is not considered in this stage of the review
- 25 • Based on title, abstract, and in some cases, full text
- 26 • Document screening process
- 27 • Two reviewers for screening process
- 28 • Some sets of articles will need to be put aside for additional decision-making
29 (i.e., should the full article be obtained?)
- 30 • If using HERO, include “tags” for each exclusion category in initial project page
31 set-up
32
33

Search Topic: Sources of Primary Data on Health Effects of Chemical X

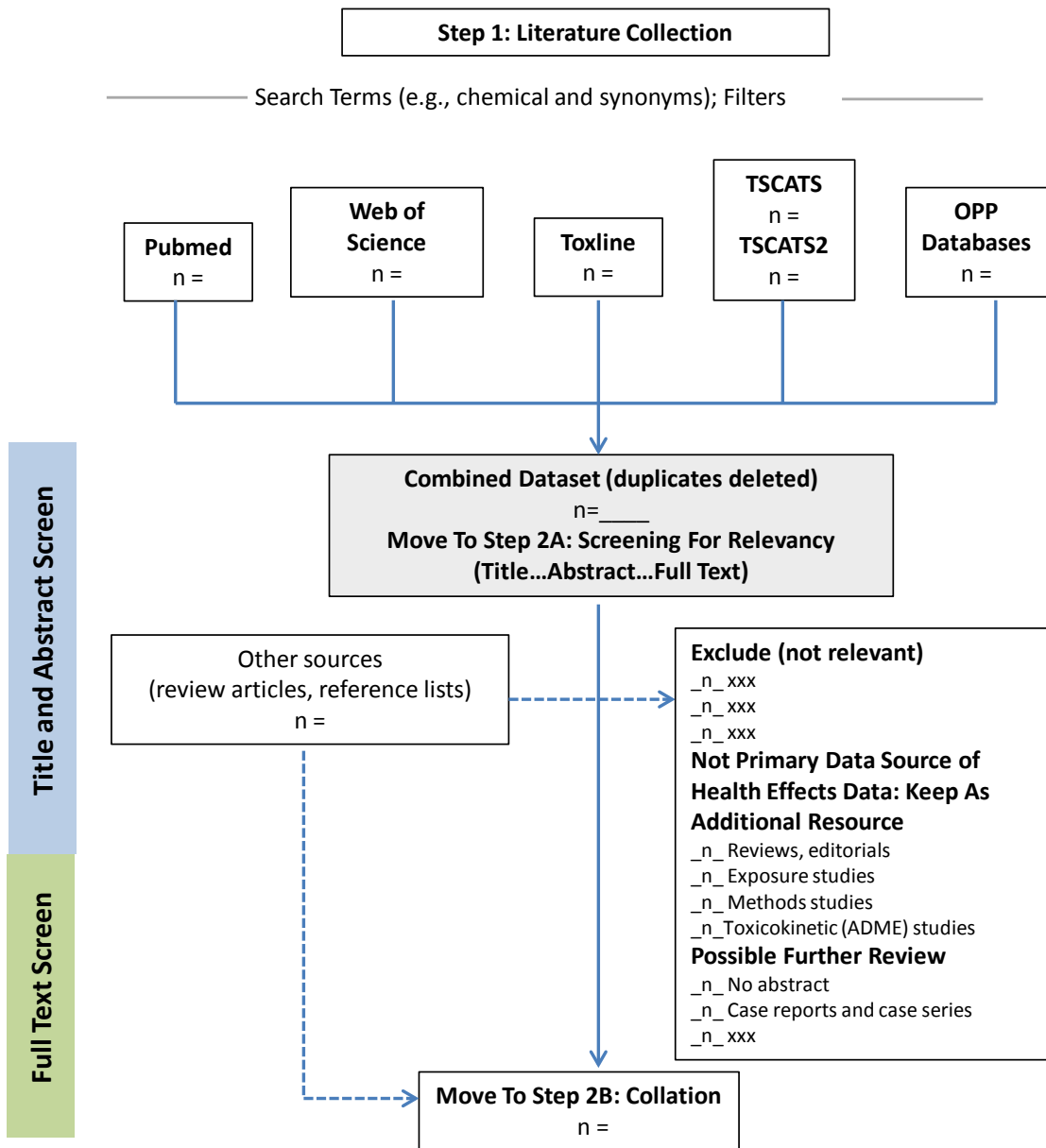


Figure F-1. Summary of Literature Collection and Screening for Relevance

1 **2B. Collation (Sorting)**
2

3 You should now be left with a (relatively) manageable number of citations that have a
4 (relatively) high likelihood of providing primary data pertaining to the question of the health effects
5 of a chemical. This collection of studies could include acute exposure animal experiments, two-year
6 bioassays, experimental chamber studies in humans, observational epidemiology studies, in vitro
7 studies, and many other types of designs. A basic organizational structure for the database is
8 needed to facilitate the evaluation of this collection of studies.

9 The optimal organization structure will depend on many factors including the number of
10 citations and breadth of topics and designs it includes. The assessment team should peruse the
11 database to get a sense of the specific question(s) addressed by the available studies and to make a
12 determination as to the optimal approach to sorting the studies. The actual process of sorting the
13 database can be done by a contractor or by an EPA staff member. The goal is to create groups of
14 studies that are of the same “type,” such that specific evaluation tools (described in the Study
15 Quality Evaluation section) can be applied. In general, the most likely divisions will be studies of
16 the chemical’s toxicity in humans, the chemical’s toxicity in animals, and mode or mechanism of
17 action (including in vitro studies). For large databases, however, additional categories or
18 subdivisions within these categories may be needed. Figure F-2 provides an example of a sorting
19 structure that may be useful for human and animal studies. For certain databases, it may be
20 necessary to provide a greater level of detail than that presented in Figures F-1 and F-2.

21
22 **Review of General Principles: Collation**

- End result is references organized into categories that make sense for next step (study quality evaluation)
- 23
24

Search Topic: Sources of Primary Data on Health Effects of Chemical X

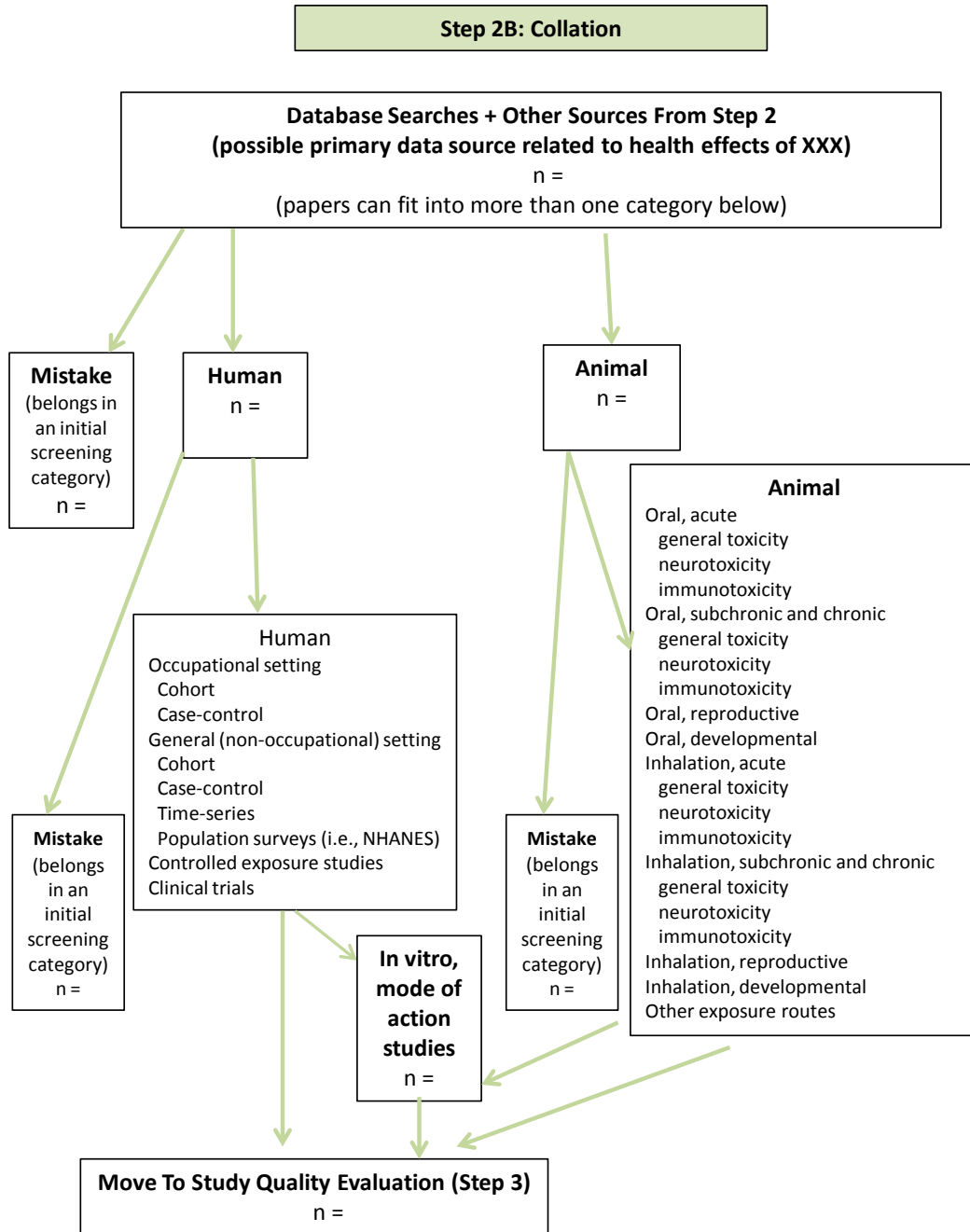


Figure F-2. Example of Collation, Human and Animal Studies

1
2

Systematic Review References

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17

Briss PA, Zaza S, Pappaioanou M, et. al. (2000) *Developing an Evidence-Based Guide to Community Preventive Services-Methods*. Am J Prev Med. 18(1S).

- [An example of a codified system for evaluating specific aspects of the design and execution of individual studies, in conjunction with the pattern of results seen across studies, for the purpose of evaluating the evidence pertaining to effectiveness of a type of intervention.]

Higgins JPT and Green S, eds. (2008) *Cochrane Handbook for Systematic Review of Interventions*. West Sussex, England: John Wiley & Sons, Ltd.

- [A guide to the content and methods of systematic reviews, as developed and applied in evidence-based medicine.]

Eden J, Levit L, Berg A, Morton S., eds. (2011) *Finding What Works in Health Care: Standards for Systematic Reviews*. Washington, DC: National Academies Press,

- [The Institute of Medicine’s report on standards for systematic reviews for comparative effectiveness evaluation for clinical practice.]

Evaluation and Display of Individual Studies

STUDY QUALITY EVALUATION

Study Quality Evaluation: Overview

- Be inclusive: it is better to include a study and evaluate effects of potential limitations than to exclude a study and eliminate any information the study could have provided
- Evaluate studies BEFORE developing evidence tables
- Series of focused questions; applied systematically to all primary data studies identified as relevant in the screening steps
- Evaluation is endpoint-specific; a given study evaluating several endpoints may have different strengths and limitations with respect to each endpoint

Study “quality,” as defined herein, is a broad term encompassing interpretations regarding a variety of methodological features (e.g., study design, exposure measurement details, study execution, data analysis and presentation). The purpose of this step in the systematic review process is not to eliminate studies, but rather to evaluate studies with respect to potential methodological considerations that could affect the interpretation of or confidence in the results. For larger databases, in particular, this evaluation can provide a transparent means to convey your assessment of a study’s methodological strengths and limitations, and thus your ability to rely on the results. The results of this systematic evaluation may also inform decisions about which studies to move forward for dose-response modeling for derivation of toxicity values.

The systematic evaluation described in this step should be conducted at an early stage of assessment development, i.e., after identifying the relevant sources of primary data but before developing evidence tables and characterizing hazards associated with exposure to a chemical. All studies identified as relevant from the literature screening process should be evaluated. Even if a deficiency in an aspect of the study is obvious, it can be useful to complete the evaluation of all of the component questions so that a full record of the evaluation can be maintained.

Examination of specific methodological features of each study can be accomplished by applying a series of focused questions. A useful starting point for generating these assessment and endpoint specific questions would be to consider the examples provided in Tables F-6 and F-7 for observational epidemiology and animal toxicology studies, respectively. Documentation of the important methodological features of a study may be an iterative process, requiring modification of an initial set of questions, as specific features of the chemical, endpoint(s), or study design(s) are discovered. It is essential that these focused questions be applied uniformly to all studies evaluated. This will allow for a comparison of the considered studies that is both systematic in design and independent of the study results. Ideally, two reviewers would independently identify the relevant methodological details, and then compare their results and interpretations and resolve any differences.

For studies that examine more than one endpoint or outcome, the evaluation process should be endpoint-specific, as the utility of a study may vary for the different endpoints.

1 The methods section of the paper will generally provide the majority of information needed
2 for this evaluation (except, of course, for considerations relating to the level of detail of the
3 reported results). In some cases, however, study details may be presented elsewhere in the
4 manuscript or report, such as the introduction or discussion sections. Identification of some study
5 details may require additional investigation, for example, by consulting other publications
6 describing the study or studies on the reliability of an assay, or by contacting the study authors. In
7 general, study quality evaluation should be independent of considerations regarding the direction
8 or magnitude of the study's results.

9 It is useful to check the citation in one of the primary databases (e.g., PubMed) to see if there
10 is any linked material, such as an erratum, supplementary or appendix material, letter to the editor
11 (and authors' reply) regarding the citation, or companion study. This kind of preliminary work can
12 prevent significant heartburn and headaches in subsequent steps.

13 It is useful to record the pertinent methodological features in an easy to read form (e.g., a
14 tabular format) so that these study details can be easily reviewed. Because observational
15 epidemiology and animal toxicology studies have fundamental differences, the documentation and
16 evaluation of these studies will differ.

17 There may be situations, most commonly when extensive literature databases exist for a
18 given chemical and effect, in which an individual study or sets of studies can be excluded from
19 further consideration. For example, acute animal toxicology studies may be excluded when
20 abundant subchronic and chronic exposure studies examining similar endpoints are available.

21 The following discussion of study quality evaluation is focused on evaluation of
22 observational epidemiology, animal toxicology, and human controlled exposure studies. This
23 approach could also be adapted for the evaluation of in vitro studies and other types of studies
24 relevant to mechanisms of action.

