



January 6, 2009
External Review Draft

U.S. Environmental Protection Agency

DRAFT

Expert Elicitation Task Force White Paper

**Prepared for the U.S. Environmental Protection Agency
by Members of the Expert Elicitation Task Force,
a Group Tasked by EPA's Science Policy Council**

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

NOTICE

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

DISCLAIMER

This document **is being** reviewed in accordance with U.S. Environmental Protection Agency peer review policy. This document, when finalized, will represent EPA's current thinking on this topic. It does not create or confer any legal rights for or on any person or operate to bind the public. The use of any mandatory language in this document is intended to describe laws of nature, scientific principles, or technical requirements and is not intended to impose any legally enforceable rights or obligations. . Mention of trade names or commercial products does not constitute endorsement or recommendation for use.

Expert Elicitation Task Force

Co-Chairs

Bob Hetes

ORD

Harvey Richmond

OAR

Chapter Leads

Chapter 1

Kathryn Gallagher, OSA

Chapter 3

Cynthia Stahl, Region 3

Lester Yuan, ORD

Chapter 6

Neil Stiber, OSA

Chapter 2

Bob Hetes, ORD

Chapter 4

Bob Hetes, ORD

Mark Corrales, OPEI

Chapter 7

Bob Hetes, ORD

Harvey Richmond, OAR

Chapter 5

Harvey Richmond, OAR

Task Force Members

Gary Bangs, OSA
David Chen, OCHP
Lisa Conner, OAR
Ila Cote, ORD
Jane Leggett, OAR
Joseph Greenblott, OCFO
Bryan Hubbell, OAR

Michael Messner, OW
Steve Nako, OPPTS
Barry Nussbaum, OEI
Marian Olsen, Region 2
Nicole Owens, OPEI
Zachary Pekar, OAR

Resha Putzrath, OSA
Zubair Saleem, OSWER
Brad Schultz, Region 5
Nathalie Simon, OPEI
David Simpson, OPEI
Holly Stallworth, SAB
Brad Venner, OECA

Science Policy Council Staff

Kathryn Gallagher & Neil Stiber

Elizabeth Lee Hofmann

Workgroup Leads for the Science Policy Council

Pai-Yei Whung
OSA

William Farland (retired) & Kevin Teichman
ORD

Rob Brenner
OAR

Table of Contents

1.0 INTRODUCTION	1
1.1 CONTEXT	1
1.2 PURPOSE OF THIS WHITE PAPER	2
1.3 ORGANIZATION OF THIS WHITE PAPER	2
2.0 BACKGROUND.....	3
2.1 WHEN DID EXPERT ELICITATION ORIGINATE?	5
2.2 WHY IS THERE INCREASED INTEREST IN EXPERT ELICITATION?	5
2.3 WHY IS EPA EXPLORING THE USE OF EXPERT ELICITATION?	6
2.4 IS EXPERT ELICITATION THE SAME AS EXPERT JUDGMENT?	9
2.5 WHAT IS THE EXPERIENCE WITH EXPERT ELICITATION AT EPA?	12
2.6 WHAT IS THE EXPERIENCE WITH EXPERT ELICITATION AT OTHER FEDERAL GOVERNMENT AGENCIES?.....	19
2.7 WHAT IS THE EXPERIENCE WITH EXPERT ELICITATION OUTSIDE THE U.S. FEDERAL GOVERNMENT?	20
2.8 WHAT EXAMPLES OF EXPERT ELICITATION ARE RELEVANT TO EPA?	21
2.9 SUMMARY	21
3.0 WHAT IS EXPERT ELICITATION?.....	22
3.1 WHAT IS EXPERT ELICITATION?	22
3.2 WHY IS EXPERT ELICITATION NECESSARY?	23
3.3 WHAT DO EXPERT ELICITATION RESULTS REPRESENT?	24
3.4 WHAT ARE SOME APPROACHES FOR EXPERT ELICITATION?	26
3.5 WHAT ARE GOOD PRACTICES FROM EXPERT PRACTITIONERS FOR ELICITING EXPERT JUDGMENT?	31
3.6 SUMMARY	42
4.0 WHAT THOUGHTS ABOUT APPLICABILITY AND UTILITY SHOULD INFORM THE USE OF EXPERT ELICITATION?	44
4.1 HOW IMPORTANT IS IT TO CONSIDER UNCERTAINTY?.....	45
4.2 WHAT IS THE NATURE OF THE UNCERTAINTIES TO BE ADDRESSED?	50
4.3 WHAT ARE OTHER METHODS TO CHARACTERIZE UNCERTAINTY?	51
4.4 WHAT ROLE MAY CONTEXT PLAY FOR AN EE?.....	57
4.5 WHAT RESOURCES ARE REQUIRED FOR AN EXPERT ELICITATION?	61
4.6 SUMMARY	64
5.0 HOW IS AN EXPERT ELICITATION CONDUCTED?.....	65
5.1 WHAT ARE THE STEPS IN AN EXPERT ELICITATION?.....	65
5.2 WHAT ARE THE PRE-ELICITATION ACTIVITIES?	67
5.3 WHAT APPROACHES ARE USED TO CONDUCT EXPERT ELICITATIONS?.....	75
5.4 WHAT POST-ELICITATION ACTIVITIES SHOULD BE PERFORMED?	79
5.5 WHEN AND WHAT TYPE OF PEER REVIEW IS NEEDED FOR REVIEW OF EXPERT ELICITATION?	85

5.6.	SUMMARY	85
6.0	HOW SHOULD RESULTS BE PRESENTED AND USED?.....	86
6.1	DOES THE PRESENTATION OF RESULTS MATTER?	86
6.2	WHAT IS THE STAKEHOLDER AND PARTNER COMMUNICATION PROCESS?	86
6.3	HOW CAN COMMUNICATIONS BE STAKEHOLDER-SPECIFIC?	87
6.4	WHAT IS IN A TECHNICAL SUPPORT DOCUMENT?	88
6.5	WHAT ARE EXAMPLES OF EFFECTIVE EXPERT ELICITATION COMMUNICATIONS?.....	93
6.6	HOW CAN EES BE TRANSPARENT, DEFENSIBLE, AND REPRODUCIBLE?	105
6.7	SHOULD EXPERT JUDGMENTS BE AGGREGATED FOR POLICY DECISIONS?	106
6.8	HOW CAN EXPERT ELICITATION RESULTS AND OTHER PROBABILITY DISTRIBUTIONS BE INTEGRATED? ..	107
6.9	HOW CAN AN EXPERT ELICITATION BE EVALUATED POST HOC?	107
6.10	SUMMARY	108
7.0	FINDINGS AND RECOMMENDATIONS	109
7.1	FINDINGS.....	109
7.2	RECOMMENDATIONS	111
	REFERENCES	117
	APPENDIX: FACTORS TO CONSIDER WHEN MAKING PROBABILITY JUDGMENTS	132

1.0 INTRODUCTION

1.1 CONTEXT

Expert elicitation (EE) is a systematic process of formalizing and quantifying, typically in probabilistic terms, expert judgments about uncertain quantities. The expert elicitation process may involve integrating empirical data with scientific judgment, and identifying a range of possible outcomes and likelihoods. An important part of the EE process includes documentation of the underlying thought processes of the experts. Expert elicitation is a multi-disciplinary process that can inform decision-making by characterizing uncertainty and filling data gaps where traditional scientific research is not feasible or data are not yet available. If performed using appropriate methods and quality standards, including peer review and transparency, EE can be a reliable component of sound science.

Expert elicitation has been used by federal agencies, the private sector, academia, and other groups. For example, in the 1980s EPA's Office of Air Quality, Planning and Standards (OAQPS) used EE to assess exposure-response relationships for lead and ozone. More recently, OAQPS used EE to analyze uncertainty in the relationship between exposures to fine particles and the annual incidence of mortality. The Department of Energy used EE to evaluate nuclear waste and other related issues. Other uses of EE by the government and academia include cost-benefit analysis, risks associated with climate change, technology development, and food safety.

1.1.1 Role of the Science Policy Council

In April of 2005 the Science Policy Council (SPC) formed the Expert Elicitation Task Force (hereafter cited as "Task Force") and charged it to initiate a dialogue within the Agency about EE and to facilitate future development and appropriate use of EE methods. The SPC was established in 1993 by EPA's Administrator to have primary responsibility for addressing and resolving cross-program, cross-media, and interdisciplinary science policy issues. The SPC is composed of senior managers from across the Agency and is chaired by EPA's Science Advisor.

Following the SPC's call for Task Force members, program and regional offices nominated staff with appropriate and relevant expertise. As a first step toward achieving its charge, the Task Force developed this White Paper that provides a framework for determining the appropriate conduct and use of EE.

This White Paper was peer reviewed internally by members of the SPC Steering Committee and by additional EPA staff. The SPC Steering Committee is composed of scientific and policy experts from EPA's program and regional offices and serves as the SPC's principal advisory group. The SPC has the authority to provide final review for this White Paper and to

approve its release for external peer review. Ultimately, the SPC can provide final approval of the document for dissemination.

1.2 PURPOSE OF THIS WHITE PAPER

The Task Force's purpose is to initiate a dialogue within the Agency about the conduct and use of EE and then to facilitate future development and appropriate use of EE methods. To that end, the Task Force facilitated a series of discussions to familiarize Agency staff with EE and to evaluate and address issues that may arise from its use. This White Paper reflects those discussions and presents issues that are pertinent to EE, including: What is EE?, When to consider using EE?, How is an EE conducted?, and How should results be presented and used?

Because input from a range of internal and external of stakeholders was not formally solicited, this White Paper does not present official EPA guidelines or policy. This White Paper may be used to facilitate the development of any future EE guidance or policy.

1.3 ORGANIZATION OF THIS WHITE PAPER

This White Paper reflects discussions about EE that were coordinated by the Task Force. Chapter 2 provides background for EE at EPA and summarizes the context of increasing interest in this approach. In addition, it reviews experiences with EE at the EPA, throughout the federal government, and with international groups. It also shares applications that are relevant to EPA issues. Chapter 3 provides the definition of EE for this White Paper and considers its advantages and disadvantages. Chapter 4 recognizes that EE is one of many tools to characterize uncertainty and examines what factors may help determine when EE is appropriate. Chapter 5 summarizes what is needed to conduct a credible and acceptable EE. Chapter 6 offers considerations for presenting and using EE results in EPA decisions. Finally, Chapter 7 presents some significant issues regarding EE. Where consensus was reached by the Task Force, the White Paper provides recommendations for further development and use of EE within EPA or by parties submitting EE assessments to EPA for consideration. The Task Force also identifies issues that may require further deliberation as part of the development of an EPA guidance or policy on EE. Recommendations are also provided for potential EPA actions, including training, that could promote the use of EE.

2.0 BACKGROUND

The EPA frequently makes decisions on complex environmental issues that require analyses from a broad range of disciplines. Among the many sources of uncertainty and variability in these analyses are estimates of parameter values and choices of models. Furthermore, in some cases, critical data may be unavailable or inadequate. Although the presence of uncertainty complicates environmental decision making, EPA must still make timely decisions. Expert elicitation is one method for characterizing uncertainty and providing estimates in data-poor situations. Thus, the influence of uncertainty must be appropriately considered and addressed in all decisions. Because EPA has recognized that uncertainty is an important aspect of risk assessment, it has developed guidance, including its *Risk Characterization Handbook* (USEPA, 2000a). Even though this handbook was developed with risk assessment in mind, it is applicable and useful for many kinds of EPA assessments, including EE.

In 1983, the National Academy of Sciences (NAS) published *Risk Assessment in the Federal Government: Managing the Process* (NAS, 1983; commonly referred to as the “Red Book”) which formalized the risk assessment process. EPA integrated the “Red Book” principles of risk assessment into its practices and, the following year, published *Risk Assessment and Management: Framework for Decision Making* (USEPA, 1984), which emphasizes making the risk assessment process transparent, fully describing the assessment’s strengths and weaknesses, and addressing plausible alternatives. Then, starting in 1986, EPA began issuing a series of guidelines for conducting risk assessments (i.e., exposure, carcinogen, chemical mixtures, mutagenicity, and suspect developmental toxicants, USEPA, 1986). Although EPA’s initial efforts focused on human health risk assessment, in the 1990s the basic approach was adapted to ecological risk assessment to address a broad array of environmental risk assessments in which human health impacts are not a direct issue. EPA continues to make a substantial investment in advancing the science and application of risk assessment through updates to these guidelines and the development of additional guidelines, as needed.

Over the next several years, the NAS expanded on its “Red Book” risk assessment principles in a series of subsequent reports, including *Pesticides in the Diets of Infants and Children* (NAS, 1993), *Science and Judgment in Risk Assessment* (NAS, 1994; commonly referred to as the “Blue Book”), and *Understanding Risk: Informing Decisions in a Democratic Society* (NAS, 1996). The purpose of the risk assessment process, as characterized by the NAS, is to ensure that assessments meet their intended objectives and are understandable. Over time, EPA risk assessment practices advanced along with NAS’s progression of thought.

In 1992, EPA provided the first risk characterization guidance to highlight the two necessary elements for full risk characterization: (1) address qualitative and quantitative features of the assessment and (2) identify any important uncertainties and their influence as part of a discussion on confidence in the assessments. Three years later, EPA updated and issued the current Agency-wide *Risk Characterization Policy* (USEPA, 1995a). To ensure that the risk assessment process is transparent, this *Policy* requires risk characterization for all EPA risk assessments. In addition, this *Policy* emphasizes that risk assessments be clear, reasonable, and consistent with other risk assessments of similar scope across the Agency. Effective risk characterization requires transparency in the risk assessment process and clarity, consistency, and reasonableness of the risk assessment product. EPA's *Risk Characterization Handbook* (USEPA, 2000a) was developed to implement the *Risk Characterization Policy*. The importance of characterizing uncertainty was re-affirmed in the recent EPA staff paper *An Examination of EPA Risk Assessment Principles and Practices* (EPA, 2004a). This staff paper identified the use of probabilistic analyses as an area in need of major improvement.

Risk assessments are often used as the basis for calculating the benefits associated with Agency regulations. Such benefits-costs analyses can be important tools for decision makers, where statutorily permitted, both in the context of regulatory reviews required under Executive Order (EO) 12866 and Section 812 of the Clean Air Act Amendments of 1990 which require EPA to assess the costs and benefits of the Clean Air Act. In its 2002 report entitled *Estimating the Public Health Benefits of Proposed Air Pollution Regulations*, the NAS emphasized the importance of fully characterizing uncertainty for decision makers and encouraged EPA to use EE in the context of expressing uncertainty associated with estimated benefits.

Guidance for conducting regulatory analyses required under EO 12866 is provided in the Office of Management and Budget's (OMB) Circular A-4 (USOMB, 2003). This guidance emphasizes that the important uncertainties connected with regulatory decisions need to be analyzed and presented as part of the overall regulatory analysis. Whenever possible, appropriate statistical techniques should be used to determine the probability distribution of relevant outcomes. For major rules involving annual economic effects of \$1 billion or more, a formal quantitative analysis of uncertainty is required. The OMB guidelines outline analytical approaches, of varying levels of complexity, which could be used for uncertainty analysis such as qualitative disclosure, numerical sensitivity analysis, and formal probabilistic analysis (required for rules with impacts greater than \$1 billion). EE is one of the approaches specifically cited in these guidelines for generating quantitative estimates (e.g., cost-benefit analysis) when specific data are unavailable or inadequate.

2.1 WHEN DID EXPERT ELICITATION ORIGINATE?

The origins of EE can be traced to the advent of decisions theory and decision analysis in the early 1950s. In 1954, Savage established the “probabilities of orderly opinions” which states that the choice behavior of a rational individual can be represented as an expected utility with a unique probability and utility measure. EE’s development also drew on the operational definitions of probability that arose out of the semantic analysis discussions of Mach, Hertz, Einstein, and Bohr (Cooke, 1991). Since the early 1970s, decision analysts in the private and public sectors have used formal EE processes to obtain expert judgments for their assessments. Section 2 presents a variety of examples from EPA, other Federal agencies, and beyond.

2.2 WHY IS THERE INCREASED INTEREST IN EXPERT ELICITATION?

There are numerous quantitative methods for characterizing uncertainty. The available types of methods are described briefly Section 4.2. While there is no consensus on a preferred method to characterize uncertainty, there is general agreement that practitioners should describe uncertainty to the extent possible with available data and well-established physical and statistical theory. However, limitations in data and/or understanding (i.e., lack of a theory relevant to the problem at hand) may preclude the use of conventional statistical approaches to produce probabilistic estimates of some parameters. In such cases, one option is to ask experts for their best professional judgment (Morgan and Henrion, 1990). Expert elicitation (which is defined and discussed in greater detail in Chapter 3) is a formal process by which expert judgment is obtained to quantify or probabilistically encode uncertainty about some uncertain quantity, relationship, parameter, or event of decision relevance.

Expert elicitation is recognized as a powerful and legitimate quantitative method for characterization of uncertainty and for providing probabilistic distributions to fill data gaps where additional research is not feasible. The academic and research community, as well as numerous review bodies, have recognized the limitation of empirical data for characterization of uncertainty and have acknowledged the potential for using EE for this purpose. In *Science and Judgment in Risk Assessment* (NAS, 1994) the NAS recognized that for “parameter uncertainty, enough objective probability data are available in some cases to permit estimation of the probability distribution. In other cases, subjective probabilities might be needed.” In this “Blue Book” report, the NAS further recognized the “difficulties of using subjective probabilities in regulation” and identified perceived bias as one major impediment; but, noted that “in most problems real or perceived bias pervades EPA’s current point-estimate approach.” In addition, the NAS stated that “there can be no rule that objective probability estimates are always preferred to subjective estimates, or vice versa.”

The utility of EE has been discussed by NAS, OMB, and EPA. In the following examples, they provide advice for the appropriate and beneficial use of EE. With respect to benefits analyses, NAS (2002) recommends,

“EPA should begin to move the assessment of uncertainties from its ancillary analyses into its primary analyses by conducting probabilistic, multiple-source uncertainty analyses. This shift will require specifications of probability distributions for major sources of uncertainty. These distributions should be based on available data and *expert judgment*.”