Study Quality Evaluation: Logistics

- Methods section of the study should provide most of the information you need; study quality evaluation should be independent of considerations regarding the direction or magnitude of the study's results
- Look for errata, supplemental files, and other material linked to the primary data citation for additional information about the study
- Published correspondence (e.g., letters to the editor, editorials) may provide additional background information on important methodological features.
- Ideally, use two *independent* reviewers, with procedures for disagreements to be reviewed and resolved

Evaluation of Observational Epidemiology Studies

37 The process of study evaluation is akin to detective work. You need to investigate specific
38 study features that directly affect the interpretation of the experimental results, including:

- exposure measures (reliability, validity, probability and level of exposure in different situations or settings)

- 1 • outcome measures (reliability, validity, prevalence in different populations, disease course,
2 relation between survival and access to health care or other socioeconomic factors)
- 3 • confounders (strong risk factors for the outcome that are also known to be strongly
4 associated with the exposure within the study)

5 These investigations may require “mini-reviews” and consultation with experts in different fields.
6 Without this background understanding, you may not be able to accurately evaluate the studies.

7 Exposure assessment is especially important in the environmental or occupational arena.
8 The ability to correctly classify “exposed” and “unexposed”, estimate quantitative measures of
9 exposure, and the range of exposure encompassed in the study is a key difference between
10 observational epidemiology and randomized clinical trials in which “exposure” (e.g., “intention to
11 treat” or type of treatment) may be less subject to measurement error and the exposure contrast is
12 less variable between studies.

13 As noted above, an inclusive approach is generally recommended: that is, it is better to
14 include a study in this systematic evaluation and examine the impact of potential limitations, rather
15 than exclude a study and thus lose any information it could have provided. For epidemiology
16 studies, to the extent possible, you want to assess not just the “risk of bias,” but also the likelihood,
17 direction, and magnitude of bias.

18 The study characteristics that inform the evaluation of observational epidemiology studies
19 are summarized in Table F-6. The first feature, the type of study design, provides a framework for
20 the subsequent evaluation; that is, the specific questions and issues will vary depending on the type
21 of study. The other features encompass aspects of the study populations, exposure measures,
22 outcome (effect) measures, and the analysis and presentation of results. Although in general your
23 evaluation is based on the information provided about study methods, review of some of the results
24 is needed, for example within the context of the evaluation of confounding, since confounding
25 depends on the strength of various relationships (i.e., between the exposure and the potential
26 confounder and between the potential confounder and the outcome).

27 A structured form may be useful for recording the key features needed to evaluate a study.
28 An example form is shown in Figure F-3; details of such a form will need to be modified based on
29 the specifics of the chemical, exposure scenarios, and effect measures under study.

31 **Study Quality Evaluation: Observational Epidemiology Studies**

- 32 • As noted in the overview, the evaluation process is inclusive in nature, is
33 conducted BEFORE developing evidence tables, uses a series of systematically
34 applied, focused questions, and is end-point specific
- 35 • Do your detective work ahead of time: investigate exposure measures, effect
36 measures, and confounders for the chemical-effect under review
- 37 • To the extent possible, assess likelihood, direction, and magnitude of bias

Table F-6. General Considerations for Evaluation of Features of Epidemiology Studies

Feature	Example Questions or Details	Useful Information
Study design	Major types, based on approach to sample selection: cohort, case-control, nested case-control, population-based survey (e.g., NHANES), times series, case-crossover	Study methods
Study population; target population; setting	Where and when was the study conducted? What is the source(s) of exposure (environmental media, consumer products, occupational, an industrial accident, or other)? What was the recruitment process? How was eligibility determined? Does the study provide information on potential vulnerable or susceptible groups? Address: Potential generalizability of study results, potential for selection bias, potential to address effect modification	Geographic area, site (occupational, etc.), time period. Age and sex distribution, other details as needed (may include race/ethnicity, socioeconomic status); recruitment process; exclusion and inclusion criteria
Participation rate; follow-up	Did rates vary by exposure (or disease) status? Were there differences between individuals who did and did not participate, or who were or were not lost to follow-up? Is it known (or possible) that participation (or loss) is related both to exposure and disease status? Is there evidence of “healthy worker” or “healthy worker survivor” effect? Are differences likely to impact the observed associations (and if so, how)? Address: Potential for selection bias	Total eligible; participation at each stage and for final analysis group; loss to follow-up, denominators used to make these calculations; length of follow-up
Comparability (exposed and non-exposed; cases and controls)	How were potential differences between groups addressed in the study design (e.g. randomization, restriction, matching) and/or analysis (e.g. stratification, multivariate methods)? How were variables associated with exposure and with outcome, or which alter the association between exposure and outcome, addressed in the study? Address: potential for confounding and effect modification	“Table 1” type participant characteristic data, by group; approach to consideration of potential confounding (if applicable); strength of associations between exposure and potential confounders and between potential confounders and outcome
Exposure measures (procedure, range)	Are exposure estimates qualitative, semi-quantitative or quantitative? How well does the exposure protocol correctly classify or rank participants with respect to exposure? What is the likelihood of systematic (differential) error? What is the likelihood of random (non-differential) error? Does the protocol adequately characterize exposure during the relevant time window? What exposure range is spanned in this study? Address: potential for exposure misclassification (either non-differential or differential).	Describe, i.e., type of biomarker(s), occupational history, lifetime consumption, evidence from validation studies, variability within and between exposure groups
Outcome measures	What is source of outcome (effect) measure? How well do the outcome(s) measures correctly classify participants with respect to the outcome? What is the likelihood of systematic (differential) error? What is the likelihood of random (non-differential) error? Address: potential for outcome misclassification (either non-differential or differential).	Describe (i.e., source, how measured/classified, incident versus prevalent disease), evidence from validation studies
Data Presentation and Statistical Analysis	Is the analysis appropriate for the data and the study question? Are aspects of the data (i.e., non-normal distributions, correlation structure) adequately accounted for? Is the rationale for inclusion of variables in a model clear and logical? Are results presented with adequate detail? Is the study population of adequate size and composition to detect a true association (of a relevant effect size) between exposure and outcome? Were stratified analyses (effect modified) motivated by a specific hypothesis? Address: ability to interpret and level of confidence in results	How groups are compared (may include t-tests, ANOVA, regression models, etc.); what results are presented in text, tables, and figures; n exposed cases (case-control studies) or N cases among exposed (cohort studies).

1
2

**Figure F-3. Example Worksheet for Recording Methodological Details of Observational Epidemiology Studies
(Note: this is a research aid; this is not expected to be included in the finalized assessment)**

Reference (primary)		
Additional reference(s)		
Study Design ___ cohort ___ case-control ___ nested case-control ___ population-based survey (e.g., NHANES) ___ times series ___ case-crossover ___ controlled exposure ___ other (describe)		
Setting	Describe, i.e., geographic area, worksite, clinic; time period...	
Study Population Age Sex Other details as relevant (e.g., socioeconomic status) Duration of exposure	Excluded if....	Descriptive Statistics (e.g., median, range, etc – what is reported will vary among studies)
Participant Recruitment	Describe process	
	Evidence that knowledge of exposure and diseases status affected participation?	
Participation Rates / Follow-up (separate data for cases and controls, exposed and non-exposed if provided)	Total eligible: Participated, any part (describe): Participated, all parts: Loss to follow-up: Length of follow-up:	Evidence of differential participation?
Comparability of Groups	Comparability between exposed and non-exposed; cases and controls (“Table 1” type sociodemographic data)	
Exposure measurement protocol	Describe, i.e., type of biomarker(s), occupational history, lifetime consumption...	
Biomarker(s) details	Sample collection (time of day, fasting?)	Assay (e.g., coefficient of variation, limit of detection, proportion < limit of detection), blinded to outcome status?)
Outcome measurement protocol	Describe (i.e., source, how measured/classified, incident versus prevalent disease, etc)	
Medical records		
Likely confounders?	Variables strongly associated both with exposure and with effect? What is strength of associations in this study? How addressed?	
Analysis and presentation of results	Approach used; assumptions made regarding distributions or shapes?	
	Standard error or confidence intervals presented for effect estimates (or could be computed)?	
Statistical power considerations <i>may differ for different effects</i>	n exposed cases (case-control studies) or n cases among exposed (cohort studies)	
Reviewer Comments	<i>May include summary evaluation of likelihood, direction, and magnitude of bias. May include usefulness and feasibility of re-analysis.</i>	

Evaluation of Animal Toxicology Studies

In contrast to observational epidemiology studies, animal toxicology studies seek, by their very nature, to control exposure and environmental conditions. Considerations relevant to the evaluation of toxicology studies include exposure and endpoint methodology, as well as control of potentially confounding variables.

Table F-7 provides a list of questions relating to study features that should be considered when evaluating in vivo animal toxicology studies, namely: exposure, test animals, study design, toxicity endpoints, data presentation and statistics, as well as the reporting of this information. These questions are based on previous approaches for evaluating toxicological data (e.g., Klimisch et al., 1997; U.S. EPA 2002, 1994). These study features reflect aspects of an experiment that have been placed into modular components to assist in the analysis and transparent documentation of decisions; however, there is some overlap among the study features and it may be useful to reorganize some of the study features or clarifying questions for a given chemical or research question. For each study feature, the table provides a primary question that an assessor should try to answer using expert judgment. The example clarifying questions are included to provide direction and suggest ways to evaluate and document the study evidence that underlies these decisions. By no means are these questions comprehensive and, for most assessments and endpoints, some of these questions will not be applicable. For example, determining whether maternal toxicity was considered by study authors will only apply to evaluations of developmental and reproductive toxicology studies.

The purpose of the questions in Table F-7 is not to exclude studies from consideration. Rather, these questions are intended to help identify and characterize features of a given study that, together, can provide a picture of how well that study informs the specific endpoint in question. These evaluations should not preclude toxicologists from looking for patterns across studies on a given endpoint, even if all of the identified studies do a relatively “poor” job at analyzing the endpoint in question.

Additionally, not all clarifying questions or considerations are of equal importance. Although the relative importance of specific criteria may vary by endpoint, chemical, or database, in general the criteria in bolded text represent some of the more important questions to examine. Evaluations of exposure quality, study design, and toxicity endpoints will generally require the greatest effort. Exposure quality refers to the characterization of the animals’ interaction with the test article, which should be specific to the chemical of interest and tightly regulated by the study director. Exposure quality is a particular concern for inhalation toxicology studies because of the inherent complexity in generating and characterizing test atmospheres. Study design refers specifically to the setup of the experiment. It includes consideration of items such as the length of exposure, the distribution of test animals into dosing groups, and the timing of endpoint(s) evaluation. Endpoint evaluation refers to the specific methods used to assess the hazard in question, including whether the protocols used to evaluate the endpoints were appropriate and complete, as well as whether said protocols are subject to modification by factors present in the study other than the chemical of interest. It is important to reiterate that all of these features should be evaluated without consideration of the magnitude or direction of the reported study results. Finally, any decisions made during the evaluation of a given study should be applied consistently throughout the database of studies on that particular endpoint.

1
2
3
4
5
6

Study Quality Evaluation: Animal Toxicology Studies

- As noted in the overview, the evaluation process is inclusive in nature, is conducted BEFORE developing evidence tables, uses a series of systematically applied, focused questions, and is end-point specific (a study may be very useful for one type of endpoint, but not for another)

Table F-7. General Considerations for Evaluation of Features of Animal Toxicology Studies

Feature	Primary Question	Example Clarifying Questions/Considerations
Exposure Quality	Are the exposures well designed and tightly controlled?	<i>General Attributes</i> <ul style="list-style-type: none"> ▪ How well was the test article identified and characterized? Are co-exposures expected as a result of test article composition? ▪ Was there a vehicle control group?
		<i>Inhalation Studies</i> <ul style="list-style-type: none"> ▪ For generation and measurement of the test article, how accurate and appropriate were the methods employed? Were analytical concentrations in the test animals' breathing zone reported (i.e., not just target or nominal concentrations)? For aerosol studies, were the mass median aerodynamic diameter (MMAD) and geometric standard deviation (GSD) reported? ▪ Was a dynamic chamber used? Static chambers are not recommended.
		<i>Oral Studies</i> <ul style="list-style-type: none"> ▪ Diet/Water: Could accurate doses be determined (e.g., was consumption measured)? Are there any expected or reported issues related to stability, homogeneity, or palatability of the test substance? ▪ Gavage: To what extent would toxicokinetic differences due to bolus dosing be expected to influence the results?
Test Animals	Are the test animals appropriate for evaluating the specified effect(s)?	<ul style="list-style-type: none"> ▪ How well are the control and exposed test animals matched in aspects other than exposure? Was information available to evaluate potential effects such as systemic or maternal toxicity that could confound interpretation of the endpoint of interest? Were there any notable issues regarding animal housing or food and water consumption? ▪ Based on what is known about the endpoint in question, how well do the species, strain, sex, age, and/ or number of test animals examined inform this evaluation?
Study Design	Is the study design appropriate for the test article and the evaluated effect(s)?	<ul style="list-style-type: none"> ▪ How well do the timing, frequency, and duration of exposures inform the effect(s) measured? For example, are critical windows of development encompassed by the exposures when assessing developmental toxicity? Were multiple exposure groups tested? ▪ If the results are expected to be subject to confounding by factors introduced as a result of selection bias, were efforts made to protect against this (e.g., control for potential litter bias in developmental studies; randomization of treatment groups)? ▪ How well do the timing and/or frequency of the endpoint evaluation(s) inform the measured effect(s)? For example, is the latency between exposure and testing expected to influence the level of confidence in the results? ▪ Was the study conducted under Good Laboratory Practices (GLP)? ▪ How well does the study conform to established guidelines (e.g., EPA, OECD)? Was it designed to specifically test the endpoint(s) in question? ▪ Did the study include other experimental conditions or procedures (e.g., surgery) that may influence the results of the toxicity endpoint(s) in question? If so, were appropriate control groups (e.g., sham) included in the study design?
Endpoint Evaluation	Are the protocols used for evaluating the endpoint(s) reliable and specific?	<ul style="list-style-type: none"> ▪ How well do the procedures used to evaluate the endpoint(s) in question conform to established protocols? If novel or uncommon, are the approaches biologically sound? ▪ What is the level of specificity of the protocols used? Did they include control experiments to discern effect-specific contributions (e.g., learning and memory) from nonspecific contributions (e.g., from motor activity) to the output measure (e.g., escape latency)? ▪ How sensitive are the protocols for a given endpoint? ▪ As appropriate, were steps taken to minimize potential experimenter bias (e.g., blinding) and sampling bias (e.g., evaluation of multiple tissue sections/ organ)?
Data Presentation	Do the results provided allow	<ul style="list-style-type: none"> ▪ Are the statistical methods and comparisons appropriate and transparent? If not, is sufficient information available for the IRIS Program to conduct its own analyses?

Table F-7. General Considerations for Evaluation of Features of Animal Toxicology Studies

Feature	Primary Question	Example Clarifying Questions/Considerations
and Analysis	one to accurately identify the direction and magnitude of the observed effect?	<ul style="list-style-type: none"> ▪ Are there any notable issues regarding presentation of the results? For example, if data were pooled (e.g., pooled exposure groups; pooled sexes) and this is expected to influence interpretation of the results for a given endpoint, are the reasons justified? ▪ Did the study report an unexpectedly high/low level of within-study variability and/or variation from historical measures that was not addressed?
Reporting	Are the methods and results well documented?	<ul style="list-style-type: none"> ▪ Are all aspects of the study described in sufficient detail such that it can be evaluated across the five study features presented above? Are any critical descriptions missing? ▪ Are group sizes and results reported quantitatively for each exposure group, time-point, and endpoint indicated as examined?

Criteria in **bolded text** represent the more important considerations.

1 Additional information on study protocols (e.g., guidelines developed by EPA and the
 2 Organization for Economic Co-operation and Development [OECD]) that may prove helpful in
 3 evaluating study features can be found in the annotated reference list. Consult EPA and OECD
 4 guidelines for recommendations on the design and interpretation of toxicology experiments.
 5 Remember, however, that these are not intended to be comprehensive protocols designed to
 6 provide in-depth analyses of all endpoints of toxicity; thus, they are not to be used as the “end-all,
 7 be-all” references for evaluations regarding study quality across every study or endpoint.
 8

Study Quality Evaluation: Animal Toxicology Studies

- Because all aspects of a toxicology study should be controlled, it is expected that the exposure causes the outcome. Anything that makes you question this is likely a study limitation.
- Decisions made during the evaluation should be applied in a consistent manner throughout a given database of studies

Evaluation of Human Controlled-Exposure Studies

17 Human controlled-exposure studies combine aspects of observational epidemiology studies
 18 and animal toxicology studies. Examples of human controlled-exposure studies include
 19 randomized controlled trials, randomized intervention studies, and chamber studies. The main
 20 distinguishing feature of controlled-exposure studies relative to observational epidemiology
 21 studies is that the exposure is determined by the investigator (similar to an animal toxicology
 22 study). Therefore, many of the considerations relevant to evaluating animal toxicology studies in
 23 Table F-7, and in particular considerations related to exposure, apply to the evaluation of human
 24 controlled-exposure studies. Many of the same study features and considerations outlined for
 25 observational studies, in particular those related to study population, are also relevant for
 26 controlled exposure studies (see Table F-6). It is also important to consider the informed consent
 27 and other human subjects research ethics procedures undertaken in these studies, relative to the
 28 ethical standards prevailing at the time the research was conducted.
 29

1 **DOCUMENTATION OF STUDY QUALITY EVALUATIONS**

2 The method for documenting information on study features that inform study quality may
3 vary depending on the size and characteristics of the epidemiology or toxicology database. For
4 example, if only a small number of epidemiology studies are available, it may be sufficient to
5 summarize methodologic details in a single table. For chemicals with a small number of animal
6 toxicology studies of generally uniform study design and quality, it may be sufficient to describe the
7 information relevant to evaluation of study quality in the text. For data-rich chemicals with a large
8 number of epidemiology or toxicology studies, however, more detailed documentation in tables is
9 recommended to allow the user to see at once the number and type of studies available, and the
10 level of information available from each. Database tools have been developed for organization and
11 management of this type of information (e.g., through LitCiter Lite or DistillerSR software), but
12 additional testing and refinement is needed to establish their usefulness for IRIS assessments.

13 Options for displaying relevant study quality information from epidemiology and animal
14 toxicology studies in tables are described in the sections that follow. These tables could be included
15 in an appendix of an IRIS assessment. These tables serve to 1) document all the studies that were
16 considered; 2) provide the means to identify and track how informative a given study was
17 throughout the assessment process; and 3) document why some studies were not further
18 considered in the assessment.

19 Once the study information has been recorded and evaluated, it may be useful to sort
20 studies into “tiers” according to the level of information they provided. The considerations and
21 judgments used to “tier” studies should be clearly and transparently documented.

22 ***Documentation of Observational Epidemiology Study Evaluations***

23 Table F-8 is an example of a summary display of relevant information for observational
24 epidemiology studies. The shading of specific cells represents those features for which a specific
25 limitation was noted.

26 In some situations, the collection of studies may be divided based on the likelihood and the
27 types of limitations or biases identified in the evaluation of study quality. A study in the top quality
28 tier would typically use an appropriate study design, have high-quality measures of exposure and
29 outcome, and use adequate methods to analyze and present results. These studies would be given
30 the greatest consideration within the context of hazard identification. Studies of lower quality are
31 limited in one or more other ways. Because the type(s) of limitation(s) noted in a study can
32 influence the direction of bias, it may be important to further classify this group based on the
33 type(s) of limitations identified. For example, one group may consist of studies where the
34 limitation(s) are likely to result in an attenuated effect measure, such as studies for which the major
35 limitation is a substantial amount of non-differential exposure misclassification. Another group
36 may include studies with different types of limitations for which it is difficult to determine the likely
37 direction (or likelihood) of bias. Another group may include studies where the limitation(s) were
38 considered to be likely to result in observation of a spurious association, such as a study that did
39 not control for a known risk factor for the disease that was also strongly related to the exposure in
40 the study.
41

Table F-8. Evaluation of Observational Epidemiology Studies of Chemical X.

Reference, Setting and Design	Participants, Selection, Follow-up	Comparability	Exposure Measure and Range	Outcome Measure	Consideration of Likely Confounding	Analysis and Presentation of Results (Estimate and Variability)	Sample Size; Power	Evaluation of Major Limitations
Lee et al., 1995 US (New York) chemical X production plant (cohort)	All men, age at baseline not reported; duration ≥ 12 months (mean 2.2 years), worked at plant 1960 – 1972 - plant operations began in 1945. Follow-up through 1990, 2% loss to follow-up, mean follow-up time 32 years	External (state mortality rates) referent; age and time-period matched (5 year groupings). Healthy worker effect seen for CVD (SMR 0.7) and all cancers (SMR 0.9). Internal referent: “no” exposure group	Exposure based on job records and personal/air monitoring; cumulative exposure calculated based on summations across all jobs (duration times average exposure)	Mortality (death certificates, ICD-8 and 9, underlying and contributing causes of death)	External comparison: use of age and time-period matched mortality rates.	SMR and 95% CI	Brain cancer: 4 obs cases	Low statistical power; not an inception cohort (had to “survive” to 1960 to be included)
Johnson et al., 1996 US (24 states) (case-control)	All deaths 1984-1986. Controls (died of causes other than cancer; frequency matched by age, sex, state and race)	Matching procedures for cases and controls	Death certificate occupation data; job exposure matrix developed to assess 11 chemical exposures	Mortality (death certificates, ICD-9), underlying cause of death)	Sex-specific odds ratios adjusted for marital status, race, socioeconomic status (3-levels), age at death	OR and 95% CI	10,540 cases, 42,160 controls	Non-differential exposure misclassification likely, particularly for women (lower quality occupation data for women)

1 ***Documentation of Animal Toxicology Study Evaluations***

2 Study quality evaluation requires an analysis and documentation of the six categories of
3 study features described above. An example tabular documentation of study quality features for
4 animal toxicology studies is provided in Table F-9. Because the delivery of exposures in inhalation
5 toxicology studies is complex, it may be advisable to develop a separate table (as shown in Table F-
6 9a) that documents exposure quality in greater detail; the overall characterization of the exposure
7 quality can be characterized in terms such as “robust” and “marginal.” The quality of the exposure
8 characterization is then incorporated into the broader evaluation of the other 5 previously
9 described study features for each study (Table F-9b). It should be noted that this example was
10 derived for an evaluation of a large and complex dataset; more simplified documentation is likely to
11 be adequate for other types of datasets. Additional separate quality tables could be developed for
12 other exposure routes or for other specific study features requiring more in-depth analyses (e.g.,
13 endpoint evaluation of neurotoxicity and respiratory pathology).

14 As previously described for epidemiology studies, a “tiering” system may be appropriate for
15 categorizing animal toxicology studies according to aspects of study design, methods, and
16 execution.

17
18

19 **Review of General Principles: Study Quality Evaluation and Documentation**

- 20 • To the extent possible, evaluation of a study is independent of consideration of
21 the direction or magnitude of the study’s results
- 22 • The goal is not necessarily to eliminate studies, but rather to understand
23 potential limitations that would affect the interpretation of the results
- 24 • Record pertinent study details: what do you need to know about how the study
25 was designed and conducted?
- 26 • “Tiering” can be useful to allow an easier flow of discussion during evidence
27 synthesis and can transparently inform weight-of-evidence considerations
- 28 • Document judgments made regarding basis for “tiering” of studies

1 **Table F-9. Example of a tabular documentation of study evaluation for a large dataset. This example includes important issues regarding inhalation exposure quality,**
 2 **and is broken into: (a) an evaluation of inhalation exposure quality; and (b) incorporation of this exposure quality analysis into a larger evaluation with the other 5**
 3 **study features. The data are generalized and the endpoint is not specified. The results of this evaluation could be used to document an expert judgment that “Smith**
 4 **et al., 1984” is likely to be a more informative study (and “Gray et al., 2012,” less) evaluating the endpoint in question.**

5
6 (a) Evaluation of **inhalation exposure** quality

Reference (Species)	Test Article Characterization	Generation Method	Analytical Method	Analytical Concentrations	MMAD (GSD)	Chamber Type	Vehicle Control
Robust Exposure Characterization			Meet a robust standard for exposure quality				
Smith et al. (1984) (Monkey)	Test article (99%) solution in water	Bubble generator	Infrared spectrophotometry	Reported	Not applicable	Dynamic whole-body	Not needed
Marginal Exposure Characterization			Studies that meet a marginal standard for exposure quality. Key exposure data are missing				
Jones et al. (1986) (Mouse)	Solid test article (98.5%) <i>Co-exposure likely</i>	Thermal depolymerization	Chromotropic acid	Reported	1.3 (1.7)	Dynamic nose-only	No
Poor Exposure Characterization			Studies that may be inadequate for exposure response but which may support other studies in informing hazard				
Gray et al. (2012) (Rat)	Not reported	Not reported	Not reported	Not reported	Not applicable	Static	No

Study deficiencies noted in **bolded text**.

7 (b) Evaluation of **all features** of animal toxicology study quality

Reference (Species)	Exposure Quality ^a	Test Subjects	Study Design	Toxicity Endpoints	Data and Statistics	Reporting
Smith et al. (1984) (Monkey)	++	++ Note: N=20	++ Note: 102 wk study	++	Not applicable	++
Jones et al. (1986) (Mouse)	+ co-exposure likely	+ N=5; variable ages at onset of exposure across groups	++ Note: 13 wk study	Potential sampling bias; No observer blinding indicated; protocols incompletely reported	+ data represents pooled sexes	+ Surgical procedures not reported
Gray et al. (2012) (Rat)	Test article and exposure methods not specified	Bacterial infection noted in animal colony; N= 3 litters; males only; overt maternal toxicity	No randomization across litters into treatment groups; testing during exposure expected to confound results; acute exposure	++	Not applicable	Results data not reported

^a Summary results from inhalation exposure quality analysis in “Table F-9a.” Criteria for the 6 categories developed based on the chemical and hazard type in question. In this example: gray box = examination of relevant study details identified potential limitations that could influence interpretations of the study's results; '+' = criteria not completely met or potential issues identified, but unlikely to directly affect study interpretation; ++ = criteria determined to be completely met. Text accompanying summary table would explain key study details informing these determinations.

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
43
44

Study Quality References

Klimisch HJ, Andreae M, Tillmann U. (1997) A systematic approach for evaluating the quality of experimental toxicological and ecotoxicological data. *Regul Toxicol Pharmacol* 25(1):1-5.

- [Presents an approach to systematically evaluating the quality of animal toxicology data and their use in hazard and risk assessment.]

NTP (National Toxicology Program). (2012) Protocol: evaluation of cancer studies in experimental animals. U.S. Department of Health and Human Services. Available online at www.ntp.niehs.nih.gov/NTP/roc/thirteenth/Protocols/PCP_animalcancer_508.pdf

- [Presents the protocol for cancer assessment of animal studies for the NTP's Report on Carcinogens Monograph on pentachlorophenol. Appendix C is a particularly useful section: Assessment of the quality of the individual animal cancer studies. Various study performance elements are described as they pertain to evaluating study quality.]

OECD guidelines and guidance documents are the standard for toxicology study quality. Available online at http://www.oecd-ilibrary.org/environment/oecd-guidelines-for-the-testing-of-chemicals-section-4-health-effects_20745788

U.S. EPA (Environmental Protection Agency). (1993) Reference dose (RfD): description and use in health risk assessments. Available online at <http://www.epa.gov/iris/rfd.htm>

- [Describes the EPA's principal approach to and rationale for assessing risk for health effects other than cancer and gene mutations from chronic chemical exposure. Section 1.3.1.1.6 (Quality of the study) provides an overview of the types of factors that are generally considered while making determinations pertaining to study quality.]

U.S. EPA (Environmental Protection Agency). (1994) Methods for derivation of inhalation reference concentrations (RfCs) and application of inhalation dosimetry. Environmental Criteria and Assessment Office, Research Triangle Park, NC; EPA/600/8-90/066F. Available online at <http://cfpub.epa.gov/ncea/cfm/recordisplay.cfm?deid=71993>.

- [Discusses criteria and information to be considered in selecting key studies for RfC derivation. Appendix F is a particularly useful section: Criteria for Assessing the Quality of Individual Animal Toxicity Studies.]

U.S. EPA (Environmental Protection Agency). (1998) Guidelines for neurotoxicity risk assessment. Risk Assessment Forum, Washington, DC; EPA/630/R-95/001F. Available online at <http://www.epa.gov/raf/publications/pdfs/NEUROTOX.PDF>

- [Summarizes the procedures that EPA uses in evaluating the potential for agents to cause neurotoxicity.]

U.S. EPA (Environmental Protection Agency). (1996) Guidelines for reproductive toxicity risk assessment. Risk Assessment Forum, Washington, DC; EPA/630/R-96/009. Available online at <http://www.epa.gov/raf/publications/pdfs/REPRO51.PDF>

- [Summarizes the procedures that EPA uses in evaluating the potential for agents to cause reproductive toxicity.]