In its Circular A-4(USOMB, 2003), OMB suggests using EE to address requirements for probabilistic uncertainty analysis:

In formal probabilistic assessments, *expert solicitation*¹ is a useful way to fill key gaps in your ability to assess uncertainty. In general, experts can be used to quantify the probability distributions of key parameters and relationships. These solicitations, combined with other sources of data, can be combined in Monte Carlo simulations to derive a probability distribution of benefits and costs.

In addition, the EPA Cancer Risk Assessment Guidelines (USEPA, 2005) provide for the use of EE in such assessments.

In many of these scientific and engineering disciplines, researchers have used rigorous expert elicitation methods to overcome the lack of peer-reviewed methods and data. Although expert elicitation has not been widely used in environmental risk assessment, several studies have applied this methodology as a tool for understanding quantitative risk. ... These cancer guidelines are flexible enough to accommodate the use of expert elicitation to characterize cancer risks, as a complement to the methods presented in the cancer guidelines. According to NAS (NAS, 2002), the rigorous use of expert elicitation for the analyses of risks is considered to be quality science.

2.3 WHY IS EPA EXPLORING THE USE OF EXPERT ELICITATION?

EPA recognizes the value of EE as a powerful tool to improve the characterization of uncertainty in risk and other types of assessments. EPA’s experience with EE (described in Section 2.5) highlights its benefits of enhancing the scientific and technical credibility of EPA assessments, and acceptability of these assessments within the scientific and technical community (NAS, 2002). However, concerns have been raised about using EE within the context of EPA decision making. These include transparency in the use of empirically versus

¹ OMB used the phrase expert solicitation rather than expert elicitation but text and references are similar to that associated with expert elicitation. For the purposes of this White paper, EPA assumes they are equivalent.

judgment derived estimates, potential for delays in rulemaking while EE is conducted, and the lack of EPA guidelines on and limited experience with EE. The American Bar Association (ABA, 2003), in its comments on the requirement for formal probabilistic analyses in *OMB's Draft 2003 Report to Congress on the Costs and Benefits of Federal Regulation*, stated

... formal probabilistic analysis will be impossible to meet rigorously in cases where the underlying science is so uncertain as to preclude well-founded estimates of the underlying probability distribution.... In such situations, the effort to generate probability distributions in the face of fundamental uncertainty through guesses derived from so-called 'expert elicitation' or 'Delphi' methods runs the risk of creating that 'false sense of precision' which OMB elsewhere cautions agencies to avoid. Accordingly, we believe such methods should be used sparingly, and we strongly endorse the recent recommendation of the National Research Council that agencies disclose all cases in which expert elicitation methods have been used.

EPA's limited experience conducting EEs has primarily been within the Office of Air and Radiation; mainly focused on improving risk and benefits assessments but not directly used to support a regulatory decision. EPA has no clear guidelines to assist in the conduct and use of such techniques for regulatory analyses or other purposes. Given these limitations, EPA would benefit from a thoughtful discussion about the conduct and use of EE to support regulatory and non-regulatory analyses and decision making. The desire for such discussion is likely to grow because EE is increasingly being identified as a method to meet various requirements for characterizing and addressing uncertainty. Any early efforts have the potential to become precedents for how EE analyses are conducted; and, more importantly, how these analyses are used to support decisions. To minimize the chance that these precedents may carry unintended detrimental consequences, early EE efforts should include a dialogue on their regulatory, legal, and statutory implications as well as their technical aspects.

This Task Force has initiated a dialogue within the Agency about these methods and plans to facilitate the development of EE guidance for EPA. In this White Paper, the Task Force considers a broad range of technical, statutory, regulatory, and policy issues including:

- When is EE an appropriate (well-suited) methodology to characterize uncertainty?
 - What are good practices based on a review of the literature and actual experience within EPA and other federal agencies in conducting an EE, considering the design objectives and intended use of the results (e.g., prioritizing research needs, input to risk assessment, input to regulatory impact analysis)?
-

- When, and under what circumstances, is it appropriate to aggregate/combine expert judgments and how should such aggregation/combination be done?
- When in the EE process is peer review beneficial?
- What type of peer review is needed to review EE methods and their use in specific regulatory actions?
- What are the implications of EPA's Quality System and Information Quality Guidelines on EE?

2.3.1 How Does this EE Activity Relate to Efforts to Develop And Promote Probabilistic Risk Assessment (PRA) at EPA?

As highlighted by a major recommendation in the EPA staff paper on *Risk Assessment Principles and Practices* (USEPA, 2004a), EPA recognizes the need for more appropriate and timely use of probabilistic assessments. As a result, EPA has several current efforts to promote the appropriate use of probabilistic analyses in support of regulatory analyses, including activities sponsored by the SPC and the Risk Assessment Forum (RAF). As described below, these major activities include a RAF Colloquium on Probabilistic Risk Assessment, its associated follow-up workgroup activities, and EPA's co-sponsorship with the Society of Toxicology (SOT) of a Workshop on Probabilistic Risk Assessment.

In April 2004, EPA's RAF held a Colloquium on Probabilistic Risk Assessment to address the following topics: identifying probabilistic techniques that can better describe variability and uncertainty, communicating probabilistic methods for the purpose of risk management and risk communication, supplementing the *Guiding Principles for Monte Carlo Analysis* (USEPA, 1997), and deciding the next steps for advancing probabilistic methods in risk assessments or improving their implementation. This Colloquium, attended by risk assessors from across the Agency and several invited external experts, was divided into three half-day sessions on human exposure, ecological risk, and human health effects. Each session included several presentations followed by a panel discussion that was open to all colloquium participants.

As a follow-up to this Colloquium the RAF formed a workgroup to address how to improve support for EPA decision making through the use of probabilistic methods. One of this workgroup's first activities addresses a major and multi-disciplinary recommendation from the 2004 RAF Colloquium: the need for a dialogue between risk assessors and decision makers (risk managers). These discussions are essential to identify specific issues of concern and determine

needs for promoting the useful application of these methods. Without these up-front understandings, probabilistic methods might be applied at high resource cost to EPA; but, provide information that is irrelevant with has little or no impact on decisions. As a priority, this workgroup also seeks to promote the exchange of knowledge between risk assessors and risk managers across program and regional offices. Such efforts may include a series of workshops or seminars along with a clearinghouse of models, methods, and experiences.

As a follow-up to the staff paper on *Risk Assessment Principles and Practices* (USEPA, 2004a), EPA co-sponsored, with the SOT and several other organizations, a workshop on Probabilistic Risk Assessment as part of the SOT's Contemporary Concepts in Toxicology workshop series. The workshop provided an opportunity for in-depth discussion of four critical topic areas: (1) exposure assessment, (2) ecological risk assessment, (3) human health risk assessment and medical decision analysis, and (4) decision analysis/multi-criteria decision analysis cost-benefit analysis. Draft white papers for each topic area, that were prepared for and discussed at the workshop, have been submitted for journal publication. Expert elicitation was discussed as a critical method for advancing probabilistic analyses in the Agency.

It should also be noted that individual program and Regional offices have, in the past and may in the future, developed their own guidance on the conduct and use of probabilistic methods. For example, EPA's Office of Solid Waste and Emergency Response (OSWER) developed guidance for this purpose – *Risk Assessment Guidance for Superfund: Volume 3 – (Part A, Process for Conducting Probabilistic Risk Assessment)*(EPA, 2001a).

The EE Task Force intends to complement the above efforts. This White Paper should serve as an initial framework for Agency discussions and development of policy or guidance on EE.

2.4 IS EXPERT ELICITATION THE SAME AS EXPERT JUDGMENT?

EPA often needs to make decisions that address complex problems many of which lack direct empirical evidence. To obtain insights, this kind of decision making requires judgment to assess the impact and significance of existing data or theory. As a result, judgment is an inherent and unavoidable part of most EPA assessments and decisions. In fact, judgment is also inherent in empirical data and highly targeted technical activities. While some try to portray traditional statistics as an objective activity, in truth subjectivity or expert judgment often plays a large role. For example, the analyst's expert judgment is critical and unavoidable when developing and selecting a study design, a sampling strategy, specific statistical tests, goodness of fit measures, and rules for excluding outliers. EPA relies on various forms of expert judgment throughout the scoping, design, and implementation of our assessments. EPA's *Risk Characterization*

Handbook (USEPA, 2000a) recognizes the role and importance of judgment and establishes standards for describing transparently when and how it is used in risk assessments and decisions.

As illustrated in Figures 2-1 and 2-2, decision making at EPA concerns complex problems that are influenced by multiple factors. Even within a particular discipline or activity, inherent expert judgment may include values or preferences (e.g., the default to promote a health-protective conservative estimate of risk). In addition to risk assessment results, social, economic, political, statutory/legal, public health and technological factors may influence a final decision. Therefore scientific (state of knowledge) and value and preferences are both included in any particular assessment.

Although values and preferences play a major role in how to weight or balance information across various disciplines, expert judgment can help to integrate this information. However, tools that rely on formal integration of values and preferences are not within the Task Force’s definition of EE.

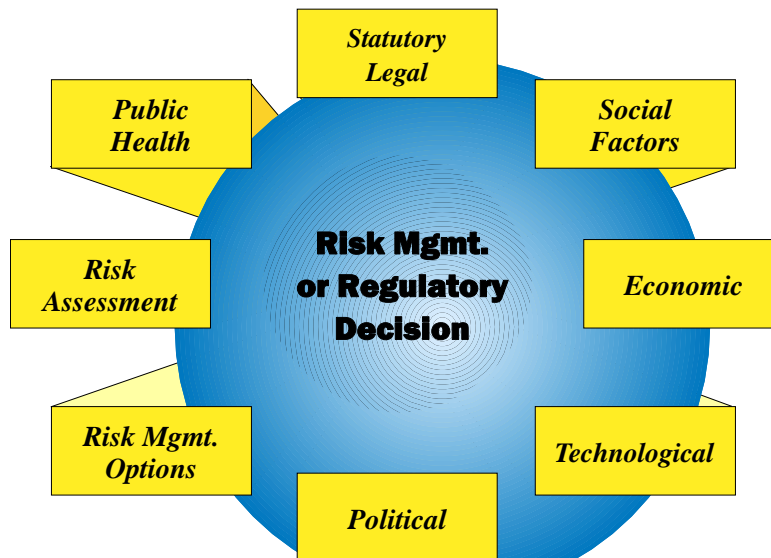


Figure 2-1. Factors that Influence Risk Management Planning at EPA

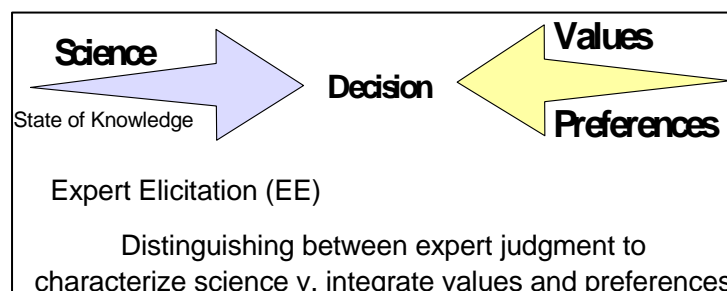


Figure 2-2. The Role of Science and Values in Environmental Decision Making

For the purposes of this White Paper the Task Force has chosen to define EE more narrowly than Expert Judgment as a method limited to characterizing the science (state of knowledge), rather than values and preferences. This definition is consistent with the field of EE and practitioners in the field generally advocate a separation of the analyses of uncertainty associated with quantitative physical parameters or empirical models from the phase involving preferences and social value judgments leading up to decisions. Experts are best suited for answering science questions, while appointed or elected decision makers (with input from risk managers and appropriate stakeholders) are best suited for providing values and preferences.

EPA also relies on various forms of expert input in formulating its judgments through external review and evaluation of the quality of its assessments. These include internal and external peer review and peer involvement such as through the EPA Science Advisory Board (SAB). In these cases, reviewers typically bring technical expertise as well as expertise. Such experts rely on peer-reviewed published literature, as well as their own knowledge (e.g., unpublished studies) and perception. The distinguishing feature of EE relative to these other forms of expert input is the use of a systematic process to formalize judgments. Expert input via these other mechanisms provides expert opinions; and while some may take procedural steps to minimize the effects of heuristics or other biases in expert judgments, they generally do not apply formal scientific protocol to address these effects.

The Task Force recognizes that a wide range of activities may fall under the term expert judgment. These activities range from the very informal (e.g., an analyst choice on a parameter value) to very formal and rigorous methods such as systematic characterization of an expert's judgment. The Task Force has chosen to restrict the usage of EE to this more formal systematic method to differentiate it from other forms of expert judgment and expert input. Furthermore, the Task Force has chosen to focus on questions related to science (state of knowledge) as the focus of EE and does not include within its definition the elicitation of personal or social values and preferences.

In considering the range of approaches to include in EE, the EE Task Force relied on the volume of decision analysis literature and history that are derived from or consistent with

Bayesian theory.² For an effort to be considered an EE, at a minimum, all of the following elements³ (as described in detail later in this White Paper) must be present:

- Problem definition -- unambiguous which meets Clairvoyance Test⁴
- Formal protocol -- required to ensure consistency in elicitation and to control for heuristics and biases
- Identification, summary, and sharing of the relevant body of evidence with experts
- Formal elicitation -- encoding of judgments as probabilistic values or distributions -- typically via interaction with objective independent party
- Output -- judgment or degree of belief is expressed quantitatively (typically probabilistically)

The Task Force's selection of these factors was influenced by the available literature and expertise on the Task Force. The Task Force did not fully consider other non-Bayesian, non-quantitative, semi-quantitative, or non-probabilistic social science encoding methods that also control for heuristics and biases (e.g., Grounded Theory, Nominal Group Technique, and variants of the Delphi method). Nevertheless, such methodologies may be appropriate for formal elicitation of expert judgments where quantitative (i.e., probabilistic) characterization of judgment is not the desired output if they fulfill the analytical requirements of the charge and are implemented with sufficient rigor and review.

2.5 WHAT IS THE EXPERIENCE WITH EXPERT ELICITATION AT EPA?

The majority of EPA's experience with EE is in its Office of Air and Radiation (OAR), Office of Air Quality Planning and Standards (OAQPS). The OAQPS first explored the use of EE in the late 1970s. In fact, NAS (2002) notes that "OAQPS has been a pioneer in the application of (use of expert judgment) approaches to estimating the health risks due to exposure to air pollutants (Richmond, 1981; Feagans and Biller, 1981; Whitfield et al. 1991; Rosenbaum et al. 1995)." As summarized below, EPA's experience includes OAQPS studies to evaluate the

² Bayesian theory, discussed in greater detail in Chapter 3 as it relates to EE, proposes that a person's belief in a proposition can be described according to probability theory.

³ As described in later chapters there are other elements which are described as components of EE. Many of the elements are not unique to EE but may be associated in some form with others modes of expert input. The ones listed here are suggested as the Task Force's minimum distinguishing operational features for EE. This does not imply that an exercise containing these minimum elements would represent a good EE. How well these elements are carried will determine the overall quality of any particular effort.

⁴ Such that an omniscient being with complete knowledge of the past, present, and future could definitively answer the question.

health effects of various criteria pollutants and other EPA offices efforts to forecast sea level rise, store radioactive waste, and support ecological model development.

2.5.1 Criteria Pollutants

2.4.1.1 1977-1978 Ozone NAAQS Review

Motivated by the statutory requirement to protect public health with an adequate margin of safety, OAQPS pursued development of probabilistic estimates. Drawing on techniques developed as part of the field of decision analysis (probability encoding), OAQPS derived probabilistic concentration-response relationships from experts for several health endpoints as part of the 1977-1978 ozone NAAQS review (Feagans and Biller, 1981). These early efforts were viewed as controversial because there was no formal protocol, no pre- or post-elicitation workshops, and little experience conducting elicitation. To review this OAQPS approach and EPA's efforts to develop probabilistic risk assessment methods, the SAB Executive Committee created the SAB Subcommittee on Health Risk Assessment in 1979.

This SAB Subcommittee held several meetings from 1979 to 1981 to review reports prepared by six teams of nationally recognized experts, additional developmental efforts by OAQPS, a literature review of probability encoding (Wallsten and Budescu, 1983), and two illustrative applications involving EE's for health effects associated with carbon monoxide. In spring of 1981, the SAB encouraged EPA to take the best elements of the illustrative applications and the original OAQPS proposed approach in carrying out an EE addressing health effects of lead (Pb) for EPA's Pb NAAQS risk assessment.

2.5.1.2 Lead NAAQS Risk Assessment (1983-1986)

Following the advice of the SAB, OAQPS sponsored a full EE on the health effects of lead to support the Pb NAAQS review. A formal protocol was developed and pilot-tested. The elicitation focused on two health endpoints (IQ decrement, hemoglobin decrement). The study and its results (Figure 2-3) received a favorable review from the Clean Air Scientific Advisory Committee Lead (Pb) Subcommittee (CASAC). The elicitation was critical in deriving and characterizing the uncertainty in policy-relevant concentration response functions beyond those available in the empirical literature. Although the SAB Health Risk Assessment Subcommittee was dissolved following its review of the elicitation project, three of its members were added to CASAC Pb Subcommittee.

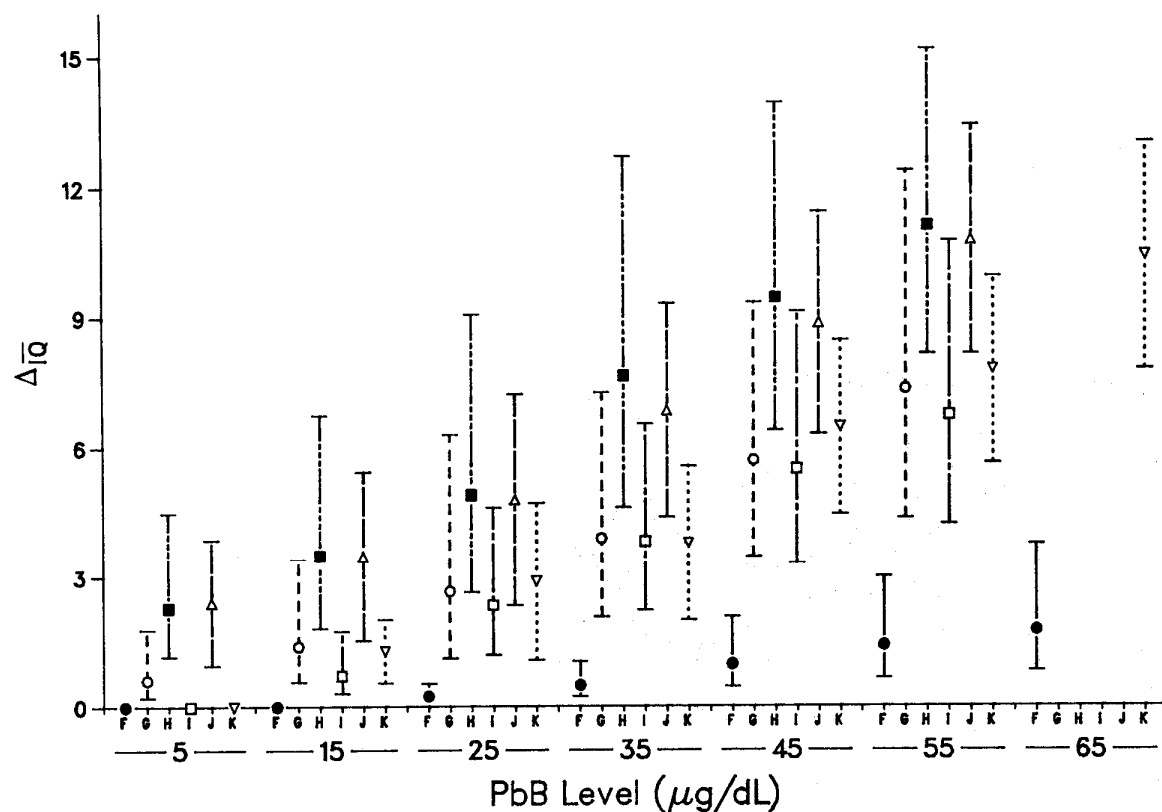


Figure 2-3. Median and 90% Credible Interval Judgments of All Experts about Lead-Induced IQ Decrements among Low-SES Children Aged 7 Years (Whitfield and Wallsten, 1989).

2.5.1.3 Ozone Chronic Lung Injury Assessment

Drawing on its experience with the Pb NAAQS assessment, OAQPS pursued an additional study to develop concentration-response functions for chronic lung injury associated with long-term exposures to ambient ozone concentrations (Winkler et al. 1995). The specific objective was to characterize scientific judgments regarding the risk of chronic lung injury to children aged 8 to 16 years and adult outdoor workers due to long-term ozone exposure in areas with exposure patterns similar to Southern California and the Northeast. Again, a formal protocol was developed and pilot tested prior to the elicitation exercise. Experts were provided with air quality information, exposure model estimates, and dosimetry model estimates. The measure of injury was the incidence of mild or moderate lesions in the centriacinar region of the lung. Probabilities of population response rates were elicited. After a post-elicitation workshop to encourage information exchange among experts, a second round of encoding was conducted (Figure 2-4).

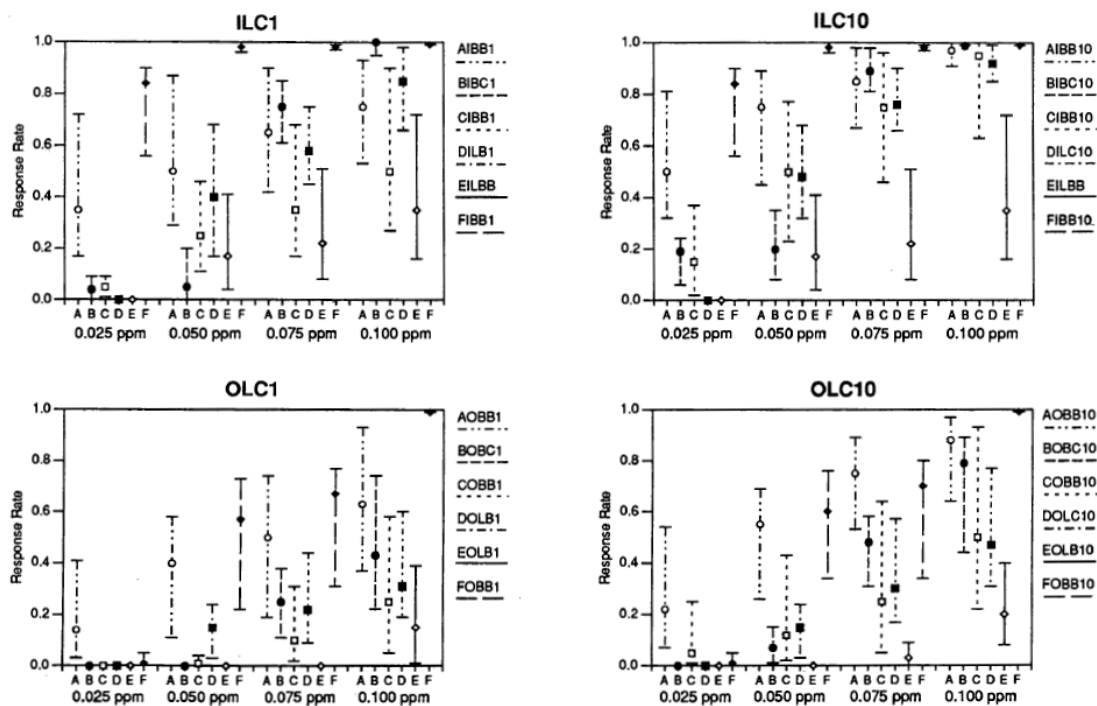


Figure 2-4. Medians and 90% Credible Intervals for Judgments about Mild and Moderate Lesions, Children, Los Angeles, and 1 and 10 Ozone Seasons (Argonne National Labs, 1991)

2.5.1.4 Lessons Learned from Early OAQPS Efforts

The use of EE in risk assessment is not necessarily straight-forward. When EE was fairly new, much of the scientific community was skeptical. Early efforts were controversial and were questioned, in part, due to the lack of experience and formal procedures. Through a lengthy collaborative effort with the SAB, OAQPS was able to improve the quality of the assessments and increase their credibility within the scientific community. While OAQPS EEs have gained acceptance, it is likely that similar collaborative efforts will be needed for EE method development and application sponsored by other EPA offices.

The results of these initial elicitations were used to inform policy makers and the public of possible health implications due to short-term and long-term exposures to specific criteria pollutants. However, these results were not used as the basis for any EPA regulatory decisions.⁵

⁵ No formal regulatory decision was made following the completion of the Pb risk assessment for reasons other than EE, and the EE was used to develop the chronic ozone lung injury assessment but was not intended to support the ozone NAAQS review decision.

2.5.1.5 PM Concentration-Response (C-R) for Mortality

As mentioned in Section 2.2, the NAS (2002) recommended that EPA improve its characterization of uncertainty in its benefits analyses by using both available data and expert judgment. They recommended that EPA build on the prior OAQPS experiences in the use of formally elicited expert judgments, but noted that a number of issues must be addressed. The NAS stressed that EPA should distinguish clearly between data-derived components of an uncertainty assessment and those based on expert opinions. As a first step in addressing these NAS recommendations regarding EE, EPA, in collaboration with OMB, conducted a pilot EE to characterize uncertainties in the relationship between ambient fine particles (PM_{2.5}) and premature mortality. This pilot EE was designed to provide EPA with an opportunity to improve its understanding of the design and application of EE methods to economic benefits analysis. The results of the pilot EE were presented in RIAs for both the Nonroad Diesel and Clean Air Interstate Rules (U.S. EPA, 2004, 2005).

The collaboration with OMB was linked to the regulatory impact assessment for the final Nonroad Diesel Rule, and thus required completion within one year. The scope of the pilot was limited to focus on the concentration-response function of PM mass rather than on individual issues surrounding an estimate of the change in mortality due to PM exposure. The limited time for completion of the pilot meant that certain aspects of a more comprehensive EE process were eliminated (e.g., neither pre-elicitation nor post-elicitation workshops were held) and some aspects of the uncertainty surrounding the PM_{2.5}-mortality relationship could not be characterized. In addition, to meet time constraints for the pilot EE, experts were selected from two previously established expert panels of the NAS.

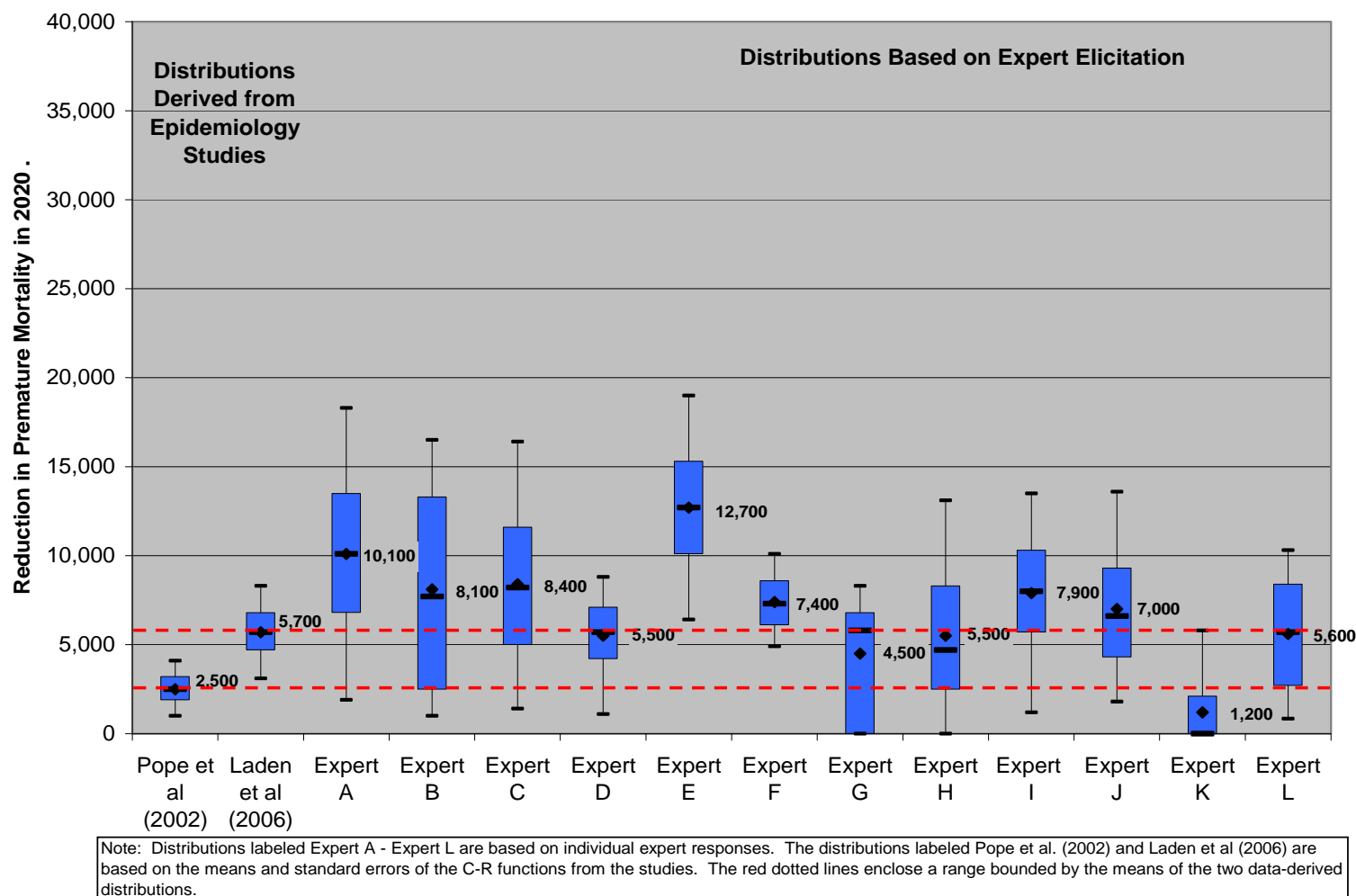
The plan for assessment and the draft protocol was initially reviewed by the Health Effect Subcommittee (HES) of the Council on Clean Air Compliance Analysis (HES Council). The protocol was pilot-tested with EPA and non-EPA PM health scientists who were not part of the final elicitation process. The project team that carried out the assessment consisted of individuals with experience in EE and individuals with expertise in PM health effects and health benefits.

EPA and OMB conducted an external peer review of the methods used in this pilot EE. In accordance with EPA's peer review guidelines (USEPA, 2005), this peer review also considered the approaches to presenting the results (particularly with respect to combining results across experts).

Based on the experience gained from the pilot EE, EPA completed a full-scale expert elicitation that incorporated peer-review comments on the pilot application. This provided a

more robust characterization of the uncertainty in the premature mortality function. The full-scale PM-mortality elicitation included an in-depth review of the protocol design, drew from a larger pool of experts using a peer-nomination process, and allowed for increased communication among experts and the project team via pre-elicitation and post-elicitation workshops. The PM-mortality elicitation was designed to evaluate uncertainty in the underlying causal relationship, the form of the mortality impact function (e.g., threshold versus linear models), and the fit of a specific model to the data (e.g., confidence bounds for specific percentiles of the mortality effect estimates). Additional issues, such as the ability of long-term cohort studies to capture premature mortality resulting from short-term peak PM exposures, were also addressed in the expert elicitation. As with the pilot EE, the full-scale PM-mortality elicitation underwent extensive review with internationally renowned PM experts, EPA management, and OMB. More details on the pilot EE and full-scale EE can be found at: <http://epa.gov/ttn/ecas/benefits.html>.

The findings from the PM-mortality elicitation were presented in the Regulatory Impact Analysis of the Final PM National Ambient Air Quality Standards (NAAQS) (U.S. EPA, 2006) and in the Proposed Ozone NAAQS (U.S. EPA, 2007). Figure 2-5 presents a depiction of results from the PM NAAQS benefit analysis, showing box plots of the distributions of the reduction in PM_{2.5}-related premature mortality based on the C-R distributions provided by each expert, as well as that from the data-derived health impact functions, based on the statistical error associated with Pope et al. (2002) and Laden et al. (2006). Each distribution is depicted as a box plot with the diamond symbol (◊) showing the mean, the dash (–) showing the median (50th percentile), the box defining the interquartile range (bounded by the 25th and 75th percentiles), and the whiskers defining the 90% confidence interval (bounded by the 5th and 95th percentiles of the distribution). The RIA also utilizes a variety of other formats for presenting the results of the elicitation, including: tables, bar graphs, and cumulative distribution functions.



Key: Closed circle = median; Open circle = mean; Box = interquartile range; Solid line = 90% credible interval

Figure 2-5. Results of Application of Expert Elicitation: Annual Reductions in the Incidence of PM-Related Mortality in 2020 Associated with the Final National Ambient Air Quality Standards for Particulate Matter

In presenting the results, EPA was sensitive to the NAS’s advice to clearly label estimates based on the results of the EE (as opposed to empirical data). In addition, EPA addressed NAS’s concerns that EE results be presented within the context of describing the uncertainty inherent in the concentration-response function. Recent RIAs have described the concentration-response functions based on the EE and considered whether these EE results should replace the primary estimated response based on analysis of the American Cancer Study

(ACS) cohort (Pope et al., 2002). This ACS cohort study was recommended by the Agency's SAB specifically for use in the 812A analyses of benefits associated with the Clean Air Act.

EPA has used or recommended EEs for other applications beyond the criteria air pollutant program, including:

- Assessing the magnitude of Sea Level Rise associated with climate change (USEPA, 1995b, Titus and Narayanan, 1996).
- *Criteria for the Certification and Re-Certification of the Waste Isolation Pilot Plant's Compliance With the 40 CFR Part 191 Disposal Regulations* [Federal Register: February 9, 1996 (Volume 61, Number 28)] -- expert judgment can be used to elicit two types of information: (1) Numerical values for parameters (variables) which are measurable only by experiments that cannot be conducted due to limitations of time, money, and physical situation; and (2) unknown information, such as which features should be incorporated into passive institutional controls that will deter human intrusion into the repository.
- Ecological model development -- EPA's National Center for Environmental Assessment (NCEA) recently undertook an effort to evaluate the utility of Bayesian belief networks (BBN) for modeling complex ecological phenomena. Because BBNs have often been cited as a promising method for modeling complex, uncertain phenomenon, NCEA undertook an effort to better understand the strengths and weaknesses of the approach. Beliefs were elicited from the panel of ecologists regarding the mechanisms by which sediment affects stream biota. Overall, the method showed promise even though a complete model was not developed in the allotted time frame.

2.6 WHAT IS THE EXPERIENCE WITH EXPERT ELICITATION AT OTHER FEDERAL GOVERNMENT AGENCIES?

Expert elicitation has been used or recommended for use by other agencies of the federal government across a broad range of applications including:

- Nuclear Regulatory Commission (NRC) – established acceptable procedures for formal elicitation of expert judgment in support of probabilistic risk assessments associated with the high-level radioactive waste program – *Branch Technical Position on the Use of Expert Elicitation in the High-Level Radioactive Waste Program* (NUREG1563) (USNRC, 1996).
-

- Army Corps of Engineers – using expert opinion to support risk studies on the assessment of unsatisfactory performance probabilities and consequences of engineered systems. *A Practical Guide on Conducting Expert-Opinion Elicitation of Probabilities and Consequences for Corps Facilities*, IWR Report 01-R-01.
- National Aeronautics and Space Administration - *Upgrading the Space Shuttle* (NAS, 1999) – Expert Elicitation should be considered as an additional formal qualitative tool in evaluating upgrades in terms of upgrades to the Space Shuttle “cost, technological readiness, contribution to meeting program goals, risks, and ability to satisfy other NASA or federal government requirements.”
- Department of Transportation / Federal Railroad Administration – use of experienced domain experts as the basis for estimating human error probabilities, -- *Human Reliability Analysis in Support of Risk Assessment for Positive Train Control*, Chapter 2, Approach to Estimation of Human Reliability in Train Control System Studies, DOT/FRA/ORD-03/15 (USDOT, 2003).
- U.S. Department of Agriculture – EEs used to support the development of a risk-based inspection program under the Food Safety and Inspection Service (FSIS) – Report to Congress, March 2, 2001. (Batz et al., 2005).