- 1 U.S. EPA (Environmental Protection Agency). (1991) Guidelines for Developmental Toxicity Risk
2 Assessment. Risk Assessment Forum, Washington, DC; EPA/600/FR-91/001. Available online at
3 <http://www.epa.gov/raf/publications/pdfs/DEVTOX.PDF>
- 4 • [Summarizes the procedures that EPA uses in evaluating the potential for agents to cause
5 developmental toxicity.]
6
- 7 U.S. EPA (Environmental Protection Agency). (2005a) Guidelines for Carcinogen Risk Assessment.
8 Risk Assessment Forum, Washington, DC; EPA/630/P-03/001F. Available online at
9 http://www.epa.gov/raf/publications/pdfs/CANCER_GUIDELINES_FINAL_3-25-05.PDF
- 10 • [Summarizes the procedures that EPA uses in evaluating the potential for agents to cause
11 cancer.]
12
- 13 U.S. FDA (Food and Drug Administration). (1982) Toxicological Principles for the Safety Assessment
14 of Food Ingredients (also known as Redbook 2000). Bureau of Foods (now Center for Food Safety
15 and Applied Nutrition), Washington, DC. Available online at
16 [www.fda.gov/food/guidancecomplianceregulatoryinformation/guidancedocuments/foodingredien
18 tsandpackaging/redbook/default.htm](http://www.fda.gov/food/guidancecomplianceregulatoryinformation/guidancedocuments/foodingredien
17 tsandpackaging/redbook/default.htm)
- 18 • [The U.S. Food and Drug Administration published this as guidance to industry and other
19 stakeholders regarding toxicological information submitted to its Center for Food Safety and
20 Applied Nutrition. The Redbook is an alternative resource wherein EPA scientists may find
21 recommendations that are useful in evaluating animal toxicology studies.]
22
- 23 WHO (World Health Organization). (2012) Guidance for Immunotoxicity Risk Assessment for
24 Chemicals. Harmonization Project Document No. 10. Available online at
25 www.who.int/ipcs/methods/harmonization/areas/guidance_immunotoxicity.pdf
- 26 • [A comprehensive immunotoxicity resource that includes useful information on aspects of
27 evaluating immunotoxicity studies in humans and animals.]
28
- 29 [Specific guidance on inhalation testing and reporting can be found in OECD Guidance Document 39
30 (GD 39):
31 [http://search.oecd.org/officialdocuments/displaydocumentpdf/?cote=env/jm/mono\(2009\)28&do
32 clanguage=en.](http://search.oecd.org/officialdocuments/displaydocumentpdf/?cote=env/jm/mono(2009)28&do
32 clanguage=en.)]

1 **REPORTING STUDY RESULTS**

2 Once a literature search has been conducted and the resulting database of primary (i.e.,
3 original research) studies have been evaluated with respect to strengths and limitations, the next
4 step is to display the results in a form that facilitates perusal, review, and synthesis. Most
5 applications to date have been with tabular display of results; however, presentations should not be
6 limited to this type of display and in some situations, particularly for large collections of data, a
7 graphical or some other type of figure may be a better choice.

8

9 ***Evidence Tables***

10 Evidence tables present information from the collection of studies related to a specific
11 outcome or endpoint of toxicity; for example, an evidence table for liver toxicity may include
12 studies which evaluated changes in liver enzyme levels or diagnosis of liver disease (epidemiology
13 studies), or increased liver weight or histopathological effects (animal studies). Included in the
14 table are the studies which have been judged adequate for hazard identification following the
15 principles outlined in Step 3 (see 4.3, *Reporting Study Results* of the IRIS *Preamble*). Evidence tables
16 display findings of informative studies evaluating a relevant exposure scenario (taking into
17 consideration route, timing, and dose). A key point is that evidence tables display the available
18 study results, and are not restricted to those which observed ‘statistically significant’ or ‘positive’
19 associations.

20 The studies considered to be informative will depend on the extent and nature of the
21 database for a given chemical, but may encompass a range of study designs and include
22 epidemiology, toxicology, and, other toxicity data when appropriate. Consequently, evidence tables
23 may be organized differently, when compared across assessments, depending on the data at hand;
24 for example, it may make sense to organize studies by route and duration of exposure, or by specific
25 endpoints within a toxicity type. If the database is extensive, the evidence tables may be organized
26 into two or more tiers based on the relevance and quality of the studies applied in the hazard
27 determination. Below are templates for evidence tables summarizing findings of observational
28 epidemiology and animal toxicology studies; as noted, these should not be considered fixed
29 structures, and may be adapted to best suit the database for a given chemical.

1 **Human Evidence Tables:**

Table F-10. Template for Reporting Results From Observational Epidemiology Study

Reference and Study Design	Results
<i>Outcome</i>	
<p>[Reference] (location)</p> <p>Study design, time period, description of study population (including sample size), Exposure assessment and estimates Outcome measure</p> <p>Related references (i.e., earlier publications of a cohort with exposure measurement details)</p>	<p>Prevalence of outcome (if applicable, i.e., cohort studies) Prevalence of exposure (if applicable, i.e., case-control studies) Effect estimates (and variability measure (e.g., Beta and standard error, odds ratio and 95% confidence interval)</p> <p>Include results from analysis of exposure as a continuous measure and a categorical measure [if applicable]]</p>

2

3 When constructing the evidence tables for human studies the following should be considered:

4

5 *Study Design*

- 6 • Order of the presentation of data in ‘Study Design’ column is flexible but must be consistent
- 7 throughout tables in document.
- 8 • Study size may be the overall number of participants, or preferably, the number in each
- 9 exposure or outcome group.
- 10 • Description of comparison groups may include population from which they were selected and
- 11 prevalence of important potential confounders relevant to the endpoint of concern (e.g., %
- 12 male, mean age, % smokers).
- 13 • Exposure estimate format will vary according to study; it is helpful to have some measure of
- 14 both average (such as median) and upper end (such as 90th percentile) in each comparison
- 15 group (such as exposed and unexposed, or cases and controls).
- 16 • If multiple dose metrics are provided (for example, both cumulative and peak exposure), all
- 17 may be presented in the table or selected metric(s) may be presented with a note that multiple
- 18 metrics were considered. It may be helpful to convert exposure metrics in order to compare
- 19 results between studies; if so, provide the conversion calculation as a table footnote.

20

21 *Results*

- 22 • If few or no quantitative results are reported, a qualitative description of results may be
- 23 provided using brief sentences or phrases.
- 24 • The effect measure(s) reported will depend on study design; for example, a mortality follow-up
- 25 study may present standardized mortality ratios, while a case-control study may present odds
- 26 ratios. Other examples include β coefficients from a regression model, risk ratios, or hazard
- 27 ratios.

- As noted above, positive and negative results should be displayed regardless of ‘statistical significance.’ If available, however, there should be some indication of the variability in the result (such as a 95% confidence interval).
- Include whether or how potential confounding variables were considered or adjusted for in the analysis.

Animal Evidence Tables:

Table F-11. Template Option 1 for Reporting Results From Animal Toxicology Studies

Reference and Study Design	Results					
<i>Descriptor of Effect</i>						
Reference	<i>[effect] (percent change compared to control)</i>					
species, strain, n /sex/group doses (converted doses)	M	0	5	10	15	20
		-	3%	7%	20%*	40%*
exposure route and details age and duration of exposure	F	0	6	12	17	23
		-	3%	7%	20%*	40%*
Species, strain, n /sex/group, (describe the chamber type; e.g., dynamic nose-only)	M	0	5	10	15	20
		-	3%	7%	20%*	40%*
Exposure regimen (e.g., 6 h/day, 5 days/wk for 13 weeks)	F	0	6	12	17	23
		-	3%	7%	20%*	40%*
Test article (substance used to generate the atmosphere)						
Analytical concentrations (in mg/m ³ ; do not report target or nominal concentrations)						
MMAD (GSD): (aerosol only)						
Other critical information (e.g., sections of nasal turbinates examined)						

Percent change compared to control = (treated value – control value) ÷ control value x 100

* Statistically significant (p<0.05) based on analysis by study authors

Table F-12. Template Option 2 for Reporting Results From Animal Toxicology Studies

Reference and Study Design	Results					
<i>Descriptor of Effect</i>						
Reference	<i>[effect] (percent change compared to control)</i>					
species, strain, n /sex/group doses (converted doses)	Male			Female		
	exposure route and details age and duration of exposure	0	-		0	-
	5	3%		6	3%	
	10	7%		12	7%	
	15	20%*		17	20%*	
	20	40%*		23	40%*	

* Statistically significant (p<0.05) based on analysis of data conducted by study authors.

Percentage change compared to control = (treated value – control value) ÷ control value x 100.

1 When constructing the evidence tables for animal studies, the following should be considered:

2
3 *Study Design*

- 4 • The organization of the information in ‘Study Design’ column is flexible (i.e., species, duration,
5 route) but must be consistent throughout tables in document.
- 6 • Details about species and number of test subjects should be presented as species, strain,
7 n/sex/group.
- 8 • In the study design column, report administered doses, as specified in the study, and converted
9 doses (when necessary) in mg/kg-d or mg/m³. Do not adjust for intermittent dosing. In the
10 results column, report converted doses only.
- 11 • Present average doses administered (converted from applied doses, using appropriate factors)
12 [e.g., 0, 1.0, 2.5, 3.9 mg/kg-d]. Do not use ‘or’ or ‘and’ before last dose (i.e., not 0, 1.0, 2.5, and 3.9
13 mg/kg-d). Use at least two significant figures, except when presenting whole numbers (e.g., 0,
14 2.5, 5, 10, 112, 1,024).
- 15 • If converted doses are different in males and females present as: 0, 1, 2, 3 mg/kg-d in males; 0,
16 4, 5, 6 mg/kg-d in females.
- 17 • Provide the dose conversion calculation as a table footnote; note if authors reported the dose
18 conversion.

19
20 *Results*

- 21 • If a study reports an effect but does not provide quantitative data, a qualitative description of
22 the observed result must be provided, as a brief sentence or phrase. For example, “treatment-
23 related histopathological changes were not observed”.
- 24 • For continuous data, report the percent change compared to control (generally, round to whole
25 percent unless one decimal point is needed).
- 26 • Provide the percent change formula as a table footnote: “(Treatment Mean – Control Mean)/
27 Control Mean”. Decreases calculated in this manner will have negative signs to sufficiently
28 describe the direction of the change in effect. Do not create confusion by including descriptions
29 such as “Decrease in...” or “Increase in...” before the numerical value.
- 30 • For quantal data, present incidence and number at risk (e.g., 0/20, 5/20, etc), percent (if
31 needed), and/or percent relative to control (if needed)
- 32 • Provide information for results that were not statistically significant but demonstrated an
33 increase or decrease that was biologically relevant
- 34 • In specifying the effect, simply name the effect (e.g., Rotorod latency). Do not qualify the ‘effect’
35 with descriptions such as “change in...,” “increase in...,” or “decrease in...” (e.g. “liver weight,” not
36 “increased liver weight”.) Also, do not use arrows to describe the direction of change in the
37 effect observed.
- 38 • For statistical tests comparing treatment groups to control, remember to:
 - 39 – state when statistical analysis designations are based on analysis conducted by study
40 authors;
 - 41 – document in study design cell when the study did not report statistical comparison to
42 control; and/or
 - 43 – State in a table footnote when statistical tests are performed by EPA.

1 **Evaluating the Overall Evidence of Each Effect**

2 Hazard identification involves the integration of evidence from human, animal, and
3 mechanistic studies in order to draw conclusions about the hazards associated with exposure to a
4 chemical. In general, evidence is integrated in the context of Hill (1965), which outlines aspects —
5 such as consistency, strength, coherence, specificity, does-response, temporality, and biological
6 plausibility — for consideration of causality in epidemiologic investigations that were later
7 modified by others and extended to experimental studies (U.S. EPA, 2005a).

8 All results, both positive and negative, of potentially relevant studies that have been
9 evaluated for quality are considered (U.S. EPA, 2002). This requires a critical weighing of the
10 available evidence (U.S. EPA, 2005a; 1994), but is not to be interpreted as a simple tallying of the
11 number of positive and negative studies (U.S. EPA, 2002). Hazards are identified by an informed
12 and expert evaluation and integration of the evidence. The sections that follow discuss evidence
13 integration for human, animal, and mechanistic data with the ultimate goal of integrating across
14 these evidence streams to answer the fundamental question of: **Does exposure to chemical X
15 cause hazard Y?**

16 **SYNTHESIS OF OBSERVATIONAL EPIDEMIOLOGY EVIDENCE**

17 Focus of this section

18 Studies in humans may include epidemiologic studies, case studies, and, more rarely,
19 controlled human exposure studies. While all of these study types may be included in an IRIS
20 assessment, epidemiology studies are the predominant source of human evidence for most IRIS
21 assessments. Therefore, this section is focused on the synthesis of evidence from epidemiology
22 studies.

23 Evaluation of epidemiologic evidence

24 The synthesis of epidemiologic evidence and conclusions regarding summary descriptors
25 focuses on whether and to what degree the collective evidence supports a conclusion that there is
26 an *association* between the exposure and a health outcome. That is, the goal is to answer the
27 question, “Is there evidence to conclude that an association or lack of an association exists between
28 an exposure and a health outcome, for which reasonable alternative explanations (e.g., reverse
29 causation, chance, bias, or confounding) are judged to be unlikely?”

30 The IRIS *Preamble* describes the framework for weighing the evidence from epidemiologic
31 studies. The *Preamble* states that, “for each effect, the assessment evaluates the evidence from the
32 epidemiologic studies as a whole to determine the extent to which any observed associations may
33 be causal.” While the *Preamble* refers to the concept of causality here, the evaluation of available
34 studies involving humans constitutes one line of evidence in the process of drawing an overall
35 conclusion regarding causality. In the context of an IRIS hazard evaluation, determinations of
36 causality involve consideration of the weight of evidence from all available sources, including
37 human, animal and mode of action (MOA) studies. Although a causal conclusion can be based on
38 human evidence alone, evidence from animal and MOA studies can add weight to a less robust set of
39 studies in humans.
40
41

1 Most epidemiologic studies used for risk assessment are non-experimental in design, in that
2 the investigator generally does not control exposures or intervene with the study population.
3 Broadly, epidemiologic studies are observational in nature and test specific hypotheses and
4 evaluate associations between exposures and health outcomes. These analytical studies fall into
5 several categories: e.g., cross-sectional, cohort and case-control studies. Each study design can
6 make an important contribution to an overall conclusion regarding an association, although any
7 particular design will have a specific interpretation with regard to individual aspects of the weight
8 of evidence evaluation. For example, a cross-sectional study may be less informative regarding the
9 temporal relationship between exposure and a health outcome, but it can be highly informative
10 about an association if the health response is immediate, rather than delayed. Case studies
11 involving one or a small number of affected individuals highlight potential toxicity of an exposure
12 but are the least informative in an overall evaluation of association. While controlled human
13 exposure studies, like clinical trials, offer advantages because of their experimental design, they
14 may be less informative for hazard evaluations focused on long-term low level exposures, or health
15 outcomes that occur many years after an exposure occurred. Properly interpreted, all types of
16 study designs may contribute to the weight of evidence concerning an association. The process of
17 weighing the evidence from human studies builds on the conclusions regarding the quality of
18 individual studies. Each study, including both those that do and do not show an association
19 between exposure and health outcome, is evaluated for study quality and considered as part of the
20 weight of evidence evaluation.

21 Aspects suggesting causality

22 This section discusses “aspects”¹⁴ of an association that suggest causality, drawn from Hill
23 (1965), elaborated by Rothman and Greenland (1998), and referred to in other risk assessment
24 documents such as those developed by the Environmental Protection Agency (U.S. EPA, 1994, 2002,
25 2005a), U.S. Surgeon General (DHHS, 2004;) and the Committee on Evaluation of the Presumptive
26 Disability Decision-Making Process for Veterans (Samet and Bodurow, 2008). The 1964 Surgeon
27 General’s report on tobacco smoking discussed criteria for the evaluation of epidemiologic studies,
28 focusing on consistency, strength, specificity, temporal relationship, and coherence (HEW, 1964).
29 These aspects of causality are briefly described in the *Preamble*, and in more detail here.

30 First, greater *strength of association* lends greater confidence that the association is not due
31 to chance or bias. However, while an association may be of small magnitude (due to factors such as
32 low potency or a low level of exposure in the study population), a widespread exposure could lead
33 to a significant public health burden, as seen for air pollution and risk of cardiovascular disease.
34 ‘Strength’ encompasses not only magnitude of the association, but statistical confidence in effect
35 measure estimates. Higher precision, as reflected by narrow confidence bounds or smaller
36 standard errors, also adds confidence in the observed association.

37 Second, *consistency* of the association across studies is another important weight of
38 evidence consideration. Observing an association in different study types, study populations, and
39 exposure scenarios makes it less likely that the association is due to confounding or other factors
40 specific to a given study, or is confined to a specific susceptible population. Characterizations of
41

¹⁴ The “aspects” described by Sir Austin Bradford Hill (Hill, 1965) have become, in the subsequent literature, more commonly described as “criteria.” The original term “aspects” is used here to avoid confusion with “criteria” as it is used, with different meaning, in the Clean Air Act.

1 consistency should distinguish between heterogeneity of findings which may be explained (e.g., due
2 to differences in populations, exposure measures, ranges of exposures, potential co-exposures, and
3 other factors specific to the exposure and health outcomes under evaluation) and unexplained
4 variability suggesting potentially spurious findings (White et al., in press). For example, one would
5 not necessarily expect to find identical results of exposure to an endocrine disruptor among those
6 exposed prenatally versus during adulthood. This difference in timing of exposure is an expected
7 source of heterogeneity in findings, rather than a signal that the findings are ‘inconsistent’. In
8 addition, a group of studies should not be characterized as ‘inconsistent’ if the results are not all
9 statistically significant, or if effect measures are of different magnitudes, but are predominantly
10 negative, null or positive.

11 The third aspect of *specificity* refers to one (or a few) causes for one health outcome. This
12 aspect draws on Koch’s postulates for infectious causes of disease, but may be less relevant in other
13 contexts. For example, many environmental exposures may have carcinogenic action, but all
14 contribute to a single health outcome. Conversely, a single exposure may be linked to a range of
15 health outcomes. Thus, specificity may lend greater confidence in an association when it exists, but
16 should not detract from an association if it does not.

17 *Temporality* is generally agreed to be the only aspect which is necessary for an association
18 to be causal. That is, the exposure must precede the health outcome. In terms of epidemiologic
19 studies, temporality is often cited as a main weakness of cross-sectional study designs. However, in
20 evaluating a body of evidence, other study designs which do inform temporality can lend strength
21 to the group of studies as a whole.

22 The *biologic gradient* or *exposure-response relationship* is another aspect which lends
23 confidence to an observed association. Observing incremental changes in the risk of a health
24 outcome which correspond to incremental changes in the exposure of interest, is a powerful
25 argument against a spurious association, since that would necessitate a third (uncontrolled) factor
26 which changes in the same manner (direction and magnitude) as the exposure of interest. Although
27 this aspect is sometimes interpreted to imply that a monotonic relationship is required, the true
28 exposure-response curve may indeed be non-linear. In evaluating a body of epidemiologic studies,
29 it may be that any one study only includes a portion of the range of exposure. Piecing together
30 evidence from multiple studies may yield a fuller understanding of the response and the shape of
31 the exposure-response curve over the full range of exposures. Similarly, an observed lack of
32 response in any one study does not imply a lack of an association between exposure and a health
33 outcome. This may be due to exposure misclassification, or the exposures in a study were below
34 some threshold for response, or that the range of exposures was too narrow to differentiate
35 between groups (White et al., in press).

36 The next group of aspects comprises biologic plausibility, coherence and analogy. These
37 were originally separate (related) aspects as laid out by Bradford-Hill, but more recently are seen
38 as variations of a common theme. Biologic plausibility, coherence and analogy are addressed when
39 weighing the totality of evidence including human, animal and mode of action. Generally, the
40 association between exposure and a health outcome should be consistent with (or not violate)
41 known scientific principles or other existing information from epidemiology, toxicology, clinical
42 medicine, or other disciplines. A difficulty in applying these aspects is the reliance on current
43 information, or the ‘state of the science.’ Associations in the epidemiologic literature may be

1 observed well in advance of experiments being performed or insight into mechanism or mode of
2 action, but confidence that an association exists is strengthened by these aspects.

3 The final aspect is the existence of *natural experiments*, occurring when environmental
4 conditions change in such a way as to mimic a controlled experiment or randomized trial—such as
5 a change in workplace standards which reduces occupational exposure, or change in medication
6 use with the introduction, or withdrawal, of a drug from the market. When such a change in the
7 exposure is followed by changes in the risk of a health outcome of interest, this provides greater
8 confidence that an association exists.

9 As discussed in the U.S. EPA Integrated Science Assessments for particulate matter (U.S. EPA
10 2009) and carbon monoxide (U.S. EPA 2010), although these aforementioned aspects provide a
11 framework for assessing the evidence, they do not lend themselves to being considered in terms of
12 simple formulae or fixed rules of evidence leading to conclusions about causality (Hill, 1965). For
13 example, one cannot simply count the number of studies reporting statistically significant results or
14 statistically nonsignificant results and reach credible conclusions about the relative weight of the
15 evidence and the likelihood of causality (U.S. EPA 2009, 2010). Rather, these aspects are taken into
16 account with the goal of producing an objective appraisal of the evidence, which includes weighing
17 alternative explanations. In addition, it is important to note that the aspects of causality cannot be
18 used as a checklist, but rather are used as a guide to help determine the weight of the evidence for
19 inferring causality. (U.S. EPA 2009, 2010) In particular, not meeting one or more of the aspects
20 does not preclude a determination of causality [(U.S. EPA 2009, 2010), see discussion in (CDC,
21 2004)]. Scientific judgment is needed to evaluate individual study quality and to weight the overall
22 body of evidence.

23 24 Evaluation of potential alternative explanation of observed epidemiologic associations

25 In evaluating epidemiologic studies, consideration of many study design factors and issues
26 must be taken into account to properly inform their interpretation and determine whether
27 observed associations are likely to represent the truth or if there are reasonable alternative
28 explanations (e.g. biases or other threats to internal validity). Such alternative explanations include
29 “reverse causality” where the health outcome precedes exposures, chance, bias (selection bias and
30 information bias) and confounding, and these alternatives are carefully considered in the
31 evaluation of the aspects of causality and of the evidence as a whole.

32 As noted earlier, a logical time sequence (temporality) is an essential aspect of causality and
33 ensures that “reverse causation” is unlikely. Chance can always be a potential explanation for the
34 results in any collection of studies but is less likely as more studies are accrued that have similar
35 observations across different settings, study designs and populations.

36 A further key consideration is evaluation of the potential effects of selection bias which may
37 occur when study groups (exposed and unexposed, cases and controls) are not sufficiently
38 comparable. Selection bias may alter epidemiologic findings when participation or follow-up rates
39 are related to the probability of exposure and to the outcome of interest. For example, effect
40 estimates that are based on a comparison of exposed workers to a general population (e.g.,
41 standardized mortality ratios) may be affected by a selection bias called the healthy-worker effect,
42 because the baseline health of workers is typically better than the baseline health of the population
43 as a whole. This type of selection bias could obscure a truly larger effect of toxicant exposure in
44 analyses based on “external” comparisons with mortality in the general population. Although this

1 type of bias would not influence analyses using “internal” or matched comparison groups, other
2 types of healthy worker effect bias should also be considered for these types of studies. Selection
3 bias can lead to either an overestimate or underestimate of risk, and the potential direction and size
4 of the bias must be considered when deciding whether individual studies are given more weight or
5 less weight for a hazard evaluation. Studies where selection bias is less of a concern are typically
6 given more weight.

7 Another key consideration is evaluation of the potential effects of measurement error
8 which can lead to information bias. One example is the uncertainty associated with using surrogate
9 exposure metrics to represent the actual exposure of an individual or population. This exposure
10 measurement error can be an important contributor to variability in epidemiologic study results.
11 Exposure measurement error can lead to misclassification (a type of information bias) that can
12 under- or over-estimate epidemiologic associations between exposures and health outcomes,
13 distort exposure-response relationships and widen confidence intervals around effect estimates
14 (i.e. decrease precision). There are several components that contribute to exposure measurement
15 error in epidemiologic studies, including the difference between true and measured concentrations
16 and the use of average population exposure rather than individual exposure estimates. The
17 importance of exposure misclassification varies with study design and is dependent on the spatial
18 and temporal aspects of the available data. For a given set of epidemiologic studies informing a
19 hazard evaluation, results from studies with more accurate exposure estimates (minimizing
20 exposure misclassification) are given more weight, barring other serious design limitations (e.g.,
21 selection bias). Generally, exposure misclassification, when nondifferential, results in a bias toward
22 the null and this is a potential explanation for relatively small effect estimates or for variability in
23 results across studies with different degrees of exposure misclassification.

24 Confounding is a type of bias that leads to “... a confusion of effects. Specifically, the
25 apparent effect of the exposure of interest is distorted because the effect of an extraneous factor is
26 mistaken for, or mixed with, the actual exposure effect (which may be null)” (Rothman and
27 Greenland, 1998). A confounder is a common cause of both the exposure and the health outcome—
28 thus, it is associated with both the exposure and the health outcome, but is not an intermediary
29 between the two. For example, confounding can occur between correlated toxicants (such as
30 pesticides used in a mixture) that are also associated with the same health outcome. Knowledge of
31 the broader literature on risk factors for the health outcome is important. Scientific judgment is
32 needed to evaluate the likely sources and extent of confounding, together with consideration of
33 how well the existing constellation of study designs, results, and analyses address this potential
34 threat to inferential validity. The ability to statistically adjust for confounding in an epidemiologic
35 study is dependent on the ability to identify and measure potential confounders. Consistency in
36 reported effect estimates across multiple studies, conducted in various settings using different
37 populations or exposures, can increase confidence that unmeasured confounding is an unlikely
38 alternative explanation for the observed associations. Such consistency also reduces the likelihood
39 of chance as an alternative explanation through the accumulation of a larger body of similar
40 evidence, as noted above. The observations of exposure-response trends across different studies
41 similarly reduce the likelihood that chance, bias, or confounding can explain the observed
42 association. Studies in which confounding is a minimal concern are typically given more weight.

43
44 Summary descriptors of epidemiologic evidence

1 The considerations described above are consistent with guidelines for systematic reviews
2 that evaluate the quality and strength of evidence. Confidence that a true association between
3 exposure and a health outcome exists is increased if the effect estimates across multiple studies are
4 judged to be consistent or when apparent inconsistencies may be explainable due to differences in
5 study designs, populations studied, exposure concentration or timing, and/or issues of potential
6 confounding, information bias and selection bias. Confidence is also increased if there is evidence
7 of an exposure-response relationship or when the magnitude of effects is considered sufficient to
8 conclude that a role of residual bias is negligible. Greater weight is given to the aspects of
9 consistency, strength of association, temporality, and biologic gradient (exposure-response
10 relationship) when assessing the epidemiologic evidence.

11 To make clear how much the epidemiologic evidence contributes to the overall weight of
12 the evidence, the assessment may include a descriptor such as “*Sufficient epidemiologic evidence of*
13 *an association consistent with causation*”, “*Suggestive epidemiologic evidence of an association*
14 *consistent with causation*”, “*Inadequate epidemiologic evidence to infer a causal association*”, or
15 “*Epidemiologic evidence consistent with no association*” to characterize the epidemiologic evidence
16 of each outcome. While each epidemiologic database is distinct and requires specific judgments be
17 made on the relative merits of those studies, some examples of the constellation of the aspects of an
18 association that suggest causality are provided below.

19 “*Sufficient epidemiologic evidence of an association consistent with causation*”

20 This descriptor is appropriate when the epidemiologic evidence is sufficient to establish an
21 association between exposure and a health outcome for which reasonable alternative explanations,
22 such as confounding, information bias and selection bias, are judged to be unlikely. Evidence of a
23 consistent finding of an association between exposure and a health outcome along with evidence of
24 an exposure-response relationship contribute considerable weight toward evidence of an
25 association. Such evidence is increased when the association is relatively strong but may not
26 necessarily be diminished when the observed associations are small in magnitude. Likewise,
27 evidence of a coherent temporal relationship allowing for disease latency (where applicable) adds
28 weight to this conclusion but the absence of such information does not necessarily detract from the
29 conclusion.

30 “*Suggestive epidemiologic evidence of an association consistent with causation*”

31 This descriptor is appropriate when the epidemiologic evidence is suggestive of a causal
32 association between exposure and a health outcome, but where there is less certainty that
33 alternative explanations such as selection bias, information bias, and confounding, have been
34 addressed. This descriptor covers a spectrum of evidence associated with varying levels of concern
35 for a health outcome. Depending on the extent of the database, additional studies may or may not
36 provide further insights.