2.7 WHAT IS THE EXPERIENCE WITH EXPERT ELICITATION OUTSIDE THE U.S. FEDERAL GOVERNMENT?

Expert elicitation has been used in governmental organizations outside of the U.S. The most notable of these is the Intergovernmental Panel on Climate Change (IPCC). The IPCC has used EE for many years to address specific components of the climate change issue (e.g., biomass, temperature gradient, thermohaline circulation, aerosol forcing) as well as generating overall estimates and predictions. Its reports characterize the science and generate estimates or projections that are used by various governmental and non-governmental entities in support of climate-related decisions. The IPCC has developed formal procedures of EE for use in various aspects of the program (e.g., IPCC, 2001) and directed guidance on issues such as how to address qualitative expressions of uncertainty (IPCC, 2005).

Other international examples of EE include:

- Uncertainty analysis of NO_x emissions from Dutch passenger cars in 1998: Applying a structured EE and distinguishing different types of uncertainty, Dutch National Institute for Public Health and the Environment (RIVM, 2003).
 - European Commission, Nuclear Science and Technology, Procedures Guide for Structured Expert Judgment (2000).
-

2.8 WHAT EXAMPLES OF EXPERT ELICITATION ARE RELEVANT TO EPA?

Academia and industry have conducted numerous EEs in scientific areas that are relevant to EPA's mission, including:

- **Environmental Transport** – Amaral et al. (1983) and Morgan et al. (1984) used expert judgment in the evaluation of the atmospheric transport and health impacts of sulfur air pollution.
- **Dose-response** – Hawkins and Graham (1988) and Evans et al. (1994a) for formaldehyde and Evans et al. (1994b) for risk of exposure to chloroform in drinking water. Others used EE to estimate dose-response relationships for microbial hazards as well (Martin et al., 1995).
- **Exposure** – Hawkins and Evans (1989) used industrial hygienists to predict toluene exposures to workers involved in a batch chemical process. In a more recent use of EE in exposure analysis, Walker et al. (2001, 2003) asked experts to estimate ambient, indoor and personal air concentrations of benzene.

2.9 SUMMARY

In the 1950s EE emerged from the growing field of decision theory as a technique for quantifying uncertainty and estimating unobtainable values to support analytic decision making. For three decades, EPA has performing and interpreted EEs as part of its regulatory analysis. EEs have been conducted and used by at least five other federal agencies and international organizations. Recently, interest in using EE has increased because of encouragement from the NAS and OMB. A wide range of activities may fall under the term expert judgment; but, this White Paper restricts the term EE to a formal systematic process to obtain quantitative judgments on scientific questions (to the exclusion of personal or social values and preferences). This process includes steps to minimize the effects of heuristics or other biases in expert judgments.

3.0 WHAT IS EXPERT ELICITATION?

This chapter provides a brief review of EE research and defines terms that will be used throughout this White Paper. Included in this chapter are discussions about the origins and foundations of EE, the general reasons why EE is conducted or should be conducted (Chapter 4 provides a detailed EPA-centric discussion of this topic), the components of an EE, and some cautions and criticisms of EE.

3.1 WHAT IS EXPERT ELICITATION?

Expert elicitation is a multi-disciplinary process that can inform decision making by characterizing uncertainty and filling data gaps where traditional scientific research is not feasible or data are not yet available. While there are informal and non-probabilistic EE methods for obtaining expert judgment, for the purposes of this White Paper, EE is defined as a systematic process for formalizing and quantifying expert judgments for an uncertainty quantity as the probability of different events, relationships, or parameters (SRI, 1978; Morgan and Henrion, 1990).

It is worth noting that the use of judgment in probability calculations results from the Bayesian approach to statistics. The other statistical approach, frequentist, applies classical statistical techniques to observed data without explicitly incorporating subjective judgment. However, in many situations, especially in environmental statistics, complete or adequate data for statistical analysis do not exist; and hence, judgment must be used to analyze existing data. In addition, probabilistic statements of belief that are essential to decision making can be provided within a Bayesian framework; but, are not permitted under approach a frequentist approach.

The goal of an EE is to characterize, to the degree possible, each expert's beliefs (typically expressed as probabilities) about relationships, quantities, events, or parameters of interest. The EE process uses expert knowledge, synthesized with experiences and judgments, to produce probabilities about their confidence in that knowledge. Experts derive judgments from the available body of evidence, including a wide range of data and information ranging from direct empirical evidence to theoretical insights. Even if direct empirical data were available on the item of interest, such measurements would not capture the full range of uncertainty. EE allows experts to use their scientific judgment to interpret available empirical data and theory. It should also be noted that the results of an EE are not limited to the quantitative estimates. These results also include the rationale of the experts regarding what available evidence was used to support their judgments and how these different pieces of evidence were weighed.

3.2 WHY IS EXPERT ELICITATION NECESSARY?

EPA and other federal regulatory agencies often are required to make decisions in the presence of uncertainty. This includes situations with significant data gaps and which lack scientific consensus. The discipline of decision analysis has developed to assist decision makers who must make decisions in the face of uncertainty. Quantitative uncertainty analysis can be useful when important decisions depend on uncertain assumptions, estimates, or model choices. Because relevant data are frequently unavailable to characterize the uncertainty of the problem at hand, decisions often rely on expert judgment through informal or formal processes. EE provides a formal process to obtain this expert judgment. The various reasons that EE might be used rather than other methods of addressing uncertainty are discussed in Chapter 4.

Among federal government agencies, the NRC has the longest and most extensive experience with the conduct of EE. The NRC's Branch Technical Position (NUREG 1563) states that EE should be considered if any of the following conditions exist:

- Empirical data are not reasonably obtainable; or, the analyses are not practical to perform.
- Uncertainties are large and significant.
- More than one conceptual model can explain, and be consistent with, the available data.
- Technical judgments are required to assess whether bounding assumptions or calculations are appropriately conservative.

These conditions, and others similar situations, have motivated many EEs over the past several decades. Additional reasons for conducting an EE include:

- To obtain prior distributions for Bayesian statistical models and to help interpret observed data.
- To provide quantitative bounds on subjective judgments. Interpretations of qualitative terms (e.g., "likely" and "rare") vary widely. EE can provide numbers with real uncertainty bounds that are more useful for subsequent analyses;
- To promote consensus among experts regarding a complex decision (heretofore *a rare application*) (Cooke and Goossens, 2000).
- To provide a decision input for the prioritization of potential research options or potential decision options.

EE can produce distributions to characterize uncertainty about parameters that lack empirical data. This is particularly helpful when better scientific information is too costly to obtain, will not be available within the time frame of the decision, or is unobservable (Van Der

Fels-Klerx et al., 2002; Merkhofer and Keeney, 1987). For example, a recent study involved the elicitation of several continuous variables related to animal health safety (Van Der Fels-Klerx et al., 2002). The emphasis in this study was to obtain one aggregated probability density function for each continuous variable based on the combined (weighted) distributions obtained from the collection of individual experts. In addition, many of the analyses to assess the accident probabilities for nuclear power plants and radiation leakage from nuclear waste disposal options have relied on EE to characterize distributions for parameters or events that lack empirically frequency data (NUREG, 1996; Merkhofer and Keeney, 1987). Sometimes the needed data would come from studies that require many years to complete (e.g., a cancer study in a target population for which observations require 20 years or more). In these cases, an EE may be more expedient. In cases where direct observations are impossible (e.g., safety of nuclear facilities or the risks of terrorist attacks), EE may provide the only available information to address a particular question (O'Hagan, 2005).

3.3 WHAT DO EXPERT ELICITATION RESULTS REPRESENT?

Most scientists are comfortable with empirical observations, view such results as objective, and understand what statistical analyses represent. However, EE results and their dependence on subjective judgment are unfamiliar to most people, including many scientists. They lack an understanding of how subjective judgment can be combined with empirical data to obtain a different type of information – one that focuses on the likelihood of the nature of an unknown quantity, event, or relationship. A useful comparison of objective and subjective probabilities is (NRC 1994):

...Objective probabilities might seem inherently more accurate than subjective probabilities, but this is not always true. Formal methods (Bayesian statistics) exist to incorporate objective information into a subjective probability distribution that reflects other matters that might be relevant but difficult to quantify, such as knowledge about chemical structure, expectations of the effects of concurrent exposure (synergy), or the scope of plausible variations in exposure. The chief advantage of an objective probability distribution is, of course, its objectivity; right or wrong, it is less likely to be susceptible to major and perhaps undetectable bias on the part of the analyst; this has palpable benefits in defending a risk assessment and the decisions that follow. A second advantage is that objective probability distributions are usually far easier to determine. However, there can be no rule that objective probability estimates are always preferred to subjective estimates, or vice versa...

Subjectivity is inherent to scientific methodologies, collection and interpretation of data, and developing conclusions. In traditional scientific research, the choice of methods may influence data, which may influence conclusions. EE is no different in this respect. However,

because EE findings contain knowledge from data combined with probability judgments about that knowledge, the subjectivity is more obvious. The remainder of section 3.3 describes in more detail what the results of an EE represent. Where appropriate, this section also highlights some areas where practitioners should pay particular attention so that EE results are described accurately and represented fairly.

3.3.1 Are Expert Elicitation Results Arbitrary?

Because EE results are based on subjective judgment, there is a concern that they may be considered arbitrary. However, EE results are based on the experts' knowledge and understanding of the underlying science. To obtain EE results, experts are asked to extrapolate from their synthesis of the empirical and theoretical literature using judgments that conform to the axioms of probability. As stated above, the EE results include quantitative estimates as well as the underlying thought process or rationale. By reviewing the qualitative discussion that summarizes an expert's reasoning, one can assess whether the expert's rationale is reasonable and consistent with available evidence and theory.

3.3.3 Do Expert Elicitation Results Represent New Data or Knowledge?

Science can be thought of as two things: (1) a description of our state of knowledge of the world – what we know and don't know (epistemic evaluation of knowledge and uncertainty) and (2) the process by which we get better information (primarily to reduce uncertainty). Only the latter involves the creation of new data or knowledge. This White Paper submits that EE results encompass only the first aspect of science (characterization of existing knowledge) because no new experimentation, measurement, or observations are conducted. Furthermore, the purpose of EE is to characterize or quantify uncertainty and not to remove uncertainty which can only be done through investigative research. However, the characterization of the knowledge could better inform a decision.

This distinction is particularly important because a common reason for conducting an EE is to compensate for inadequate empirical data (Keeney and Winterfeldt, 1991) (Meyer and Booker, American Statistical Association,(eds.), 2001) (O'Hagan, 2005). In contrast, it has been suggested that EE judgments themselves be treated like data (Meyer and Booker, American Statistical Association,(eds.), 2001). However, while the results of EE can be used in ways similar to data (e.g., model inputs), one should ensure that the distinction between experimental data and EE results is maintained and the pedigree of data is clear. Users of EE results are cautioned to understand the differences between EE results and experimental data and to be aware of the role of expert judgments in EE results. For these reasons, NAS recommended that EPA identify clearly which analyses are based on experimental data and which are based on

expert judgment. This distinction should be maintained and communicated to the decision maker (NAS, 2002).

EE reflects a snapshot of the experts' knowledge at the time of their responses to the technical question. Because of this, users of EE should expect that the experts' judgments will change as the experts receive new information. An alternative approach to EE is to use experts to develop principles or rules that generalize the data so that very sparse data can be used in a broader way, i.e., provide additional certainty for sparse data. In one study, decision rules for making health hazard identifications were elicited from national experts (Jelovsek et al., 1990). The authors concluded: (1) many experts must be consulted before determining the rules of thumb for evaluating hazards, (2) much human effort is needed in evaluating the certainty of the scientific evidence before combining the information for problem solving, (3) it is still not known how experts look at uncertainty in their areas of expertise, and (4) the knowledge elicited from experts is limited but workable for medical decision making.

3.3.4 Are Expert Elicitation Results Equivalent to a Random Statistical Sample?

EE results should not be treated as a random statistical sample of the population being studied. In contrast to a valid survey that randomly samples the study population to obtain a representative sample, an EE seeks to reflect the range of credible scientific judgments. If experts are selected from multiple legitimate perspectives and relevant expertise, the EE will indicate of the range of plausible opinions. Consequently, the selection of experts is critical to the success of an EE.

3.4 WHAT ARE SOME APPROACHES FOR EXPERT ELICITATION?

This section describes the advantages and disadvantages of the two general approaches to the EE process: individual and group techniques. Chapter 5 presents how expert judgments are formally captured and lays out the process and specific steps needed to control rigorously for potential biases that may arise during elicitations.

3.4.1 Individual Elicitation

EPA's early EE efforts have primarily utilized individual elicitation techniques. This is the expert elicitation approach recommended by the NRC (1996). In it, individual experts are elicited separately using a standardized protocol. Often, the intent of these individual elicitations is to characterize uncertainty rather than defining a "best" estimate or consensus position. One advantage of individual elicitations is that the broadest range of factors contributing to overall uncertainty can be identified explicitly by each expert. In an EE involving elicitation of individuals, one can assess which parameters (quantitative or qualitative, e.g., model choice) has

the greatest impact on uncertainty. Furthermore, using individual elicitation eliminates the potential biases that arise from group dynamics.

While relying on a collection of individual elicitation does provide the most robust picture of uncertainty, it does not necessarily promote harmony or represent consensus. By encouraging a diverse spectrum of responses, some may think individual elicitation obfuscate rather than illuminate decisions. However, there are decision analytical techniques to evaluate the impact of this diversity on decisions (Clemen, 1996). Chapter 5 presents more detail on how individual elicitation are conducted.

3.4.2 Group Elicitation

A second EE approach is a group process in which experts evaluate data interactively and determine their collective judgment (Ehrmann and Stinson, 1999). By sharing data and judgments, group interactions can identify a “best” estimate or consensus opinion given the current state of knowledge. Group processes typically generate data interpretations that are different from those obtained by individual experts. These group processes include the Delphi method, nominal group techniques, group nomination, team building, and decision conferencing.

While group processes have the advantage that they can often obtain consensus, they are potentially limited by the influence of group dynamics (e.g., strong and controlling personalities). Therefore, if group techniques are used, the effect of group dynamics must be considered in addition to the general heuristic biases (Section 3.5.5). In addition, group processes that promote consensus may not characterize the full range or extent of the uncertainties. Chapter 5 includes a more detailed discussion about conducting group elicitation.

3.4.3 Combining Individual Experts Judgments

EE results from multiple experts often produce insights without combining the experts’ judgments. However, there are many circumstances where aggregated results are desired. Because EE results are often used as model inputs or information for decision makers, it may be desirable to aggregate or combine multiple expert judgments into a single metric. Section 5.4.3 provides a discussion on the advantages and cautions of combining expert judgments. There are a number of methodologies that aggregate individually elicited expert judgments. This process is different from obtaining collective judgments via a group process (Section 3.4.2). Section 3.4.3.1 presents methodologies for aggregation of individual expert judgments to produce a combined result. Section 3.4.3.2, discusses the aggregation of expert judgments by consensus processes.

3.4.3.1 Mathematical and behavioral approaches for combining individual judgments

A number of approaches have been proposed and used to combine individual expert judgments. Mathematical aggregation methods involve processes or analytical models that operate on the individual probability distributions to obtain a single combined probability distribution. Mathematical approaches range from simple averaging using equal weights (Keeney and Winterfeldt, 1991) to a variety of more complex Bayesian aggregation models. While the Bayesian aggregation methods are theoretically appealing, difficult issues remain concerning how to characterize the degree of dependence among the experts and how to determine the quality of the expert judgments (e.g., how to adjust for such factors as overconfidence). Clemen and Winkler (1999) reviewed both mathematical and behavioral approaches for combining individual judgments along with empirical evidence on the performance of these methods. Using mathematical methods to combine expert opinions relies on an assumption that the individual expert opinions are independent (O'Hagan, 1998). Behavioral aggregation approaches “attempt to generate agreement among the experts by having them interact in some way” (Clemen and Winkler, 1999). Chapter 5 provides additional discussion of the use of group processes for EE.

Based on their review of the empirical evidence evaluating both mathematical and behavioral aggregation methods, Clemen and Winkler (1999) found both approaches tended to be similar in performance, and that “simple combination rules (e.g., simple averaging) tend to perform quite well.” They also indicated the need for further work in the development and evaluation of combination methods and suggest that the best approaches might involve aspects of both the mathematical and behavioral methods. In the meantime, they express the view that simple mathematical averaging will always play an important role given its ease of use, robust performance, and defensibility in public policy settings where it may be difficult to make distinctions about the respective quality of different expert judgments.

Cooke (1990) recognized that all individuals do not possess equal skill in generating or thinking in probabilistic terms. Given similar technical knowledge, some experts are more adept at providing higher quality estimates (see section 3.5.4 about what makes a good judgment). Therefore, Cooke advocates assessing an individuals' ability to provide “statistically” robust probability estimates. To assess individual probabilistic abilities, he uses seed questions that are similar in nature to the EE's questions of interest but for which answers are known. Their performance is characterized by their statistical calibration (i.e., their ability to capture the correct proportion of answers within stated bounds) and their informativeness (i.e., the degree to which probability mass is distributed relative to background, the narrowness of the bounds). An expert's performance is gauged relative to other experts in the EE exercise and weighted

accordingly. Cooke has shown that such weighted combinations are superior to equal weighting and citation-based weighting. However, the success of Cooke's approach hinges on the quality of the seed questions in terms of their clarity (i.e., ability of the experts to correctly understand and respond to them) and their relevance to the specific problem area and question of interest. To overcome these obstacles, significant time and effort may be needed to develop, review, and evaluate these seed questions.