37 An example of an aggregation of suggestive epidemiologic evidence might include apparent
38 unexplained inconsistency of risks across studies with varying strength of the association but
39 multiple studies reporting exposure-response relationships and a coherent temporal relationship
40 allowing for disease latency. Another example of a constellation of suggestive epidemiologic
41 evidence might include repeated observations of increases in risk across studies, especially for high
42 exposures, but only a relatively modest overall strength of the association and limited evidence of
43 an exposure-response relationship from one or more high quality studies.

1 *“Inadequate epidemiologic evidence to infer a causal association”*

2 This descriptor is appropriate when the epidemiologic evidence is judged inadequate for
3 describing an association. An example of inadequate epidemiologic evidence might include
4 explained heterogeneity of the observed increases in risk across studies with the majority of
5 studies having relatively poor quality exposure assessment methodology and reporting null results
6 contrasted with a single large high quality study with clear evidence of an exposure-response
7 relationship.

8 Additional high quality studies generally would be expected to provide further insights.
9 Additional supportive evidence demonstrating exposure-response relationships might lend more
10 confidence that associations reported in epidemiologic studies are not due to alternative
11 explanations, while additional evidence from other high quality studies showing a lack of exposure-
12 response or that previous findings may be due to confounding might tip the balance in another
13 direction.

14 *“Epidemiologic evidence consistent with no association”*

15 This descriptor is appropriate when the available data are considered robust for deciding
16 that there is no basis for human hazard concern. An example of evidence suggestive of no
17 association would include a consistent pattern of results indicating a lack of an association across a
18 large number of studies that had adequate statistical power spanning different exposure patterns
19 and exposure ranges, including high exposure levels, and evidence of the absence of exposure-
20 response relationships even at high exposure levels.

21
22 **SYNTHESIS OF ANIMAL TOXICOLOGY EVIDENCE**

23 In IRIS assessments, human data are generally preferred for hazard identification because
24 these data are more relevant in the assessment of toxicity to human health and avoid the
25 uncertainty associated with potential interspecies differences when using animal data. However,
26 many chemical databases contain little or no human data; thus, IRIS assessments frequently rely on
27 available animal data in order to determine potential chemical hazards. In the absence of human
28 data, well-conducted animal toxicology studies can support the identification of hazards. Animal
29 data are used under the assumption that toxicity is conserved across species, in that effects
30 observed in animals would be expected to occur in humans (U.S. EPA, 1998c; 1996; 1991). This
31 section discusses how to approach synthesis of evidence from animal toxicology studies and
32 focuses on whether and to what degree the collective evidence supports a conclusion that there is
33 an association between chemical exposure and an effect.

34 In contrast to observational epidemiology studies that do not control exposures or
35 intervene with the study population, experimental animal toxicology studies are designed to
36 control exposure and environmental conditions. These studies permit the use of study design to
37 control the number and composition (age, gender, species) of test subjects, the levels of doses
38 tested, and the measurement of specific responses. Use of a designed study typically leads to more
39 meaningful statistical conclusions than an uncontrolled observational study where additional
40 confounding factors must also be considered for their impact on the conclusions. Thus, the
41 observed responses in animals are expected to be due to chemical exposure. However, dose-
42 response relationships observed in animal toxicology studies are often at much higher doses than
43 would be anticipated for humans.

1 Animal toxicology studies fall into two broad categories: (1) general toxicology studies
2 designed to evaluate a comprehensive array of endpoints following varying durations of exposure
3 (e.g., chronic, sub-chronic, short-term, or acute) or (2) toxicology studies designed to evaluate a
4 specific type of toxicity: e.g., neurotoxicity, immunotoxicity, reproductive toxicity, and
5 developmental toxicity.

6 As noted in the *Preamble* and described in the Synthesis of Epidemiology Studies section,
7 several aspects of causality discussed by Hill (1965) are pertinent to the interpretation of animal
8 evidence: consistency of response, exposure-response relationship, strength of response,
9 specificity of response, biological plausibility and coherence, and temporality (U.S. EPA, 2005a,
10 2002, 1994). These considerations, as they relate to synthesizing animal toxicology evidence in
11 animals, are further described below.

12 ***Principles and Considerations for Writing a Synthesis of Animal Evidence***

13 *NOTE: In general these considerations apply to both human and animal data; however, for purposes*
14 *of providing a simple example, this section is focused on animal evidence.*

15
16 For each health effect, the evidence from animal experiments is evaluated to determine the
17 extent to which this evidence indicates a potential for effects in humans. The starting points for a
18 synthesis of the data for a given health effect (e.g., hepatic, immune system, cancer) are the
19 following: (1) the evidence table(s) as developed in the Reporting Study Results subsection in the
20 Evaluation and Display of Individual Studies section, (2) the actual papers or reports captured in
21 the evidence tables, (3) information on study quality as documented in the Evaluation and Display
22 of Individual Studies section, and (4) other information not summarized in the evidence table that
23 contributes to the evidence of an association between exposure to the chemical and the given
24 health effect. This other information could include short-term and acute experimental animal
25 studies, and data from studies using routes other than oral, inhalation, or dermal. Keep in mind that
26 while the evidence tables provide a useful framework for starting the evaluation, they are not
27 sufficient for completing the synthesis. You will need the additional content provided in the papers
28 and reports themselves to prepare the synthesis.

29 In general, the evidence table and accompanying synthesis text should be complementary,
30 and provide a comprehensive and critical evaluation of the animal evidence. The synthesis should
31 not be a text version of the information contained in the evidence table. For example, the synthesis
32 text should not repeat study design details provided in the evidence tables, but should discuss
33 strengths and limitations with studies (identified and documented in the Evaluation and Display of
34 Individual Studies section) that would influence interpretation of the study results. To the extent
35 possible, the information in the evidence table and text should be presented in the same order.
36 However, it is important to remember that the synthesis should be a discussion from the
37 perspective of the evidence for particular effects *across* studies, not by study, with the caveat that
38 you will generally discuss evidence following oral, inhalation, or dermal exposures of chronic
39 durations (i.e., more relevant to estimating potential toxicity to humans following chronic exposure
40 to chemicals) prior to describing results in shorter duration studies. However, developmental and
41 reproductive toxicology studies may provide pertinent evidence resulting from short-term
42 exposures during a critical period of development.

43 There is no formula for writing a synthesis of the animal evidence. The approach for
44 organizing the information will depend on the nature and extent of the literature for a given

1 chemical and health effect. Potentially relevant studies have been evaluated for quality (as
2 identified and documented in the Evaluation and Display of Individual Studies section), and studies
3 of higher quality are given more weight than those of low quality. This is mirrored in the
4 development of evidence tables, which also capture findings from the most pertinent and higher
5 quality studies without consideration of the presence or absence of an effect. All results, both
6 positive and negative, are considered (U.S. EPA, 2002) and discussed.

7 In comparing and contrasting results across studies, evidence evaluations of study quality
8 are further considered and discussed due to potential impacts on the interpretation of results (e.g.,
9 may explain differences in the results for a given endpoint). For example, could differences in test
10 article preparation or delivery vehicle between studies account for differences in the reported
11 results? Was the study adequately powered to identify an effect associated with a chemical
12 exposure? Could co-exposures alter the response in one study versus another? Issues with study
13 power, design, or conduct may limit the ability to draw conclusions about chemical-related effects
14 when a single study is considered in isolation; when considered in the totality of studies that
15 examine a given health effect these flawed studies can still add qualitative evidence for an effect
16 associated with chemical exposure. Additionally, historical background levels of effects (if
17 available) should be considered. While comparisons to concurrent controls are preferred when
18 identifying effects, the use of appropriate historical control data may be informative when a
19 particular effect is rare

20 The write-up should address the **consistency (including any lack of consistency) of the**
21 **results** across studies. Consistent results across species, strains, sexes, life stages, routes of
22 exposure, and exposure regimens and durations increases confidence that similar results would
23 occur in humans. While consistency across higher quality studies is preferred, consistency of an
24 effect across studies of varying quality and statistical power may provide qualitative information
25 about a given effect. Inconsistency of effects among studies and/or species that cannot be explained
26 by differences in timing and/or magnitude of exposure or toxicokinetics/metabolism can decrease
27 confidence. As discussed in the *Preamble* (Section 5.2), distinguishing between **conflicting evidence**
28 (that is, mixed positive and negative results in the same sex and strain using a similar study
29 protocol) and **differing results** (that is, positive results and negative results in different sexes or
30 strains or using different study protocols) is also important. Ask yourself why valid results are
31 inconsistent and include information that could reconcile the differences in your evaluation. For
32 example, did the “negative” study use an exposure range that was too low (e.g., were the highest
33 exposures in the “negative” study similar to the range that produced no exposure-related response
34 in the “positive” study)? Where you have a positive and negative study for a specific endpoint (e.g.,
35 neurotoxicity), investigate whether the negative study was adequately designed to look for that
36 endpoint. Can differences in response be explained by differences in toxicokinetics across species?
37 Refer to Agency guidance, where available, for additional information on evaluating specific health
38 effects.

39 Your discussion should also indicate whether effects showed an **exposure-response**
40 **relationship**, i.e., whether the incidence and/or intensity of response changes in an orderly
41 manner as a function of exposure. Note, however, that the exposure-response relationship need not
42 be monotonic. U-shaped (or inverted U-shaped) exposure-response functions are not uncommon in
43 toxicology. In addition, information on the **strength (or magnitude) of the response** (in general
44 terms) and the **exposure range where effects are first observed** should be provided. Confidence

1 in an association between a chemical exposure and a given health effect is increased when an
2 exposure-response relationship is demonstrated and when the magnitude of effect is large. Also, be
3 precise in characterizing the strength of the association between chemical exposure and effect. For
4 example, the word “demonstrates” indicates a relatively strong association and should be used with
5 caution. Words such as “suggests” or “indicates” are appropriately used when the evidence for an
6 association is not as strong (e.g., an association based on a small number of studies or less
7 consistent results).

8 If related effects in a target organ are observed (e.g., changes in serum enzymes that are
9 markers of liver damage, increased liver weight, and liver histopathology), it is worthwhile to note
10 the **coherence** of these related effects as well as **characterize the exposure ranges at which**
11 **these effects** were observed. Coherence of the exposure ranges for related effects strengthens the
12 biological plausibility for a given effect as well as provides a more complete picture of the toxicity
13 associated with exposure to a chemical. For example, changes in liver enzymes are likely to occur at
14 earlier time points and/or at lower exposures than histopathologic changes of the liver.

15 Another criterion that is important in interpreting data is the **temporal relationship**
16 between exposure and effect. Temporality is generally assumed in animal toxicology studies since
17 the exposure precedes measurement of effects. However, temporal considerations also need to be
18 evaluated in the context of the observed effects. That is, the exposure should precede the effect at
19 an interval that is consistent with what is known about the toxicokinetics and mode of action of the
20 chemical. It may be the case, however, that higher exposures produce a shorter latency to effect
21 than do lower exposures. Additionally, exposure-response relationships may vary due to temporal
22 considerations. For example, if a study’s dose groups result in premature mortality at higher doses,
23 you may not observe an effect with increased frequency (or severity) at the higher doses because
24 the animal died prior to the time needed to develop an effect with a long latency. Similarly, when
25 considering cancer effects, the number of adenomas may decrease with increasing dose as they
26 progress into larger tumors or progress to carcinomas. In some cases, initial effects may disappear
27 as the pathogenesis of a lesion evolves or resolves (and early effects may not even be observed
28 depending on the doses used and the resulting exposure-response relationship). How different
29 studies illustrate the development of a lesion, weaving in considerations of the exposure-response
30 relationship and temporality can increase or decrease the biological plausibility of a given effect.

31 **Biological plausibility and coherence** are also evaluated; although these aspects are best
32 considered when integrating evidence across human, animal, and mechanistic evidence streams.
33 Several types of information should be considered (e.g., toxicokinetics/metabolism, similarity of
34 effects, exposure-response relationships, mode-of-action, and temporal relationships) when
35 determining the likelihood of the occurrence of effects in humans based on observations in animals.
36 All of this information must be weighed in light of the known heterogeneity of the human
37 population versus the relatively inbred status of laboratory animals used in toxicology studies and
38 housed under carefully controlled environmental conditions (U.S. EPA, 2002). These concepts are
39 more fully described in the discussion of overall integration of evidence.

40 Additionally, in writing a synthesis of animal evidence, there are a couple of other
41 considerations. Keep in mind that the evaluation is not a study by study summary of the literature.
42 When discussing “significance,” be clear as to whether you mean statistical or biological

1 significance¹⁵. Biological and statistical significance are both considered when making a judgment
2 about the adversity of an observed effect. Where possible, emphasize biological significance over
3 statistical significance, or be clear that biological significance is not well understood for a particular
4 endpoint.

6 ***A Practical Example for Synthesizing Animal Toxicology Evidence***

7 The following text example highlights some key considerations in writing synthesized text
8 for a specific endpoint. Compare the text in “Draft 1” with the revisions made in “Draft 2.”

9 **Draft 1:** In several studies in rats and mice, decreased sperm count, motility, and
10 production, and an increase in morphologically abnormal sperm have been observed.
11 Decreased epididymal sperm count (approximately 50% at 1 mg/kg-day) and sperm
12 motility (approximately 20% at 1 mg/kg-day) were observed in mice exposed by gavage to
13 doses \geq 1 mg/kg-day for 42 days prior to mating to unexposed females (Mohamed et al.,
14 2010). This study also demonstrated transgenerational impacts on sperm parameters, as
15 these endpoints were also decreased in the F1 and F2 generations produced from treated
16 F0 males. Decreased epididymal sperm count (25%) and a slight increase in abnormal
17 sperm morphology were observed in rats treated with 5 mg/kg-day benzo[a]pyrene by
18 gavage for 84 days (Chen et al., 2001). A decrease in sperm motility (approximately 30%)
19 and an apparent (but not statistically significant) decrease in epididymal sperm count
20 (approximately 15%) were also observed in rats treated by gavage at 0.01 mg/kg-day for
21 90 days (Chung et al., 2011).

22 Similar effects on sperm parameters have been observed in short term oral studies and
23 inhalational studies. Significantly decreased sperm count, number of motile sperm, and
24 daily sperm production (~40% decrease from control in each parameter) were observed
25 following 10 days of gavage exposure to 50 mg/kg-day benzo[a]pyrene in rats (Arafa et al.,
26 2009). In addition, decrements in sperm parameters (specifically sperm motility, sperm
27 count, and percent morphologically normal sperm) were observed following inhalation
28 exposure to benzo[a]pyrene in rats for 60 days to 75 $\mu\text{g}/\text{m}^3$ (Archibong et al., 2008;
29 Ramesh et al., 2008). In addition, decreased sperm motility, but not sperm count, was found
30 to be decreased in rats exposed by inhalation to benzo[a]pyrene for 10 days at \geq 75 $\mu\text{g}/\text{m}^3$
31 (Inyang et al., 2003).

32
33 **Draft 2:** In several studies in rats and mice, decreased sperm count, motility, and
34 production, and an increase in morphologically abnormal sperm have been
35 reported. Alterations in these sperm parameters have been observed in different
36 strains of rats and mice and across different study designs and routes of exposure.

37 Decreases in epididymal sperm (25 to 50% compared to controls) counts have been
38 observed in SD rats and C57BL6 mice treated with 1- 5 mg/kg-day benzo[a]pyrene
39 by oral exposure for 42 or 90 days (Chen et al., 2011; Mohamed et al., 2010).

¹⁵ Biological significance is the determination that the observed effect (a biochemical change, a functional impairment, or a pathological lesion) is likely to impair the performance or reduce the ability of an individual to function or to respond to additional challenge from the agent (U.S. EPA, 2002).

1 Additionally, a 15% decrease in epididymal sperm count was observed at a dose two
2 magnitudes lower in Sprague Dawley rats exposed to benzo[a]pyrene for 90 days
3 (Chung et al., 2011). However, confidence in this study is limited as authors dosed
4 animals with 0.001, 0.01, and 0.1 mg/kg-day benzo[a]pyrene but only reported on
5 sperm parameters at the mid-dose. A short term study in mice and a subchronic
6 inhalation study in rats lend support for the endpoint of decreased sperm count
7 (Arafa et al., 2009; Archibong et al., 2008; Ramesh et al., 2008). Significantly
8 decreased sperm count and daily sperm production (~40% decrease from control in
9 each parameter) were observed following 10 days of gavage exposure to 50 mg/kg-
10 day benzo[a]pyrene in mice (Arafa et al., 2009). In addition, decrements in sperm
11 count were observed in rats following inhalation exposure to 75 $\mu\text{g}/\text{m}^3$
12 benzo[a]pyrene for 60 days (Archibong et al., 2008; Ramesh et al., 2008).

13 In addition to effects on sperm count, both oral and inhalation exposure of rodents
14 to benzo[a]pyrene has been shown to lead to decreased epididymal sperm motility
15 and altered morphology. Decreased motility of 20-30% compared to controls was
16 observed in benzo[a]pyrene-exposed C57BL6 mice ($\geq 1\text{mg}/\text{kg}\text{-day}$) and SD rats
17 ($0.01\text{ mg}/\text{kg}\text{-day}$) (Chung et al., Mohamed et al., 2010). The effective doses spanned
18 two degrees of magnitude; however, as noted above, confidence in the study
19 observing effects at $0.01\text{ mg}/\text{kg}\text{-day}$ benzo[a]pyrene (Chung et al., 2011) is limited
20 by poor reporting. A short term oral study in mice also reported significantly
21 decreased number of motile sperm (~40% decrease) following 10 days of gavage
22 exposure to $50\text{ mg}/\text{kg}\text{-day}$ benzo[a]pyrene in mice (Arafa et al., 2009). In addition,
23 decreased sperm motility was observed following inhalation exposure to $75\ \mu\text{g}/\text{m}^3$
24 benzo[a]pyrene in rats for 60 days (Archibong et al., 2008; Ramesh et al., 2008) and
25 $\geq 75\ \mu\text{g}/\text{m}^3$ for 10 days (Inyang et al., 2003). Abnormal sperm morphology was
26 observed in Sprague Dawley rats treated with $5\text{ mg}/\text{kg}\text{-day}$ benzo[a]pyrene by
27 gavage for 84 days (Chen et al., 2001) and in rats exposed to $75\ \mu\text{g}/\text{m}^3$
28 benzo[a]pyrene by inhalation for 60 days (Archibong et al., 2008; Ramesh et al.,
29 2008).
30

31 Note the following elements of the Draft 2 compared to Draft 1:

- 32 • Rather than providing the results of each study in a sentence (as in Draft 1), the Draft 2 text
33 pulls together studies on the same effect (decreased sperm count, decreased motility,
34 altered morphology) to provide a more integrated analysis of the results from multiple
35 studies simultaneously.
- 36 • Information on the magnitude of the effect (or range of magnitudes of effect) is provided.
- 37 • Studies are generally organized by duration, with longer duration studies described first
38 (within the limits of the available data).
- 39 • Study quality considerations are included in the discussion; in one instance, confidence in
40 the findings is limited because the authors only reported effects in one mid-dose group.

41
42 While no single example could capture all the elements of a synthesized summary of evidence of
43 animal toxicity, the above text highlights major considerations when writing up your synthesis of
44 the animal toxicology data.

1 **MECHANISTIC CONSIDERATIONS IN ELUCIDATING ADVERSE OUTCOME PATHWAYS**

2 Mechanistic data contribute to the hazard evaluation of empirical evidence from human and
3 animal studies by informing the following:

- 4 • The biological plausibility of a causal interpretation in humans
- 5 • The biological plausibility that animal experimental data is generalizable to humans
- 6 • The susceptibility of certain populations or lifestages

7
8 Evaluating mechanistic considerations is a critical part of weighing the evidence for hazard
9 identification. The focus of this evaluation is on adverse outcome pathways (AOPs) that encompass
10 both:

- 11 1) the toxicokinetic processes of absorption, distribution, metabolism, and excretion (ADME)
12 that lead to the formation of the active agent and carry it through its distribution to the
13 target cell, and
- 14 2) the toxicodynamic processes in the mode(s) of action (MOA[s]), leading to the adverse
15 outcome.

16
17 While not a prescribed process—the database for every chemical and endpoint will be
18 unique—the following steps may be informative in conducting the evaluation.

19 For each endpoint, the evaluation of AOPs begins by identifying:

- 20 • Information that may help identify the toxic moiety and the target site, and how the toxic
21 agent is delivered to that site. Note that the target site at which the initial biological
22 interaction occurs is not necessarily the site of the adverse effect.
- 23 • Information that may help identify key events in the hypothesized MOA(s).

24
25 This information may include both experimental and observational evidence specific to the
26 chemical and endpoint, as well as additional evidence such as:

- 27 • Information on compounds that are similar in structure, function, and/or metabolism
- 28 • Information on how the chemical may disrupt normal biological processes or interacts with
29 background aging or disease processes
- 30 • Interactions with other chemicals and/or mixtures
- 31 • Factors affecting biological susceptibility

32
33 Using this evidence, AOPs are described as sequences or networks of steps, from exposure
34 to the chemical, formation of the active agent and delivery to the target site, and the key events
35 leading to the adverse outcome.

36 Based on the evaluation of the available information, one or more of the following
37 determinations may be possible:

- 38 • Whether there is sufficient information available to specify AOP hypotheses with respect to
39 (1) and (2), above. In many cases the answer will be “no” to one or both of these due to lack
40 of data.
- 41 • Whether the ADME and/or MOA data add to the biological plausibility of the hazard being
42 evaluated.
- 43 • Whether a hypothesized AOP(s) is(are) sufficiently supported.

- 1 • Whether observed animal responses are generalizable to humans (i.e., human relevancy).
2 • Whether differences are anticipated in responses among humans, including susceptible
3 subpopulations and lifestage-specific sensitivities.
4

5 While these determinations are qualitative in scope, the evaluation of AOP(s) should flag
6 important quantitative information that may be carried over to dose-response analysis. These may
7 include:

- 8 • Dosimetry for route-to-route extrapolation.
9 • Quantitative inter- or intraspecies differences in dosimetry.
10 • Quantitative inter- or intraspecies differences in response susceptibility.
11 • The shape of the dose-response relationship below the POD that may inform the choice of
12 linear or non-linear extrapolation.
13

14 **INTEGRATION OF EVIDENCE EVALUATION**

15 *[Note: EPA is investigating the use of standard descriptors to characterize the overall weight of*
16 *the evidence for effects other than cancer. The NRC will hold a workshop in March 2013 on this*
17 *topic, and EPA will follow up with a workshop to further develop this topic. In the meantime,*
18 *the Preamble cites descriptors from EPA's 2005 Cancer Guidelines and, for effects other than*
19 *cancer, the descriptors from EPA's Integrated Science Assessments.]*
20
21

1 Dose-Response Analysis

2 SELECTING STUDIES FOR DERIVATION OF TOXICITY VALUES

3 For each health effect for which there is credible evidence of hazard, a group of studies has
 4 been identified and evaluated as part of the hazard identification (See Section on Evaluating the
 5 Overall Evidence of Each Effect). Once these studies have been identified, the basic criterion for
 6 selecting a subset for the derivation of toxicity values is whether the quantitative exposure and
 7 response data are available to compute a NOAEL, LOAEL or benchmark dose/concentration. When
 8 there are many studies, the assessment may focus on those that are more pertinent or of higher
 9 quality.

10 The relative merits of deriving toxicity values for each endpoint will depend on the size of
 11 the relevant database and various preferences as stated in the IRIS *Preamble* Section 6 as well as
 12 specific considerations appropriate for each chemical and health endpoint. All studies of sufficient
 13 quality (as evaluated in the Section on Evaluation and Display of Individual Studies) with data
 14 suitable for deriving toxicity values are considered (see the Benchmark Dose Technical Guidance
 15 sections 2.1.3, 2.1.5; U.S. EPA, 2012). All aspects of study quality evaluation for Hazard
 16 Identification are also important for dose-response (Tables F-6 and F-7). This section discusses
 17 **additional** considerations that are specific to dose-response analysis.
 18

Table F-13. Attributes used to evaluate studies for derivation of toxicity values.		
Aspects of Study	Data Characteristic	Considerations
Species studied	Human studies	Human data are preferred to reduce interspecies extrapolation uncertainties.
	Animal studies	Animal data are considered as supporting studies when adequate human studies are available, and as principal studies when adequate human studies are not available. Results from experiments using mammalian laboratory animals are favored over those conducted using non-mammalian species.
Relevance of exposure paradigm	Exposure route	Studies by a route of human environmental exposure are preferred, although a validated toxicokinetic model can also be used to extrapolate across exposure routes.
	Exposure durations	When developing a chronic toxicity value, chronic or subchronic studies are preferred over studies of acute exposure durations. There are exceptions, such as when a susceptible population or life stage is more sensitive in a particular time window (e.g., developmental exposure).

Table F-13. Attributes used to evaluate studies for derivation of toxicity values.		
Aspects of Study	Data Characteristic	Considerations
	Exposure levels	Studies with multiple exposure levels are preferred to the extent that they provide information about the shape of the exposure-response relationship (BMDTG 2.1.1).
Potential selection bias	Representativeness of the study sample to the target population and the potential for selection to be based jointly on both exposure status and disease status	In both cohort studies and case-control studies, higher participation rates are preferred. In cohort studies, higher follow-up rates are preferred. With lower participation (or follow-up rates), evidence for or against the potential for differential selection (e.g., greater participation of diseased among exposed compared with non-exposed) should be considered.
Potential confounding	A confounder is a common cause of both the exposure and the health outcome—thus, it is associated with both the exposure and the health outcome, but is not an intermediary between the two.	Studies with a design (e.g., matching procedures) or analysis (e.g., procedures for statistical adjustment) that adequately address the relevant sources of potential confounding for a given outcome are preferred.
Measurement of exposure	Standardized exposure assessment tools; validity and reliability	Studies are preferred that evaluate exposure during a biologically relevant time window for the outcome of interest, using higher quality exposure assessment methods that reduce measurement error. Measurement of exposure at the level of the individual is preferable to group-level exposures. Measurements of exposure should not be influenced by knowledge of health outcome status.
Measurement of health outcome	Standardized outcome assessment methods: validity and reliability	Studies that evaluate outcomes using generally accepted, standardized tools (e.g., disease classification systems, neuropsychological evaluation questionnaires) are preferred. Measurement or assignment of the outcome should not be influenced by knowledge of exposure status.
Power and precision	Numbers of test subjects and doses; experimental design	Preference is given to studies using designs reasonably expected to have power to detect responses of suitable magnitude. ³ This does not mean that studies with substantial responses but low power would be ignored, but they should be interpreted in light of a confidence interval or variance for the response.

Table F-13. Attributes used to evaluate studies for derivation of toxicity values.		
Aspects of Study	Data Characteristic	Considerations
<p>NOTES:</p> <p>1 USEPA (2002), A Review of the Reference Dose and Reference Concentration Processes, EPA/630/P-02/002F (page 4-11).</p> <p>2 Eliminating studies for which responses were not statistically significant will lead to bias toward larger effects. However, responses can be evaluated and weighted using standard errors or confidence intervals for the responses during hazard evaluation.³ A judgment about endpoint and study ‘sensitivity’ or protectiveness can be made <u>after</u> dose-response modeling⁵, in light of the range of candidate RfVs⁵ and their precision and quality.</p> <p>3 Power is an attribute of the design and population parameters; it cannot be inferred post-hoc using data from one experiment (Hoenig & Heisey, 2001, The American Statistician 55:19-24). Power is an ensemble property (based on a concept of repeatedly sampling a population) and is not a property of an individual study.</p>		

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

CONSIDERATIONS FOR COMBINING DATA FOR DOSE-RESPONSE MODELING

For most IRIS assessments, each POD has been derived based on data from a single study dataset. This is because in most cases, datasets are often **expected** to be heterogeneous for biological or study design reasons. Sources of potential heterogeneity include:

- Laboratory procedures used
- Population, species, and/or strain studied
- Sex
- Route of exposure

However, there are cases where one may consider conducting dose-response modeling after combining data from multiple studies, resulting in a single POD based on multiple datasets. For instance, this may be useful to increase precision in the POD or to quantify the impact of specific sources of heterogeneity.

Note that deriving toxicity values based on combining the results of multiple datasets **subsequent** to dose-response modeling of **each** dataset is discussed separately (see section 7.6 of *Preamble*, and associated draft *Handbook* text).

Examples of preliminary considerations as to whether are **potentially** suitable to derive a POD based on combining multiple datasets include the following:

- a. *Sufficient quality for deriving PODs (see section 6 of Preamble, and associated draft Handbook text).* Note that statistical precision should not be a quality consideration for this question, as it can be automatically accounted by statistical weighting. Indeed, one of the reasons for considering combining datasets may be to increase overall precision.
- b. *A common endpoint of concern reported.* Note that here “common endpoint” refers to the same specific outcome measurement, not just a common target site.
- c. *A common measure of dose available.* PBPK models may be useful for estimating a common (internal) dose measure, particularly across routes of exposure.