3.4.3.2 Consensus processes for combining individual judgments

Alternatively, individual judgments can be combined through a consensus process. This approach differs from group elicitation (Section 3.4.2) where the entire elicitation was conducted as a group. Here, the experts are elicited individually and then, as a second step, their judgments are combined via a group process. In this iterative approach, experts are allowed to discuss their original opinions and to arrive together at a collective opinion (i.e., group EE) (Meyer and Booker, American Statistical Association (eds.), 2001; Gokhale, 2001).

Under this approach the aggregation of individual expert judgments requires the experts to adjust their judgments and move toward consensus. By defining the quantitative issues of interest and removing ambiguous judgments, this process can help experts to refine their understanding of the problem and potentially narrow their differences. Thus, when used interactively, EE can aid in moving experts toward greater consensus on science-relevant problems that can not be directly measured (Cooke and Goossens, 2000; Meyer and Booker, American Statistical Association (eds.), 2001). Although not commonly used, this is a potentially useful approach to EE, particularly where the goal of the assessment is to obtain consensus views.

3.4.4 Problems Combining Expert Judgments

In individual elicitation, each expert supplies judgments on the same set of questions and combining these judgments is left to post hoc analysis. A typical EE obtains judgments from at least three experts because diversity is more likely to reflect all relevant knowledge. However, in addition to their knowledge, each expert brings different biases (and experience) to the question of interest. Therefore, EE practitioners must be cautious about aggregating expert judgments and presenting combined conclusions about EE results. Combining expert judgments may present several pitfalls, including the potential for misrepresenting expert judgments, drawing misleading conclusions about the scientific information, and adding biases to the conclusions (Hora, 2004; Keith, 1996; O'Hagan, 2005). As discussed above, individual expert judgments can be combined to provide an aggregated single average value or can be aggregated

through subsequent discussion and consensus. The chosen approach should be part of the EE study methodology and agreed-upon procedures.

According to Keith (1996), combining judgments could be problematic because the methodological assumption is that the experts chosen for the elicitation represent the entire continuum of “truth” with respect to the technical question. He cautions that the “fraction of experts who hold a given view is not proportional to the probability of that view being correct.” (Keith, 1996). This results, in part, from how the experts are selected. As mentioned in section 3.1, expert opinions are not necessarily evenly distributed across the entire spectrum of potential opinions. Furthermore, prior to interviewing experts, it may not be possible to determine the range of expert opinion on a particular question. Consequently, depending on which experts are selected (and agree) to participate in the EE, the fraction of experts used for the elicitation cannot be assumed to be proportional to the probability of that view or opinion being correct. In addition, if all else is equal and since a true value cannot be known, there is no objective basis to value the opinion of any one expert over any other.

Resolving differing expert views can be done by combining individual judgments with a mathematical method or via consensus building. Combining expert judgments requires the relative weighting of individual expert judgments to each other. They may be weighted equally or in some differential manner – for example, by social persuasion (as might occur in Delphi consensus building methods), by expert credentials, or by some form of calibration or performance assessment (Cooke, 1990). Keith (1996) argues that equal weighting of expert judgments is generally inappropriate since it is not possible to obtain a sufficiently large sample of in-depth EEs so as to ensure that all possible expert views are represented. Others have argued that equal weighting is often as effective as more sophisticated differential weighting approaches (Clemen and Winkler, 1999).

It is also possible to combine expert judgments via consensus building. Unlike the combination of individual expert judgments, which can be performed without the presence of the experts, consensus building is often a key part of a group EE process. Group elicitation provides a forum for experts to interact and exchange information, with the intention of ensuring that all experts make their judgments using the same set of prior knowledge, which they update through the deliberative process. The combination of judgments is often accomplished implicitly, by eliciting a single judgment from the entire group. Group elicitation also offers an advantage by allowing experts to collaborate and learn from each other, producing a common definition of the problem; and often, a common judgment. Allowing experts to interact could help mitigate the problem of expert selection in a particular elicitation by providing an opportunity for a wider range of opinions to be articulated and explored among the expert group than they may have

individually expressed on their own. A disadvantage of this type of group elicitation is that the social dynamics and interaction may lead to an overly narrow uncertainty characterization, especially if minority views that express a broader range of uncertainty are swamped by the goal of reaching a consensus judgment. It is therefore important that minority opinions and their rationale also be presented to decision-makers.

Other EE practitioners also urge caution about combining the individual judgments (Wallsten et al., 1997). In a methodology similar to Jevolsek et al (1990), Wallsten et al. (1997) proposed a model developed by his team that considers both “the structure of the information base supporting the estimates and the cognitive processes of the judges who are providing them.” Wallsten et al. determined where experts agree and derived rules that satisfy those conditions. The resulting model avoids some of the criticisms of combining expert judgments when subjective inputs for data and the processes used by the experts in the elicitation are not considered.

Other formal methods have also been devised that combine individual and group elicitation (e.g., Delphi). In all of these cases, one can expect that the final elicited judgments will vary with the methods that are selected. Therefore, care must be exercised to use elicitation methods that are most appropriate for a particular problem.

3.5 WHAT ARE GOOD PRACTICES FOR ELICITING EXPERT JUDGMENT?

Having established the utility of EE, this White Paper will now concentrate on good practices for eliciting expert judgment based on a literature review and actual experience within EPA and other federal agencies.

3.5.1 Why is a Probabilistic Approach Needed to Elicit Judgments?

EEs provide insights into parameter values, quantities, events, or relationships and their associated uncertainty in support of decision making. A common mathematical language is needed to provide a rigorous, credible, and transparent assessment to support decisions. For EE assessment, the common language that is most effective at ensuring usability of results and comparability across experts is probability. Although subjective terminology (e.g., “likely” or “unlikely”) can convey probabilities, numerous studies have shown that a natural language approach is inadequate because:

- The same words can mean very different things to different people.
 - The same words can mean very different things to the same person in different contexts.
-

- Important differences in expert judgments about mechanisms (functional relationships) and about how well key coefficients are known can be easily masked in qualitative discussions.

Wallsten et al., (1986) documented that individual interpretations of words can differ dramatically if they are presented without context. In this study, they evaluated ten qualitative descriptions of likelihood (almost certain, probable, likely, good chance, tossup, unlikely, improbable, doubtful, and almost impossible). For each description the participants expressed an associated probability. The range varied considerably between participants, including overlap across words such that some were indistinguishable.

Similarly, Morgan (1998) presented the results of an exercise in which he queried the members of the SAB Executive Committee at a time when EPA was considering moving toward a more qualitative description of cancer hazard. He asked the committee members about their interpretations of the terms “likely” and “not likely.” This exercise found that the minimum probability associated with the word “likely” spanned four orders of magnitude, the maximum probability associated with the word “not likely” spanned more than five orders of magnitude, and most importantly, there was an overlap of the probability associated with the word “likely” and that associated with the word “unlikely.” Because interpretations of qualitative descriptions have such high inter-individual variability, a quantitative framework is needed for experts to provide comparable and tractable expressions of belief. Probability can provide this framework; and in particular, the subjectivist approach to probability is ideally suited for this application. Subjective probability is a formal expression of the degree of belief about some unknown quantity; and therefore, is ideal for quantifying uncertainty in expert beliefs.

Because the general population is continually exposed to and familiar with probabilistic information, it provides a useful and consistent framework from for eliciting expert judgments. For example, people are fairly comfortable interpreting weather forecasts such as one that calls for a 30% chance of precipitation. However, most people, scientists included, do not think in terms of fully formed probability distributions and providing probabilities for unknown events, quantities, relationships, or parameters is a non-trivial exercise. Therefore, specialized techniques (as discussed in Chapter 5) are available to facilitate obtaining such estimates from experts.

3.5.2 What are Problems with Probability? Are there Alternatives?

Shackle (1972a) states that “probability is...a word with two quite opposed uses. It is at the same time the name of a kind of knowledge, and an admission of a lack of knowledge.”

When subjective probabilities sum to one, it implies omniscience and that all alternative hypotheses have been considered. However, it is often the case that statisticians using subjective probability do not know all the hypotheses and cannot set up a statistical problem that satisfies this major premise (Shackle, 1972b). Consequently, Shackle asserts that when statisticians use statistics to draw conclusions about data, they need to be mindful that the statistics may be derived from an incomplete set of potential hypotheses.

While a useful tool, the field of statistics involves simplification of real world problems and subjectivity is intrinsic in its methodologies. Pretending that statistical approaches are objective may result in misplaced confidence in data and conclusions. In EE, the expression of expert judgment as probabilities assumes that experts understand all alternatives so that their judgments can be compared. This may be true for a binary outcome where the expert is asked for the probability of occurrence (or non-occurrence). However, in most situations, the expert is asked to make judgments about the probability of one event occurring compared with the occurrence of multiple other events, some of which may be unknown. This challenge is further complicated by “self-reinforcing” or inherently evolutionary systems (Shackle, 1972b). Evolutionary systems have elements of unpredictability (some events are completely unrelated to previous events) that make it unacceptable for using probability to describe the system because probability contains an inherent assumption of stability within a given system.

While expert judgment is commonly expressed solely in probabilistic terms, there are other feasible approaches. Meyer and Booker (2001) define expert judgment as “data given by an expert in response to a technical problem.” Using such a definition, it is possible to obtain expert judgments in a variety of non-probabilistic forms. Expert judgment is often used where data cannot be collected practically or are too expensive to assemble. Quantitative but non-probabilistic methods for expressing expert judgment have been commonly used in decision analysis. Such approaches tend to use pair-wise comparisons and stated preferences among the pairs. Doing so does not require the expert to formally give probability estimates for an event, parameter, or relationship. This method has been particularly useful for supporting decisions in which values and preferences (not just scientific evidence) are considered. However, it can also be used to elicit expert judgment about a specific technical problem or scientific interpretation as well. Analysts who are considering the use of EE but are reluctant because of concerns about probabilistic approaches may find these alternative methods for expert judgment more suitable.

3.5.3 What Makes a Good Expert?

The intent of an EE is to characterize the state of knowledge by integrating available evidence with scientific judgment to provide as complete a picture as possible of the relevant

knowledge regarding a particular question. Elicitation of expert judgment provides a useful vehicle for combining formal knowledge, as reflected in the published literature, with expert judgment. Hence, there are two aspects that define “good” experts. The first is an understanding of the body of literature for the problem of interest. However, this technical knowledge alone does not define a good expert. Experience and judgment, including intuition, and the ability to integrate information and theories beyond the reported data are also critical. Technical knowledge, experience, and judgment ability play critical roles in obtaining good expert judgments.

3.5.4 What Constitutes Good Expert Judgment?

A well-conducted EE should reflect accurately the selected experts’ judgments and capture the “truth” within the range of expert judgments. EE goes beyond empirical observation, which, in general, can not capture the true estimate of uncertainty. Therefore, good expert judgments should consider more than just the statistical confidence limits from empirical studies.

A good judgment properly captures the range of uncertainty; but, it still should be reasonable. Some individuals are more capable of formulating and expressing their judgments probabilistically than others. Cooke (1990) identified characteristics of good probability judgment:

- ***BEING CALIBRATED OR STATISTICALLY ACCURATE*** – a good probability judgment is one that mimics the underlying probability of predicting the “truth” if it were known. In other words, the credible intervals presented by the experts should capture the “true” value within the expressed credible intervals, i.e., 90% confidence intervals should include 90% of the true values. Furthermore, the estimates should be balanced. For example, 50 percent of any “true” values should be above, and 50 percent should be below an expert’s estimated median values.
 - ***INFORMATIVENESS*** – a good probability judgment is one where the probability mass is concentrated in a small region (preferably near the true value) relative to the background rate.
-

As illustrated in Figures 3-1 and 3-2, ideally one prefers experts whose judgments are unbiased and precise. Building on that premise when it comes to expert judgments one would like to be as unbiased as possible where the central mass is closest to the true value. In addition, the credible

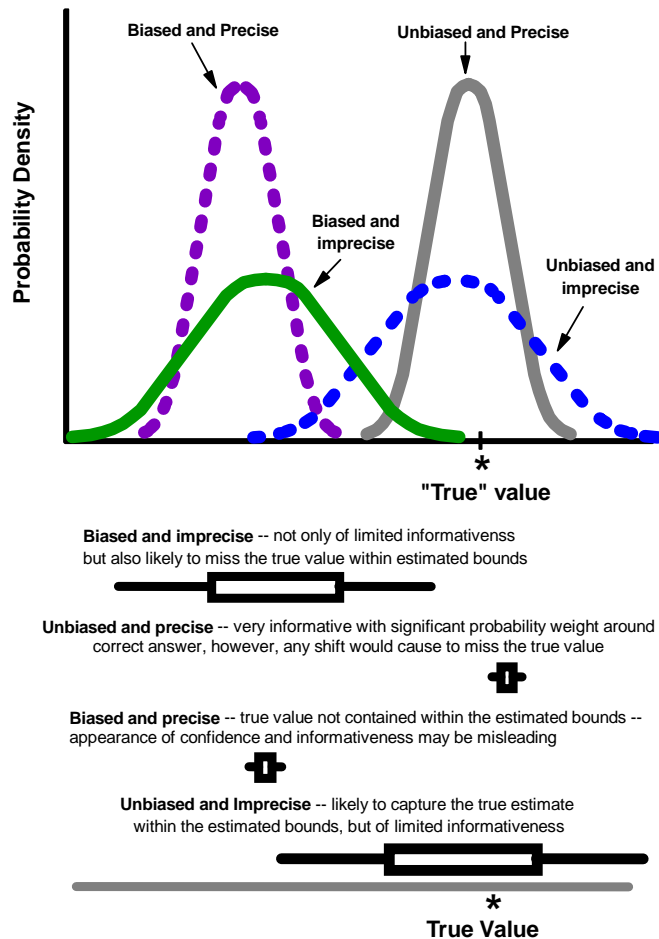


Figure 3-1. Distributions to Illustrate Bias and Precision

limits should be sufficiently broad so as to include the proper proportion of “true” values. However, a good expert should not have bounds that are too broad so as to reduce the mass or confidence around the true value. In addition to expressing probabilities that are statistically robust it is also important that experts describe clearly the information they used to support their opinions. Experts who can express the basis for their judgments are strongly preferred.

3.5.5 Where can the Elicitation Process Go Awry? Why?

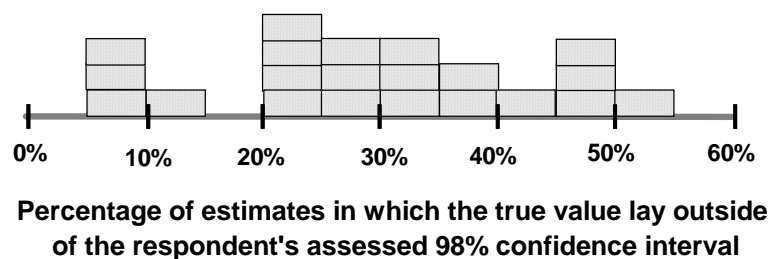
Most lay people, and even experts, do not have well-formed probability distributions for quantities of interest *a priori*. Instead, a process is necessary to conceptualize and develop these

probability values. During this process, experts use existing experience and knowledge to infer the probabilities that are elicited. To accomplish this task, most people (including subject matter experts) make use of simple rules of thumb called “cognitive heuristics.” In some instances, these heuristics can lead to biases in judgment. The psychometric literature (e.g., Kahneman, Slovic, and Tversky, 1982) describes the heuristic biases that most impact EE judgments:

- Overconfidence
- Availability
- Anchoring and adjustment
- Representativeness bias
- Motivational bias

3.5.5.1 *Overconfidence*

One consistent finding across all elicitation techniques is a strong tendency toward overconfidence (Morgan and Henrion, 1990). Figure 3-2 provides a histogram summary of the results of 21 different studies on questions with known answers and the observed surprise indices from a wide range of studies of continuous distributions. The surprise index is the percentage of time the actual value would fall out of the elicited credible range (e.g., 98 percent credible interval). Ideally the actual value would occur outside of the elicited credible range 2 percent of the time. However, as Figure 3-2 illustrates the surprise index is almost always far too large, ranging from 5 to 55 percent instead of 2 percent.



Summary of data provided in Morgan and Henrion (1990)

Figure 3-2. Summary of “Surprise Index” from 21 studies with known answers (Morgan and Henrion, 1990)

Given this tendency toward overconfidence there has been significant interest in whether training potential experts, using trial tasks such as encyclopedia questions, can improve performance (Morgan and Henrion, 1990). Table 3-1 summarizes the impacts of several training experiments and their impact on overconfidence. Some studies attempted to reduce the overconfidence by explaining prior performance and exhorting the experts to increase the spread of their estimates. Results showed modest decreases in overconfidence. Other studies that

provided comprehensive feedback showed significant improvement on discrete elicitations; however, improvement was marginal for continuous distributions. The nature of the feedback is critical. In particular, it may be important to include personal discussion in feedback.