- 1 d. *Comparable durations, given the nature of the endpoint.* Note that this may include
2 exposure duration as well as observation duration (e.g., follow-up for cancer
3 epidemiology).
- 4 e. *Evidence for homogeneous responses to dose.* Species and sexes often differ in response
5 to dose, so convincing evidence would be needed to consider combining. A hypothesis
6 test of no difference would not be convincing unless it has high power to detect a
7 difference that matters (e.g., of the same magnitude as standard error of the mean).
- 8 f. *There is no one study that is clearly preferred.*

9 If potentially suitable datasets are available, then a statistician needs to be consulted to
10 evaluate in more detail whether the datasets are appropriate for combining, and if so, what
11 modeling approaches are appropriate to employ. Specific criteria for such evaluations will depend
12 on the design of the underlying studies and the sources of potential heterogeneity.

13 14 **CONDUCTING DOSE-RESPONSE MODELING**

15 ***[Note: EPA has guidance addressing this topic. The draft Handbook will eventually contain
16 more detailed information that summarizes the implementation of EPA's guidance in IRIS
17 assessments.]***

18 19 **DATA MANAGEMENT AND QUALITY CONTROL FOR DOSE-RESPONSE MODELING**

20 The IRIS Program has developed tools and approaches to manage data and ensure quality in dose-
21 response analyses. The objectives, described in more detail below, are to minimize errors, maintain
22 a transparent system for data management, automate tasks, where possible, and maintain an
23 archive of data and calculations used to develop assessments.

24 25 **A. Objectives**

26 **1. Minimize Errors**

27 Data (and metadata) should be entered into a database as early as possible in the
28 process, verified, and "locked" to prevent accidental changes. Verification should be
29 done either by double inspection or by double entry followed by machine comparison.
30 Data should be entered once, before use in evidence tables (which require
31 computations), and a subset of the same data will then be moved forward for calculating
32 PODs (using dose-response analyses or tabulation of LOAELs and NOAELs) and for use
33 in exposure-response arrays. Work after initial data entry and quality assurance (QA)
34 should not involve any cut and paste operations. No calculations will be made except
35 those recorded and retained transparently in the database. Later data entry or revision
36 will be subject to the same QA process. Initial QA and later changes will be recorded
37 and identified as to person responsible for data entry and data QA.

38 39 **2. Transparent From Source to Result**

40 Data entry is only one source of errors. Conversions and calculations, even simple ones
41 can introduce errors; tracing such errors can be time-consuming.

42
43 The objective is to have data entered as reported by the source and then verified.
44 Subsequent conversions and other calculations will be made transparently in a
45 database. For example, if the source reported inhalation concentrations in ppm and

1 exposures of 6 hours/day, 5 days/week, for 78 weeks, then for cancers, (a) an average
2 exposure would be calculated for a standard rodent lifetime as $\text{ppm} \times (6/24) \times (5/7) \times$
3 $(78/104)^3$, and (b) ppm would be converted to mg/m^3 using molecular weight and
4 other quantities in a “dosimetry tool”.

5
6 Also, the database should cite the source (reference) and the page(s) or table(s) or
7 figure(s) from which the endpoint data were extracted.

8
9 This approach will enable ready verification (and quick revision when necessary).

10 11 **3. Automate Tasks While Reducing Errors**

12 Data should be maintained in a modern database management system (dbms) designed
13 specifically for IRIS assessments to handle the types of endpoint data typically required.
14 The dbms allows users to prepare data for dose-response analysis (including choice of
15 BMRs), to execute analyses using Benchmark Dose Software (BMDS), and to marshal
16 and organize the results for review and model selection. The dbms allows users to
17 prepare custom reports and Microsoft (MS) Word tables and reports required for IRIS
18 assessments in the new streamlined formats. The dbms allows users to prepare
19 exposure-response arrays (figures) and import these into MS Word reports.

20
21 The dbms automates data processing, runs BMDS models, delivers results quickly, and
22 requires minimal human intervention (after initial setup and QA of modeling choices
23 and BMRs).

24 25 **4. Accessibility: Retain and Archive Working Files and Data**

26 Data in a dbms is easy to review and update. It is simple to re-run modeling for selected
27 endpoints or an entire set of endpoints. Metadata should identify sources and data
28 within sources unambiguously.

29
30 The dbms allows saving working files and data in a project “folder” when a project is
31 suspended or completed, making it relatively simple to renew work or revise previous
32 work. Project files can be shared among staff members working on a project (with
33 locking of files in current use). Completed project files will serve as an archive to
34 document work, including modeling decisions.

35 36 **B. Current Database Management System and Excel Tools**

37 The IRIS Program currently uses several tools for data management. These software tools
38 are still being refined to satisfy all of the objectives outlined above. These software tools
39 have been applied to several IRIS assessments, including dioxin, arsenic, and several other
40 chemicals.

41 42 **1. “BMDS Wizard”**

43 The BMDS Wizard is an MS Excel-based tool that was designed to facilitate benchmark
44 dose modeling when developing IRIS assessments. It handles one endpoint at a time in
45 each MS Excel workbook. It expedites setting up BMDS modeling and automates
46 running the various BMDS models. It includes forms for data and some other
47 information about the study.

48
49 The BMDS Wizard expedites setting up modeling choices. A user selects models from a
50 menu and the Wizard creates a table showing options for each of these models. One
51 type of model can be selected several times, each for a different BMR value. The user

1 then reviews the options and BMRs in the table before asking the Wizard to
2 automatically create the BMDS session file and BMDS model option files in a folder
3 specific to the data set. This expedites setting up a BMDS run of multiple models; the
4 table layout ensures better QA of modeling, option choices, and BMRs.
5

6 The Wizard then runs the chosen models and collects results in a single worksheet. It
7 also reports a number of warnings and flags that can be used to review models and
8 make a final selection of one model. Pop-ups reveal BMDS dose-response plots and
9 tables of estimates and residuals. The warnings and flags are based on an included
10 “logic” worksheet that can be modified by the user. The default logic worksheet
11 includes model selection criteria recommended in EPA’s Benchmark Dose Technical
12 Guidance (U.S. EPA, 2012). The results worksheet makes it easy to compare models side
13 by side and to document (in a comment column) any special reasons for rejecting and
14 accepting models and any unusual situations.
15

16 Wizard also allows the user to request MS Word tables and plots for selected models,
17 providing a well-organized summary of modeling results and model selection criteria.
18 MS Word templates are provided for this purpose, and these can be modified readily as
19 IRIS streamlined reporting requirements evolve.
20

21 2. “Dragon”

22 Dragon is a custom database management system (dbms) built in MS Access. Dragon
23 can be used for all data pertaining to an IRIS assessment. Dragon is still being improved
24 using feedback from IRIS users. Dragon works with several other software tools: BMDS
25 Wizard, Dosimetry Tool, and Exposure-Response Array software.
26

27 Dragon allows for: data entry; QA, review; data transformations; dosimetry calculations
28 (using the Dosimetry Tool); BMDS modeling for selected data (using the Wizard) and
29 collection of modeling results and model selection decisions; creating a variety of
30 reports and MS Word tables; and generating “skeleton” chapters and appendices, with
31 modeling results, for an IRIS assessment. It is also designed to make MS Word reports,
32 suitable as study summary tables and evidence tables. A summary of Dragon features
33 follows:
34

- 35 • Data are entered and viewed using forms.
- 36 • QA/QC of data entry is integrated into the tool.
- 37 • Data entry: Study quality, dose-response data, design and other metadata.
- 38 • Intermediate results: Dose conversions (dosimetry tool), BMD modeling (Wizard).
- 39 • Review of results: model and endpoint selection by user.
- 40 • Final results: Summary tables, Figures, MS Word reports (Tox Review and
41 Appendix).
- 42 • Customized for IRIS assessment requirements.
- 43 • Easy to set up model runs: creates BMDS session and option files for the model run.
- 44 • Organizes results for review: one model per line, all stats; flags problems.
- 45 • Identifies best model(s) using established criteria (user-modifiable).
- 46 • Writes IRIS assessment tables (Ch. 2 *Dose-Response Analysis* and modeling
47 appendices in supplemental information of the IRIS assessment), Exposure-
48 Response Arrays.
- 49 • Greatly reduces time to complete and report dose-response analyses.
- 50 • Will incorporate controlled nomenclature for endpoints.
- 51 • Flexible import and export capability, allowing data exchange with other software.

1
2 **DRAGON stores the following information:**

- 3 • Chemical-specific information—including name, molecular weight, etc.
- 4 • Study-specific information—including citation, HERO ID, study quality, etc.
- 5 • Dose-Protocol information—including species/strain/sex, route of exposure, dosing
6 protocol, etc.
- 7 • Dose information—including doses from PBPK modeling and dosimetric
8 conversions.
- 9 • Endpoint information—including NOAEL/LOAEL and statistical significance of each
10 dose-group.
- 11 • BMD information—including output from the BMDS Wizard.

12
13 **3. Dosimetry Tool**

14 A Dosimetry Tool was developed using MS Excel to overcome several challenges to
15 reliably making dosimetric conversions:

- 16 ■ Simple calculations, but numerous values to keep track of
 - 17 • Default values are in multiple guidance documents
 - 18 • Equations are specific to endpoint and study type
- 19 ■ Transparency and documentation
 - 20 • Calculation methods vary by author
 - 21 • Reporting format is not consistent
- 22 ■ Missing study data
 - 23 • Sources of body weight or food and water consumption are difficult to
24 locate.

25
26 **Dosimetry Tool Capabilities:**

- 27 ■ Makes dose conversions, consistently and transparently, that are easy to document.
28 The tool determines the correct formulas and defaults based on user inputs
- 29 ■ Used for Provisional Peer Reviewed Toxicity Values (PPRTVs)
 - 30 ■ Used by multiple study authors in ~100 PPRTV assessments
 - 31 ■ Sent as a compact deliverable showing all the inputs and results
- 32 ■ Consistency across multiple authors
 - 33 ■ Everyone uses the same equations, defaults, and reporting formats
- 34 ■ Easy QA of inputs
 - 35 ■ Summary tables show all input data and defaults
 - 36 ■ Equations show step-by-step calculations
- 37 ■ All conversions for an assessment can be saved in one workbook
 - 38 ■ Formatted tables stand alone as supporting documentation
- 39 ■ Will make default conversions (oral and inhalation)
- 40 ■ For more complex dosimetry calculations, converted doses can be entered into the
41 tool (for reproductive studies, etc.)
- 42 ■ Can use study-specific information on body weight, inhalation rate, food and water
43 consumption rates.

44
45
46 **EXTRAPOLATION TO LOWER DOSES AND RESPONSE LEVELS**

47 *[Note: EPA has guidance addressing this topic. The draft Handbook will eventually contain*
48 *more detailed information that summarizes the implementation of EPA's guidance in IRIS*
49 *assessments.]*

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

CONSIDERING SUSCEPTIBLE POPULATIONS AND LIFESTAGES

[Note: EPA has guidance addressing this topic. The draft Handbook will eventually contain more detailed information that summarizes the implementation of EPA’s guidance in IRIS assessments.]

DEVELOPING CANDIDATE TOXICITY VALUES

[Note: EPA has guidance addressing this topic. The draft Handbook will eventually contain more detailed information that summarizes the implementation of EPA’s guidance in IRIS assessments.]

CONSIDERATIONS FOR SELECTING ORGAN/SYSTEM-SPECIFIC OR OVERALL TOXICITY VALUES

The assessment derives or selects an organ/system-specific toxicity value for each organ or system affected by the agent. The assessment explains the rationale for each organ/system-specific toxicity value (for example, based on the highest quality studies, based on the most sensitive outcome, or based on a clustering of values). By providing these organ/system-specific toxicity values, IRIS assessments facilitate subsequent cumulative risk assessments that consider the combined effect of multiple agents acting at a common site or through common mechanisms (U.S. EPA, 2002).

Given multiple candidate toxicity values for a particular organ or system, each candidate value should be evaluated with respect to the multiple considerations:

- **Strength of evidence of hazard for the health effect or endpoint.** All other considerations being equal, effects and endpoints with stronger evidence of a causal relationship are preferred.
- **Attributes previously evaluated when selecting studies for deriving candidate toxicity values.** These include the study population/species, exposure paradigm, and quality of exposure and outcome measurement (see Section for Selecting Studies for Derivation of Toxicity Values). All other considerations being equal, studies of higher quality when evaluated according to these attributes are preferred.
- **Basis of the POD.** All other considerations being equal, a modeled benchmark dose (BMD) is preferred over a NOAEL, which is in turn preferred over a LOAEL. Additionally, when there is sufficient knowledge of toxicokinetics and the active toxic agent for the effect, a POD based on an internal dose metric would be preferred over one based on administered dose.
- **Other uncertainties in dose-response modeling.** These include the uncertainty in the BMD (e.g., reflected in the BMD/BMDL ratio) and uncertainty due to poor model fit.
- **Uncertainties due to other extrapolations.** All other considerations being equal, toxicity values for which other extrapolations are less uncertain are preferred. Note that the size of the composite uncertainty factor may **not** be a good indication of the remaining uncertainty, because some “uncertainty factors” overlap with aspects already addressed separately above (e.g., study population/species, use of a LOAEL as opposed to a NOAEL). Therefore,

1 to avoid double-counting, the remaining uncertainties that are discussed should be
2 explicitly enumerated.

3 Based on the results of this evaluation, the organ/system-specific toxicity value may be:

- 4 • Based on selecting a single candidate value considered to be most appropriate for
5 protecting against toxicity in the given organ or system.
- 6 • Based on deriving a “composite” value supported by multiple candidate toxicity values that
7 protects against toxicity in the given organ or system. One should carefully document how
8 the supporting candidate toxicity values are selected and how the composite value is
9 derived.

10 The assessment then selects an overall reference dose and an overall reference
11 concentration for the agent to represent lifetime human exposure levels where effects are not
12 anticipated to occur. This is generally the most sensitive organ/system-specific toxicity value,
13 though consideration of study quality and confidence in each value may lead to a different selection.
14

15 **CHARACTERIZING CONFIDENCE AND UNCERTAINTY IN THE TOXICITY VALUES**

16 *[Note: EPA has guidance addressing this topic. The draft Handbook will eventually contain*
17 *more detailed information that summarizes the implementation of EPA’s guidance in IRIS*
18 *assessments.]*

20 **SELECTING FINAL TOXICITY VALUES**

21 *[Note: EPA has guidance addressing this topic. The draft Handbook will eventually contain*
22 *more detailed information that summarizes the implementation of EPA’s guidance in IRIS*
23 *assessments.]*

24