	Number of assessed distributions <i>N</i>	Interquartile index (ideal = 50)		Surprise index (ideal = 2)	
		before	after	before	after
<i>Alpert & Raiffa (1969)</i>					
Groups 2 & 3	1,670	33	44	39	23
Group 4	600	36	43	21	9
<i>Schaefer & Borcharding (1973)</i>					
Fractiles	396	23	38	39	12
HFS	396	16	48	50	6
<i>Pickardt & Wallace (1974)</i>					
Group 1 (5 sessions)	?	39	49	32	20
Group 2 (6 sessions)	?	30	45	46	24
<i>Schaefer (1976)</i>					
(5 sessions)	660	27	34	25	14
<i>Lichtenstein & Fischhoff (1980)</i>					
(Training on discrete tasks)	924	33	37	41	40

Table 3-1. Impact of Training Experiments on Reducing Overconfidence

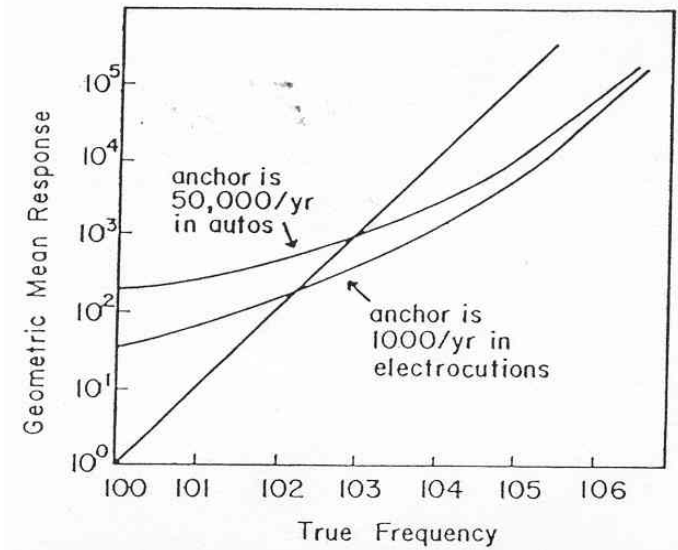
3.5.5.2 Availability

Probability judgment is also driven by the ease with which people can think of previous occurrences of the event or can imagine such occurrences. In Figure 3-3, the results of elicited judgments about annual deaths rates from various causes are plotted against actual death rates. The judgments tend to overestimate the occurrence of rare events and underestimate the rate of more common causes of death. The explanation may be that rare events receive relatively more publicity and are more readily recalled. Therefore, they are believed to occur more frequently than they do. By comparison, deaths from common causes receive minimal publicity and therefore elicited judgments may underestimate their true rates.

3.5.5.3 Anchoring and Adjustment

Anchoring occurs when experts are asked a series of related probability questions. Then, when forming probabilities, they may respond by adjusting values from previous questions. For example, if an expert is first asked to state the median probability of an event, this stated probability for the median may become an “anchor.” Then, when responding to subsequent questions about other quantiles (e.g., 10th percentile) of the distribution, the responses may be influenced by the anchor.

Probability judgment is frequently driven by the starting point which becomes an “anchor.” For the example shown in Figure 3-4, people were asked to estimate annual death rates from different causes. They



They were provided with either the true annual deaths rate for autos (50,000 per year) or the true annual deaths for electrocutions (1,000 per year). The given death rate provided an anchor and influenced the results. As can be seen from the graph, estimated death rates were shifted by which reference value was given. This given value becomes an anchor and subsequent estimates are made relative to it by adjustment.

Figure 3-3: Availability Bias (Lichtenstein et al., 1978)

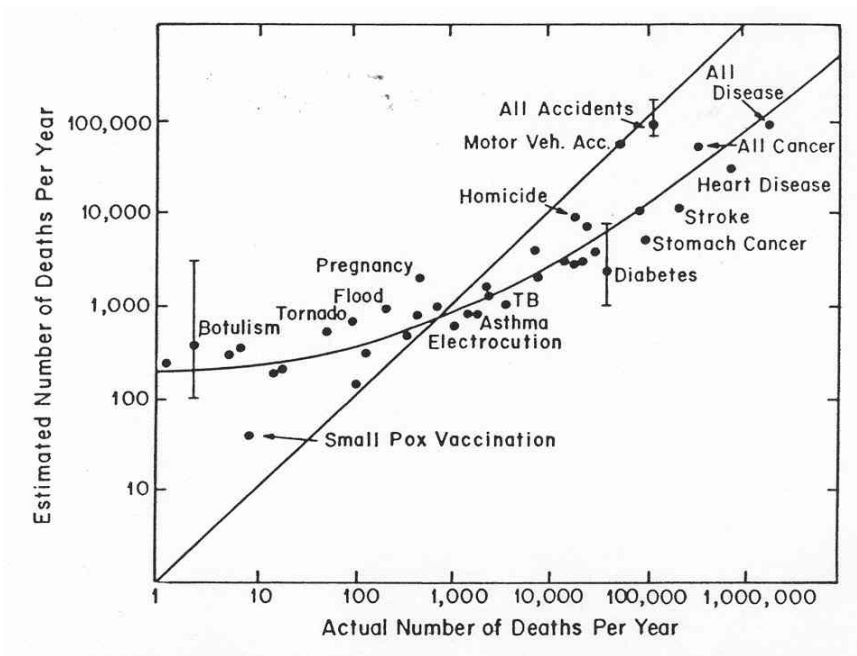


Figure 3-4: Anchoring for estimating of death rates (Lichtenstein et al., 1978)

In an observation of flawed probabilistic assessments, Bruine de Bruin et al. (2002) found that respondents to some probability questions tend to have an “elevated frequency” of 50%

responses. This disproportionately high number of “50%” responses does not reflect true probabilistic beliefs. Rather, it is “caused by intrusion of the phrase ‘fifty-fifty,’ which represents epistemic uncertainty, rather than a true numeric probability of 50%.” These non-numeric 50%s may be an artifact of the question format and the elicitation methodology that are compounded by the elicitee’s understanding of probability. Because treating these non-numeric 50%s as true 50%s could lead to erroneous conclusions, Bruine de Bruin et al. (2002) presents two redistribution techniques to mitigate this difficulty.

3.5.5.4 Representativeness Bias

People also judge the likelihood that an object belongs to a particular class based on how much it resembles that class. This phenomenon can be illustrated considering the following example. Suppose one flips a fair coin ten times. Which of the following two outcomes is more likely?

Outcome 1: T, T, T, T, T, H, H, H, H, H

Outcome 2: T, T, H, T, H, T, T, T, H, H

Both sequences are equally likely; but, the second may appear more likely because it seems to better represent the underlying random process. By contrast, the first sequence gives the appearance of a non-random pattern. In general, people tend to underestimate the occurrence of patterns or sequences that appear to be non-random.

3.5.5.5 Motivational Bias

Frequently, experts may have direct or indirect interests in the outcome to the question at hand. Hence, whether consciously or unconsciously, their judgments may be influenced by motivational bias. In some cases the stakes may be clear (e.g., when the outcome of a question may impact employment or investments). For other cases, motivational bias may be subtler. For example, the professional reputation of a particular expert may be associated with a particular point of view or theory, making it difficult to express an alternative perspective (Morgan and Henrion, 1990).

The regulatory process is characterized by complex multi-factor problems where decisions are influenced by technical analyses as well as social, economic, and political considerations. Furthermore, this process involves multiple stakeholders each with their own positions, frames, and agendas. As a result, motivational biases are among the more elusive and yet critical biases to consider for EEs that support regulatory decisions. The existence of motivational bias may be difficult to demonstrate; but, the adversarial nature of regulatory decisions suggests that motivational biases may exist. In part, this is due to the fact that people

tend to trust that their judgments are less prone to bias than those of others. One explanation of this misperception may be that people tend to rely on introspection for evidence of bias in themselves but on lay theories when assessing bias in others. As a result, people are more inclined to think they are guilty of bias in the abstract than in specific instances. Also, people tend to believe that a personal connection to a given issue is for them a source of accuracy and enlightenment but for others it is a source of bias (Ehrlinger, et al., 2005). Because of the importance of transparency and appearance in the regulatory process, motivational bias may be an important consideration whether or not it actually influences an assessment.

3.5.5.6 Cognitive Limitations of Experts

Human limits on cognition restrict the complexity of the relationships that can be attempted with EE. Hence, eliciting conditional probabilities relating three or more variables can be very difficult (O'Hagan, 2005). Poor performance when the number of variables increases can be caused by many different factors. Among the major factors are the degree of correlation among the variables (a higher degree of correlation produces more confusion), human information processing capacity, and barriers to learning (Fischhoff, 2003). In addition, Hamm (1991) argues that expert performance in EE can be compromised when the expert does not make numerical judgments carefully because he does not understand the model or is unfamiliar with the elicitation language that is different from how he expresses himself within his field (Hamm, 1991). Even when they are involved in model construction, experts tend to think about a few familiar cases rather than consider all applicable cases. This well-documented phenomenon has motivated the development of many decision support tools (Chechile, 1991; Saaty, 1990).

Another limitation is that experts may not update their beliefs when new information becomes available (Meyer and Booker, 2001). Judgments are contingent on the consideration of all possibilities and the assessment of a relative belief in one possibility compared with the other possibilities. Shackle (1972b) argues that when experts do not know all alternative possibilities, the premise of subjective probabilities is violated.

3.5.5.7 Experts and Controversial Issues

Mazur (1973) found that, when issues are controversial (e.g., nuclear power and water fluoridation), experts performed as non-experts. He noted that many conflicts were not disagreements among experts but rather, arguments about different points, i.e., a different understanding/definition of the problem. The conflicts resulted from divergent premises, usually resulting from poor communication between adversaries. When scientific controversies contain subtle perceptions and nuances, each side's proponents must simplify the conclusions in a manner that results in one side choosing to accept as a conclusion what the other side regards as

an unproven hypothesis. Then, each side seeks to gain acceptance of its view through non-scientific persuasion. The real uncertainty of the science is “removed” by proponents of each side as they state their own case with increasing certainty, not necessarily supported by the scientific information. In addition, the desire to “win” the argument may heighten conflicts.

Experts can take different positions along the uncertainty/ambiguity continuum and are subject to the same biases as lay people. For both experts and lay people, controversy heightens emotion and subjugates ambiguities. Mazur concludes that the more controversial the issue, the more experts tend to behave like non-experts. Therefore, when the issue is controversial and the science/scientific analysis is ambiguous or uncertain, the value of that science/scientific analysis may become more questionable. The process of EE, if conducted professionally, may reduce the emotion and ambiguity in an assessment by breaking down a problem into distinct questions and carefully and explicitly defining the uncertain quantities, events, relationships, or parameters of interest.

3.5.5.8 Additional limitations of eliciting experts

Three additional cautions are offered for the decision of whether or how to use EE. The first pertains to an expert’s carefulness when translating qualitative thinking to quantitative judgment. Hamm (1991) states that if experts are careless, they may unintentionally respond without adequately understanding the implications of their responses. Sometimes, the carelessness is an artifact of experts failing to understand how their responses will be used. Furthermore, when the decision context is either not shared with or understood by the expert, factors that might have been considered important when responding may be forgotten or ignored.

The second caution pertains to whether and to what degree models correspond with how experts think of the problems presented to them. By design and necessity, models are simplifications. They provide approximations of reality; but, in the process abstraction, may alter the problem being analyzed.

The third caution pertains to the need to balance analysis and experience (Hammond et al., 1987). Hammond states that judgment, even by experts, includes subjectivity. This does not mean that the judgment is arbitrary or irrational; however, the use of scientific information requires interpretation of its meaning and significance. Because the facts do not speak for themselves, there is a need for experts to provide judgments about the information. When experts from the same discipline disagree, it is important to evaluate where and why the disagreements arise. The analysis of disagreement can be the source of important insights.

3.5.6 How can the Quality of Expert Judgments be Improved?

These limitations do not preclude the use of EE methods, rather they provide targets for their improvement. Heuristic biases are well-recognized by EE practitioners and methods have been developed to control those biases. For example, anchoring biases can be mitigated by first eliciting extreme values of distributions (e.g., the highest and lowest possible values) and then the median. Also, the same information can be elicited redundantly to help validate values. In any case, EE practitioners are urged to be mindful of these biases when designing elicitations and interpreting results. An understanding of these biases highlights the need for rigor in the design and implementation of an EE protocol. This rigor can produce a more credible analysis. Chapter 5 describes many of the approaches to reducing the impact of such biases and heuristics.

3.5.7 How can the Quality of an Expert Elicitation be Assessed?

In most circumstances, the true accuracy of an EE can not be quantified. Uncertainty exists in the experts' understanding of the process and in the process itself. Disentangling these two sources of uncertainty is difficult. Furthermore, experts are not value-free (Shrader-Frechette, 1991; Slovic et al., 1988) and bring their own biases to the elicitation (Renn, 2001; Renn, 1999; Slovic et al., 1988). While experts' training and experience add credence to their judgments, it is still a judgment and incorporates values. Expert judgments are a fusion of science and values. Garthwaite et al. summarizes by stating that a successful elicitation "faithfully represents the opinion of the person being elicited" and is "not necessarily 'true' in some objectivistic sense, and cannot be judged that way" (Garthwaite et al., 2005). In the case of EE results being used as missing data, it is important for EE practitioners to understand that they are not obtaining traditional experimental data. If expert judgments are being combined (as discussed in section 3.5.3, Garthwaite et al.'s caution become even more important as the combination of judgments presumes that the "truth" lies within the spectrum of those judgments.

3.6 SUMMARY

For the purposes of this paper, EE is defined as a formal process for developing quantitative estimates of the probability of different events, relationships, or parameters using expert judgment. The use of EE may be of value where there is missing data that either can't be obtained through experimental research, lack of scientific consensus, and/or the need to characterize uncertainty. However, as with many other analytical tools, there are also significant theoretical and methodological cautions for its use. As with experimental data, users of the elicited expert judgments must be aware of the biases incurred due to the choice of methods and to biases that are due to the fact that experts are human too. Expert judgment does not exist in a vacuum apart from sociological and personality influences. Therefore, the most defensible uses

of expert judgments obtained through probabilistic EE are those that consider the impacts of these biases on the final results and provide good documentation for how the expert judgments were obtained and how they will be used.

4.0 WHAT THOUGHTS ABOUT APPLICABILITY AND UTILITY SHOULD INFORM THE USE OF EXPERT ELICITATION?

This Task Force recognizes that EE is one of many tools to characterize uncertainty and/or address data gaps. Many factors influence whether an EE would be helpful and should be used by EPA, including: (1) the purpose and scope of the EE (e.g., to estimate missing data or characterize uncertainty), (2) the nature of available evidence and the critical uncertainties to be addressed, (3) the nature of the overarching project or decision that the EE will support, (4) the potential time and resource commitment required for conducting EE, and (5) the impact on the decision in the absence of the data provided by the EE.

The advantages and disadvantages of EE should be evaluated in light of the particular application for which it is being considered. An EE may be advantageous because it uses experts to help characterize uncertainties, allowing analysts to go beyond the limits of available empirical evidence. This is especially important when additional data are unavailable or unattainable. An EE may be disadvantageous where the perceived value of its findings is low and/or the resource requirements to properly conduct an EE are too great. Given resource constraints, it may be necessary to balance the time, money, and effort needed to conduct a defensible EE against the requirements of other forms of expert input such as external peer review.

As discussed previously, EEs are conducted typically to address unresolved uncertainties and/or to fill gaps when additional data are unattainable within the decision time-frame. However, even when these needs are clearly articulated, determining whether or not to conduct an EE requires specific consideration of several factors. This chapter reviews some questions and issues that influence the decision to conduct an EE.

The process for considering formal EE to characterize uncertainty and/or address data gaps can be divided into three steps:

1. How important it is to quantitatively characterize major sources of uncertainty or address a critical data gap, in a particular case?
2. Is EE well-suited for characterizing the uncertainty or for providing estimates to address a particular data gap?
3. Is EE compatible with the overall project needs, resources, and timeframe?

Each of these will be discussed in turn, below.

4.1 HOW IMPORTANT IS IT TO CONSIDER UNCERTAINTY?

To support its decisions, EPA often conducts complex assessments that draw on diverse expertise. In many cases, empirical data are unavailable on the outcome of interest. Therefore, EPA may rely on extrapolations, assumptions, and models of the real world. These abstractions of reality are all sources of uncertainty in the assessment. When risk or cost information are presented to decision makers and the public, the findings have often been reduced to a single numerical value or range. This approach may provide insufficient information to the decision maker and has the hazard of conveying undue precision and confidence.

In their text *Uncertainty*, Morgan and Henrion (1990) present five criteria that define when considering uncertainty is important:

- When people's attitudes towards uncertainty are likely to be important (e.g., if uncertainty itself is likely to be an argument for avoiding a policy option)⁶.
- When various sources of information need to be reconciled or combined, and some are more certain than others (i.e., where more weight should be given to the more certain information).
- When deciding whether to expend resources to collect more information (e.g., prioritizing possible areas of research or data collection, or deciding whether to seek additional data).
- When the "expected value of including uncertainty" (EVIU)⁷ is high, such as when the consequences of underestimating would be much worse than for overestimating (e.g., underestimating the time needed to get to the airport may have severe consequences, so uncertainty about traffic has to be taken into account).

4.1.1 What is EPA's Position on Characterizing Uncertainty?

EPA has long recognized the importance of characterizing uncertainty. To that end, it has established Agency-wide policy and guidance that give significant attention to how uncertainty is characterized and dealt: EPA's risk characterization policy (EPA, 1992; EPA, 1995a), *Principles of Monte Carlo Analysis* (EPA, 1997), *Risk Characterization Handbook*

⁶ Numerous studies have shown that many people will choose a course of action that has a more certain outcome over one with a relatively uncertain outcome, even if the expected net benefits are somewhat higher in the uncertain case. They are willing to give up the additional expected benefit, just to avoid the uncertainty. This behavior is called "risk aversion" and the "risk premium" is the amount of benefit they are willing to give up in exchange for avoiding the risk of loss (i.e., avoiding the uncertainty).

⁷ EVIU, as defined by Morgan and Henrion (1990), refers to the quantitative impact that uncertainty analysis can have on a decision. From a decision-analytic perspective, the EVIU is a measure of how much will the expected value outcome of a decision will increase if uncertainty is included in the analysis. If considering uncertainty can lead to a decision with a higher expected value outcome, then the EVIU is high.

(EPA, 2000a), *Risk Assessment Guidance for Superfund (RAGS) Volume 3* (EPA, 2001a), and the *Risk Assessment Staff Paper* (EPA, 2004a).

EPA's *Risk Characterization Handbook* (USEPA, 2000a) makes clear that "it is generally preferred that quantitative uncertainty analyses are used in each risk characterization." Furthermore, it states that "even if the results are arrived at subjectively, they will still be of great value to a risk manager." EPA's *Guidelines for Preparing Economic Analyses* (USEPA, 2000b) presents a tiered, practical approach:

In assessing and presenting uncertainty the analyst should, if feasible: present outcomes or conclusions based on expected or most plausible values;... [and as an initial assessment,] perform sensitivity analysis on key assumptions."... "If, however, the implications of uncertainty are not adequately captured in the initial assessment then a more sophisticated analysis should be undertaken... Probabilistic methods, including Monte Carlo analysis, can be particularly useful because they explicitly characterize analytical uncertainty and variability. However, these methods can be difficult to implement, often requiring more data than are available to the analyst.

EPA practice often involves a "tiered approach" to conducting uncertainty analysis. Hence, EPA often starts as simply as possible (e.g., with qualitative description) and sequentially employs more sophisticated analyses (e.g., sensitivity analysis to full probabilistic). These additional analyses are only added as warranted by the value added to the decision process (USEPA, 2004a).⁸ The *Risk Characterization Handbook* (USEPA 2000a) provides examples of the appropriate way to characterize uncertainty. This approach focuses on the need to balance limited resources, time constraints, and analytical limitations against the potential for quantitative uncertainty analysis to improve the analysis and regulatory decision.

4.1.2 How can Expert Elicitation Characterize Uncertainty and Address Data Gaps?

In general, EPA has used EE to address data gaps and/or characterize uncertainty surrounding estimates of important quantities, such as the health impacts of a specified change in air quality. Because EE can provide subjective probability distributions that quantify uncertainty estimates, it is often suggested as an analytic method worth considering. For example, the

⁸ While it may be important to *consider* each source of uncertainty in an analysis, it may not make sense to *quantify* every uncertainty. This is an important distinction because quantifying uncertainty is sometimes very difficult and not always very useful. Some sources of uncertainty can be adequately addressed with a qualitative discussion, and a judgment that quantitative analysis is not merited.

National Academy of Sciences (NAS, 2002) recommended that EPA consider greater use of EE to quantify uncertainties estimates in its health benefits analyses.

In general, EE can be useful when:

- Acceptable quantitative estimates of uncertainty can not be made adequately with additional data collection (e.g., cannot be observed, such as future oil prices), can not be observed *directly* (e.g., effects of a new substance on human health), or the events are so rare that data are very limited (e.g., risk of nuclear plant accident). Statistical methods can not address this type of data limitation. When empirical data are essentially impossible to obtain, EE is a viable approach to quantification.
- Uncertainty estimates using other techniques will not be quantified adequately because of the time frame for a decision or decisions about available resources. Situations may arise where data collection would require more time than analyses based on expert judgment, where data collection is not technically feasible, or where the benefits of additional data collection in terms of improved confidence may not justify the cost and/or time.

As defined in this document, EE is a formal process for developing quantitative estimates for the probability of unknown events, relationships, or parameters using expert judgment (SRI, 1978; Morgan and Henrion, 1990). EE goes beyond empirical data and allows experts to integrate across various lines of evidence. When data are unavailable or unattainable, EE may be used to either fill data gaps and/or to characterize uncertainty.

4.1.3 What are the Alternatives Methods for Expert Judgment?

As described in Chapter 2, EE is one of many expert judgment methods, including activities that range from informal to formal. Other expert judgment methods include public comment and peer review. These methods vary in their level of rigor and the degree to which they control for heuristics and biases. They also differ in the range of questions that they can address, the level of effort and resources required, and the degree of public acceptability. Table 4-1 presents basic descriptors for expert judgment methods that should be considered when determining if an EE should be conducted. One should consider whether such an activity is compatible with the timeline, available resources, and the overall nature of the issue and decision-making process. Table 4-1 compares various methods of expert judgment in terms of resource needs. It should be noted that these forms of expert judgment are not necessarily comparable in terms of purpose and their ability to provide information to decision makers and stakeholders. If detailed quantitative characterization of uncertainty is necessary, then EE could

be compared to a range of uncertainty methods. The estimates provided in Table 4-1 focus on EE and should be compared to other methods to characterize uncertainty.

Table 4-1. Illustrative Comparison of EE and Other Methods for Expert Judgment

	Public Comments	Limited (letter) peer review	Formal FACA peer review⁹	Expert elicitation
Problem addressed	Broad, no limit, defined by commenter	Broad, but defined by charge	Broad, but defined by the charge	Narrow, specific and well-defined
Timing	Typically 45 days	1-4 months	4-12 months	8 months – 2 years
Resource needs	limited	~\$25K	~\$250K	~\$250K - \$2M
Role of Public/ Stakeholders	Open to all to provide comments	Formal selection process	Public nomination, selection process, open public process	Nominations by peers and limited involvement of public/stakeholders
Evidence considered	No limit	No limit	No limit	No limit, but must be formally shared with all experts to evaluate
Acceptance	Publicly acceptable	Familiar though not transparent to public	Generally accepted, recognized	Some wary of method (i.e., concerns about perceived bias)
Selection of experts	None	Formal selection process	Formal and public nomination process	Systematic process usually involving nomination by technical experts

⁹ Note: Review by the National Academy of Sciences (NAS) could also be included in this category. If so, time and resource requirements may be substantially increased over an EPA-led FACA review such as with the SAB.

4.2 WHAT IS THE NATURE OF THE UNCERTAINTIES TO BE ADDRESSED?

Many different sources of uncertainty may arise when models and analyses are used to support policy analysis and decision making. This section presents different sources of uncertainty and discusses how EE could be used to address them. As discussed below, EE has the potential to be helpful in characterizing uncertainty regardless of its source.

4.2.1 What are the Categories of Uncertainty?

Based on the literature (Cullen and Frey, 1999; Finkel, 1990; Hattis and Burmaster, 1994), sources of uncertainty can be classified into four categories:

- ***Input (or parameter) uncertainty:*** Models and assessments utilize a wide range of parameters and other inputs to generate estimates. Typically the values for these inputs are not known with confidence. Among the factors that can introduce uncertainty into model inputs are random error (including lack of precision in measurement), systematic errors (i.e., bias), lack of empirical data, and lack of representativeness of empirical data.
 - ***Model uncertainty:*** All models include uncertainty about the appropriate modeling approach (e.g., which model best represents reality, including how inputs and constants should be combined in equations based on an understanding of the real world). Because models are simplified representations of the real world, uncertainty can result from imperfect knowledge about the appropriate conceptual framework, specific model structure, mathematical implementation, detail (precision/resolution), boundary conditions, and extrapolations, as well as choices among multiple competing models.
 - ***Scenario uncertainty:*** Decisions related to the overall design of the scenarios modeled in the analysis (e.g., selection of receptor populations, chemicals, exposure sources, and study area delineations) can be a source of uncertainty for analyses.
 - ***Decision rule uncertainty:*** Decisions on the types of questions asked to support policy making and the theoretical framework used to make those decisions can introduce uncertainty. These areas of uncertainty include: (a) the design of the decision framework used in guiding policy making (e.g., acceptable risk levels and the types of risk metrics used such as individual- versus population-level risk) and (b) global protocols used in the analysis (e.g., use of standard risk reference doses (RfDs) and associated hazard quotient (HQ) values as the basis for non-cancer risk assessment, versus the use of epidemiologically-based disease incidence estimates).
-

4.2.2 Which Types of Uncertainty are Well-Suited for Expert Elicitation?

If the problem statement and questions can be formulated clearly and consensually, then EE may be used to address any category of uncertainty. If an adequate knowledge base exists and there are qualified experts, then their judgments can form a credible basis for judgments that can provide insight about any type of uncertainty.

Analyses that support EPA decisions typically involves numerous components, such as risk assessments that include toxicity assessments, emissions or discharge estimates, air or water quality modeling, exposure assessment, economic impact analyses, and so on. For each of these steps, EE may be valuable; but, EE may be more useful or appropriate for some steps than others. Therefore, it is important to identify a very specific problem statement and questions when deciding whether to use EE in an analysis. When EE is considered for use in an analysis, EE may be more or less appropriate depending on what specific questions the EE would be used to address (Morgan and Henrion, 1990).

4.3 WHAT ARE OTHER METHODS TO CHARACTERIZE UNCERTAINTY?

In addition to EE, other methods are available to characterize uncertainty. Some methods can both characterize the uncertainty of particular parameters and propagate this uncertainty through the model. For example, EE can be used to develop subjective probability distributions; characterize the uncertainty of specific parameters, events, quantities, or relationships; and estimate the overall uncertainty of a modeled process. The context for applying these methods to characterize uncertainty depends greatly on the user's perspective. The focus could be on a single modeling parameter. For example, an EE could be conducted to estimate the magnitude of a cancer slope factor (including that number's uncertainty). However, the carcinogenic process for that chemical could be viewed as a complex multi-element process with multiple modeling steps. In this case, the estimates from an EE may be used to propagate uncertainty through an entire model (e.g., the cancer process).

Methods for probability-based uncertainty characterization can be divided into five broad categories: (a) statistical/frequentist (e.g., Monte Carlo and Latin Hypercube Simulation), (b) judgmental/subjectivist (e.g., EE, Bayesian), (c) scenario analysis, (d) other (e.g., interval, probability bounds, fuzzy logic, and meta analysis) and (e) sensitivity analysis techniques. These categories and methods are discussed briefly below:

- **Statistical/frequentist:** These uncertainty characterization methods are based on the frequentist paradigm and hence require empirical data to establish a probabilistic characterization of uncertainty. These approaches treat probability as an objective measure of likelihood based on frequencies observed in data that are subject to
-

- sampling error, measurement error, and other random processes. For example, wind speed and its associated uncertainty could be described by reporting the range, mean, and 95th percentile of historical measured values. Common methods that are founded on the frequentist paradigm include numerical uncertainty propagation, bootstrap, and response surface methods. Some of these methods, such as bootstrap, can quantify uncertainty even with very small sample sizes. In general, they are less capable and used less often to characterize uncertainty about model choice or causality. In addition, they typically can not address uncertainty arising from data that are not representative of the value to be estimated (e.g., an epidemiological study that focused on a population that is very different from the one to be analyzed in a risk assessment where no data on the population differences are available).
- ***Judgmental/Subjectivist:*** These methods are based on the concept that probability is an expression of the degree of confidence in some parameter, quantity, event, or relationship. In addition, they are based on the concept of logical inference—determining what degree of confidence an expert may have, in various possible conclusions, based on the body of evidence available. Common subjectivist methods include Bayesian analysis, EE, and Generalized Uncertainty Likelihood Estimation.
 - ***Scenario analysis:*** Uncertainty can be characterized through presentation of alternative scenarios that are thought to span the range of plausible outcomes. Scenario analysis is useful to evaluate groups of variables or assumptions that are correlated and/or vary together (e.g., worst-case scenario), to predict future conditions, and to assess model uncertainty.
 - ***Other methods:*** In addition, there is a group of diverse methods that do not depend heavily on subjective judgment (as does Bayesian analysis) and can be applied in contexts where uncertainty characterization is limited by inadequate empirical data. These methods occupy a middle ground between the frequentist and subjective methods. They include interval methods, fuzzy methods and meta-analysis.
 - ***Sensitivity analysis techniques:*** These methods assess the sensitivity of the results to choices of inputs, assumptions, or model. However, sensitivity analysis does not necessarily quantify the probability of those alternatives choices. Methods for sensitivity analysis include local methods (these examine the impact of individual inputs in relative isolation on model outputs); combinatorial methods (varying two or more inputs simultaneously while holding all other inputs constant and determining the impact on model output); and global methods (these generate output estimates by varying inputs across the entire parameter space and determine contribution of individual inputs to overall uncertainty).
-

4.3.1 How Does Expert Elicitation Relate To The Other Methods?

EE falls within the judgmental/subjectivist category of methods. These methods have the advantage that they can provide a robust characterization of uncertainty without requiring as much data as frequentist approaches.

For a particular analysis, the uncertainty characterization is not limited to a single method; multiple methods may be used. For example, EE may be employed to generate a probability distribution for an input parameter while a frequentist approach is used for other input parameters. Alternatively, EE can be applied to assess the appropriateness of specific model choices while a frequentist approach can be drawn on to address uncertainty in the inputs to those models.

4.3.2 How Relevant And Adequate Are Existing Data?

As mentioned above, EE may be useful when empirical data are severely limited or are contradictory. An assessment of the adequacy of any empirical data and theory should be part of a decision to use EE requires. Data may be limited in quantity, quality, or both; adequate data may be difficult or even impossible to obtain. Quality problems may arise from relevance of the data or problems with imprecision or bias in the data.

4.3.2.1 What types of evidence are available?

Rarely is there direct empirical data that are specific to the quantity of interest within the exact context of interest. In other words, EPA rarely has direct observations of the impacts of environmental pollutants at concentrations encountered in the environment within the specific exposed population. As a result, EPA often makes inferences based on lines of evidence (Crawford-Brown, 2001). These five types of evidence can be ordered by relevance as shown in Figure 4.1

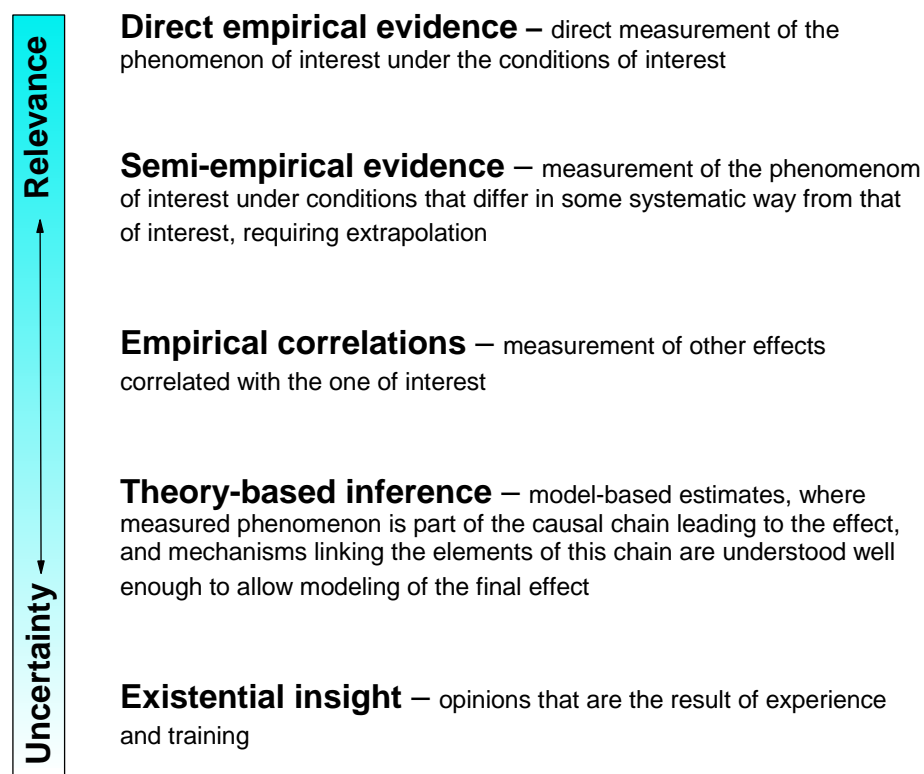


Figure 4.1. Five Types of Evidence

A key consideration for assessing the adequacy of empirical data, and the possible utility of EE, is to evaluate the representativeness of existing data (e.g., studies are limited to animal data, non-U.S. populations, or unique subpopulations). EPA’s *Report of the Workshop on Selecting Input Distributions for Probabilistic Assessments* provides useful methods for assessing and handling data with respect to its representativeness (USEPA, 1999b). The option of working with real data, given their limitations, should be considered as an alternative to EE. One should consider whether there are suitable approaches to adjust data so that it is sufficient to support a specific decision. Furthermore, EE may be a useful approach allowing experts to provide judgments based on adjusting empirical data that are not entirely relevant, such as in determining how data from another population may be used to represent the U.S. population.

4.3.2.2 What is the quality of the available information?

The type and relevance of data are not sufficient to evaluate the need for EE. A second dimension of data is its quality. One may have a study that represents direct empirical evidence but the quality of the study may be poor (e.g., poor level of detection). In some cases, EPA has only a very limited database of information on a critical element of an analysis (e.g., very small

sample size measuring a poorly-characterized population). For example, there may be only five known measurements of a certain emissions rate and additional measurements may be very costly to obtain. Even assuming the data points are representative, which is not necessarily the case, such a small sample provides very limited information. The average value can be estimated, and the degree of uncertainty over that average can be estimated as well. Separating variability from uncertainty in this case would present another challenge.

There are powerful statistical methods for estimating distributions that describe parameter uncertainty and variability. They should be considered vis-à-vis EE. EPA has utilized and documented these methods in a variety of analyses. The report *Options for Development of Parametric Probability Distributions for Exposure Factors* (USEPA, 2000c) provides a good discussion of such methods. Techniques exist that can provide useful information about uncertainty using even very small data sets (i.e., <20 data points) e.g., Cullen and Frey (1999). These methods should be explored as options before launching into an EE to solve the problem of data limitation.

4.3.2.3 Are there critical data gaps?

The many uncertainties within an assessment are likely to have varied impacts on the overall assessment or confidence in that assessment. Sensitivity analysis can identify the most critical uncertainties and these should be the focus of uncertainty analysis.

As described in Section 5.2.1, an EE will typically be of most value when it is possible to define a very limited number of critical uncertainties. Because defensible EE efforts can be very resource intensive to provide insights about a particular event, quantity, parameter, or relationship, it is best to identify those areas that would most benefit from this concentrated analysis. Various techniques are available to identify the relative contribution of various uncertainties (Renn, 1999; Wilson, 1998; Warren and Hicks, 1998; and Stahl and Cimorelli, 2005). One of these approaches could be used to help identify appropriate targets for an EE.

4.3.2.4 Is the state of uncertainty acceptable?

An inference from data and the quality of that inference are based on the entire body of evidence. Typically, multiple categories of evidence are available and it may be difficult to determine how to combine these bodies of evidence. This combination can be accomplished by the judgment of a scientific expert who will evaluate the categories of evidence, the weights given to those categories, and the quality of evidence within each particular category.

For EPA, decision making under uncertainty is inherent and unavoidable. Overall uncertainty affects the correctness and acceptance of decisions. By using uncertainty analysis,

we can improve our understanding of evidence, its impact on the overall estimate, and our confidence in that estimate.

In general, the acceptability of uncertainties is judged by decision makers, scientists, and stakeholders after the data is placed in the context of a decision or used for a decision. This is because the extent to which data uncertainty impacts the availability and value of decision options is not known until the context of the decision is clear. Hence, the acceptability of uncertainty is necessarily a contextual decision. In some circumstances, existing uncertainty is acceptable and in others it is not. Further discussion about the role of uncertainty in a decision context is available elsewhere (Jamieson, 1996; Renn, 1999; Wilson, 1998; Warren-Hicks and Moore, 1998; Harremoes et al., 2001; Stahl and Cimorelli 2005).

4.3.2.5 Is knowledge so limited that EE would not be credible?

In some situations with scant empirical data, expert judgments along with sound theoretical foundations can form the basis for EE. However, where knowledge is scant, the EE should not be seen as a proxy for creating “data” where none exist. In such circumstances, it may be appropriate to use EE to support scenario analysis for ranking rather than a formal quantification of options. Quantifying options may give the appearance of greater precision than is defensible.

Within complex decisions, debate and opposition are often related to specific concerns of stakeholders and the diverse ways that they frame issues. Because individuals react to risks on multiple levels, including analytical, emotional, and political, differing approaches can lead to unnecessary conflict. This is supported by recent brain imaging research that is exploring how people react to risk and uncertainty. Hsu (2005) has shown that neural responses are different when reacting to risk (risk with a *known* probability, based on event histories, relative frequency, or accepted theory) vis-à-vis ambiguity (risk with *unknown* probability, or uncertainty about risk levels – meager or conflicting evidence about risk levels or where important information about risk is missing). Attempting to quantify uncertainties may be at odds with how people perceive and react to a situation. Therefore, this may affect the credibility of any such attempt. In fact, the International Risk Governance Council (IRGC, 2005) recognized the importance of this factor in the development of its integrated framework, in which knowledge is categorized to allow distinction between “simple,” “complex,” “uncertain,” and “ambiguous” risk problems.

This does not imply that EE is inappropriate for these purposes; rather that care must be taken to conduct a decision-making process in which stakeholders will view EE as a credible process. In these circumstances, it is worth considering whether the major stakeholders would view EE as credible. For example, such an approach was used successfully to forecast the

potential microbial impact of a Mars landing (North, 1974). This problem involved an area with no history or expectation of empirical data. In this case, the success of the EE exercise may have been because it was a technical exercise and had limited stakeholders. However, EPA tends to be involved with problems that are complex and involve numerous diverse stakeholders (e.g., regulated community, environmental NGOs, community members, and independent scientists). EE can be used to support scenario analysis (such as focusing on mental models); but, as it becomes more quantitative it may become challenging. This challenge will be particularly critical when data is scant, theoretical foundations on which to base judgments are limited, and/or multiple stakeholders are involved in the discourse. In addition, the use of uncertain data tends to be more accepted when stakeholders know how and for what purpose that data will be used. Placing uncertain data within a particular decision context may limit the use of that data and increase stakeholder acceptance for using that data for that limited purpose.

4.3.2.6 Can additional data be collected to reduce uncertainty?

The possibility of seeking additional data should be considered carefully in any decision about using EE. In addition, one should also evaluate the option of working with imperfect data (e.g., not representative), by using suitable approaches to adjust it. As for any analysis, the proposed use of EE data should be clearly defined in the context of the assessment. A final option is using techniques to obtain useful uncertainty assessments for very small data sets (e.g., less than 20 data points). These methods for using imperfect data should be considered before the data limitation is established as the rationale for using EE. However, one should bear in mind the time and resources needed to augment or improve data, which itself can be a lengthy process. In many circumstances, EE may provide information more promptly.

4.4 WHAT ROLE MAY CONTEXT PLAY FOR AN EE?

EE is suitable for many EPA activities, including identification of research needs, strategies, and priorities; risk assessments for human or ecological health; and cost-benefit analyses to support major regulatory decisions. The context of each potential use, including the level of scientific consensus, the perspectives of anticipated stakeholders, and the intended use of results may indicate whether it is appropriate to use EE.

4.4.1 What Is The Degree Of Consensus Or Debate?

Another consideration about whether to rely on existing data or to conduct an EE is the degree of consensus in the scientific community. One of EE's strengths is that it provides the carefully considered and fully described views of several highly-respected experts who are affiliated with diverse institutions and perspectives. Obtaining these cross-institutional viewpoints may be preferable to relying on the views of an in-house expert, judgments from an

advisory committee, or otherwise limited data. When evaluating the status of scientific consensus or debate, the following factors may indicate that EE is applicable:

- Conflicting empirical evidence and lack of consensus on selecting analytical options.
- No clear consensus exists and there is substantial debate among experts.
- The problem concerns an emerging science challenge and/or the scientific controversies include model selection or use and/or data selection or use.
- The range of views are not easily articulated or captured by EPA's existing professional judgment processes (e.g., analytical approaches and external peer review).
- Problems are complex and multidisciplinary and hence need methodical deliberation by a group of experts to become tractable.

4.4.2 Will Stakeholders View Expert Elicitation as Credible?

Given the novelty of EE to many potential user communities, stakeholder concerns should be considered early in the process. One potential concern is that subjective assessments may be viewed as unreliable and unscientific. Other stakeholders may have concerns regarding potential bias and manipulation of expert judgments (e.g., ABA, 2003, NRDC, 2005). Consequently, if EE is to be used in regulatory decisions, transparency is critical for credibility. The presentation of results should take into account what is known about effective risk communication (Bloom et al, 1993; Johnson, 2005; Morgan et al., 2001; Slovic et al., 1979; Slovic, 1986; Thompson and Bloom, 2001; and Tufte, 1983; USEPA, 1994; 1998; 1999b,c; 2001a,b). This topic is discussed further in Chapter 6.

As part of this communication strategy, stakeholders may need to be shown that subjective judgment (the linking of facts and judgments) is a component of many environmental analyses, not just EE. Nevertheless, some stakeholders may be unfamiliar with the central role that subjective judgment plays in EE. Therefore, one should consider stakeholder perspectives as part of the overall process of deciding whether and how to conduct an EE. Early interactions and effective communication with stakeholders may help to garner support and satisfy their desire to play a role in this process.

OMB's *Guidelines for Ensuring and Maximizing the Quality, Objectivity, Utility, and Integrity of Information Disseminated by Federal Agencies* (USOMB, 2002) defined "quality" as encompassing objectivity, utility, and integrity. Objectivity is a measure of whether information is accurate, reliable, and unbiased, and whether it is presented in an accurate, clear, complete, and unbiased manner. Furthermore, the guidelines also highlight the need for reproducibility in that "independent analysis of the original or supporting data using identical methods would

generate similar analytic results, subject to an acceptable degree of imprecision or error.” These OMB guidelines recognize that some assessments are impractical to reproduce (e.g., ones that concern issues of confidentiality) and hence do not require that all results be reproduced. In any case, a carefully conducted and transparent EE can meet these goals for objectivity and reproducibility. To demonstrate the feasibility of reproducibility, some research (Wallsten et al., 1983) showed that expert judgments are stable within a reasonable period.

In general, OMB has been very supportive of EE as a means for EPA to improve uncertainty characterization, as indicated by the following statements and actions:

- “In formal probabilistic assessments, expert solicitation is a useful way to fill key gaps in your ability to assess uncertainty.”(*Circular A-4*, USOMB, 2003c)
- “Judgmental probabilities supplied by scientific experts can help assessors obtain central or expected estimates of risk in face of model uncertainty.” (*Draft OMB Risk Assessment Bulletin*, USOMB, 2006).
- OMB played an active role in the non-road diesel EE pilot to demonstrate how one might comply with Circular A-4.

4.4.3 What is the Nature of Review or Dialogue of the Overarching Activity?

The nature of the decision-making process and its context can affect whether EE is appropriate. Many EPA decisions are multi-factorial, including technical, social, political, and economic elements that involve multiple stakeholders. EE can be used to address specific scientific or technical issues; but, if the decision is inherently political, then EE may not provide helpful insights. When values rather than science are critical, other decision analytic methods may be preferred. Some other methods that facilitate decision making and promote stakeholder interaction include the Analytical Hierarchy Process (AHP), Multi-Criteria Decision-Making (MCDM), and Multi-Attribute Utility Theory (MAUT).

The analytic approaches that are used to support Agency decisions differ in their degree of stakeholder involvement and participation. These processes include negotiated rulemaking, extensive dialogue, and detailed technical analyses and review. Although the EE process can be transparent and offers stakeholder review, it generally lacks an opportunity for active stakeholders input. If direct stakeholder interaction is desired, it can be added as a supplemental activity to the basic EE process.¹⁰

4.4.3.1 Is There a Perception of a Major Bias Among Stakeholders?

¹⁰ Although stakeholders may nominate experts, review the protocol, and review the final EE, they may seek a more participatory process than is available via EE.

As with any peer or expert judgment activity, the credibility of an EE may be influenced by the overall reliability of the experts involved and whether those experts are perceived as biased. The psychological literature has shown that people tend to see others as more susceptible to a host of cognitive and motivational biases (Pronin et al., 2004). In general, people may be predisposed to believe that their own personal connection to a given issue is a source of enlightenment; but, for those with opposing views, such personal connections may be a source of bias (Ehrlinger et al., 2005). Therefore, stakeholders have a natural tendency to be cautious when evaluating the stated preferences of experts with opposing views.

The success or acceptance of an EE may depend on the openness of stakeholder communications, efforts taken to address potential bias, and transparency. As described in Chapter 5, a well conducted EE attempts to mitigate motivational biases. Although EE results should reflect a range of valid opinions, one should not select biased experts. If the process is transparent, any biases should be evident. When deciding whether to conduct an EE, for a topic with potentially biased experts, one consideration may be the difficulty of obtaining unbiased credible experts.

4.4.3.2 What Is The Role Of Peer Review in the Overall Effort?

EPA's *Peer Review Handbook* (3rd Edition) outlines the process for determining when to conduct a peer review and for selecting appropriate mechanisms and procedures to conduct peer reviews (USEPA, 2006). Formal peer review is used frequently for assessments or portions of assessments that support EPA regulatory decisions. According to OMB guidelines, highly influential information should undergo external peer review to ensure its quality.

Peer review of an EE should include subject matter experts and experts in the use of EE. As with any peer review, it may be challenging to obtain peer reviewers if the technical domain is small and the pool of relevant experts is very limited. It is important to note that an expert who participates in an EE becomes part of the analysis and would be unable to conduct a peer review of the resulting product. This may make it especially difficult to find sufficient number of experts for both the EE and peer reviewers for the final product. See Chapter 5 for a more detailed description of the process and criteria that might be used to select experts for an EE.

4.4.4 How Will Expert Elicitation Results Be Used?

EPA faces many decisions, with varying levels of quality requirements,¹¹ for which EE is potentially applicable. For example, EE may be relevant for the following decisions: identify

¹¹ OMB Information Quality Guidelines, 2002 states "We recognize that some government information may need to meet higher or more specific information quality standards than those that would apply to other types of government information. The more important the information the higher the quality standards to which it should be held."

research needs, develop research priorities, make regulatory decision, make major regulatory decision (greater than \$100 million impact) – with increasing importance and therefore increasing requirements for information quality. EEs can be costly and resource intensive undertakings. In an effort to reduce cost and time requirements, one may wish to eliminate EE elements that control for biases and heuristics (see Section 4.5). This may reduce the overall quality of the EE. Whether the diminished quality is critical depends on the planned use of the EE’s results. Hence, when planning an EE, it is important to consider the use of its results.

As with all analytic activities, the use of results should be guided by the design of the protocol and the purpose for which they were developed. When considering a secondary use of results, one should consider the rigor of the protocol design and whether the necessary elements were elicited. For example, if the EE was developed for internal deliberative purposes (e.g., to identify research needs), depending on design protocol, it may be inappropriate to use the results for a regulatory decision that has higher information quality requirements. If it is expected that the results may be used for other purposes (especially uses with higher quality requirements), this may be considered during protocol design. It is prudent to consider whether demands to use results beyond the intended purpose may exist, the potential impact for any misuse, and whether additional resources are needed to ensure the robustness of results.

4.5 WHAT RESOURCES ARE REQUIRED FOR AN EXPERT ELICITATION?

Expert elicitation as defined by this White Paper are generally ambitious undertakings. As described in more detail in Chapter 5, careful attention should be given to the design and conduct of any EE effort in order to minimize the impact of heuristics and biases. The cost of conducting an adequate defensible EE includes both EPA resources and time and in some cases contractor support. Table 4.2 provides a general outline of the various portions of the EE and considerations regarding time and effort.

Predicting resource needs for EE is not straightforward. The resources for an EE, like any complex analysis, depends on the design, scope and rigor desired. As previously discussed, a well-designed EE controls for heuristics and biases, thereby, elevating the quality and credibility of results. Controlling for such heuristics and biases usually requires extensive planning and protocols. Although numerous methodological adjustments can be implemented that would lower the level of effort and hence resource needs, such adjustments can affect the overall quality and/or acceptability of results.

This section briefly describes the resource implications and considerations for conducting an EE (see Chapter 5 for additional discussion).

Table 4.2. Potential Pre-Expert Elicitation Activities, Responsibilities, and Time.

Activity	Responsibility	Time
Defining the Problem		
Problem definition/scope	EPA Project Manager Senior Managers	2 months
Structuring/decomposing the problem	EPA Project Manager	2-6 months
Identify expertise needed for EE	EPA Project Manager	2 months
Contracting		
Contract planning, secure funding, contract capacity, expertise in EE, write SOW, bids, selection	EPA Project Manager and contracting staff	3-5 months
Financing contractor	EPA Contracting Officer	1-2 years (entire project)
Selecting Experts		
Development of selection approach and criteria, identification and recruitment	EPA Project Manager	1-3 months
Review of nominated experts, selection of experts, evaluation of conflicts of interest	EPA Project Manager and/or Contractor's Project Officer	1-2 months

4.5.1 How Long Does an Expert Elicitation Take?

A well-conducted and rigorous EE that adequately controls for biases and heuristics can be a lengthy process and require several years. Table 4.2 provides a general picture of the time needed for the individual steps. If the EE includes more complex analyses, additional resources may be required. The time estimates in this table are part of a continuum of activities, from simple to sophisticated. The level of effort for any given EE is contingent on the type of complexity of the assessment, which is influenced by how the information will be used.

4.5.2 What Skills are Needed to Conduct an Expert Elicitation?

Conducting an EE requires a breadth of activities and skills, each with its own resource implications. The necessary skills as summarized here and described in more detail Chapter 5, can be broken down into two categories: organizational and technical. First, the conduct of an EE involves an expenditure of EPA resources (work years, contracting funds, etc.). As a result, the EE project manager identified in Table 4.2 is involved in many steps of the project. Secondly, significant technical skills are needed to support an EE effort. In addition to expertise about the subject matter, there should be experience in the EE process itself. To properly capture and document expert judgment requires both a thorough understanding of the state-of-science for the particular discipline and expertise in the EE process (e.g., cognitive psychology). Whether these expertise can be met internally or require contractor support depends on the skills and availability of EPA staff. Furthermore, additional skills will be needed by the overall project team to define clearly the problem, develop the protocol, assemble the body of evidence to be presented to the experts, and review the results. The steps outlined in Chapter 5 require involvement of both EPA staff and specialized contractors. As described in Chapter 5 a team approach may be most effective to conduct the actual elicitations.

4.5.3 How Much Does an Expert Elicitation Cost?

As is detailed in Table 4.2, a well-conducted rigorous EE may require a large resource commitment. Past EPA experiences indicates that such efforts can range from \$200K to \$2M, depending on the level of effort and rigor. This is generally consistent with external cost estimates for conducting an EE with individual face-to-face elicitations that report a range of \$100,000 to \$1,000,000 (Hora and Jenssen, 2002; and Moss and Schneider, 1996, respectively). Adjustments can be made to the process that may provide cost savings. As discussed in Section 5.3.4, alternatives to face-to-face elicitation may reduce costs. However, such adjustments will typically lessen the rigor that controls for heuristics and biases, and/or reduce transparency, thereby diminishing the overall quality of the results.

4.6 SUMMARY

As described above, there are many technical, administrative, political, and procedural factors that influence whether to conduct an EE. This section summarizes circumstances where EE might or might not be appropriate. In most cases, EE is but one of several methods which can be used to characterize or address critical uncertainties or data gaps. For these situations, this chapter showed how various factors may be evaluated to select the preferred method. In such cases, the decision to conduct an EE is not clear and may be influenced by many factors. For a given project, analysts and decision makers need to integrate the numerous factors discussed above to facilitate a decision on whether to conduct an EE.

4.6.1 What Conditions Favor Expert Elicitation?

The following conditions tend to favor EE:

- The problem is complex and more technical than political
- Adequate data (of suitable quality and relevance) are unavailable or unobtainable in the decision time framework.
- Reliable evidence or legitimate models are in conflict.
- Appropriate experts are available and EE can be completed within the decision timeframe.
- Necessary financial resources and skills are sufficient to conduct a robust and defensible EE.

What Conditions Suggest Against Expert Elicitation?

The following conditions tend to suggest against EE:

- The problem is more political than technical
 - A large body of empirical data exists with a high degree of consensus.
 - The findings of an EE will not be considered legitimate or acceptable by stakeholders.
 - The information that the EE could provide is not critical to the assessment or decision.
 - The cost of obtaining the EE information is not commensurate with its value in decision-making.
 - Available financial resources and/or expertise are insufficient to conduct a robust and defensible EE.
 - Other acceptable methods or approaches are available for obtaining the needed information that are less intensive and expensive.
-

5.0 HOW IS AN EXPERT ELICITATION CONDUCTED?

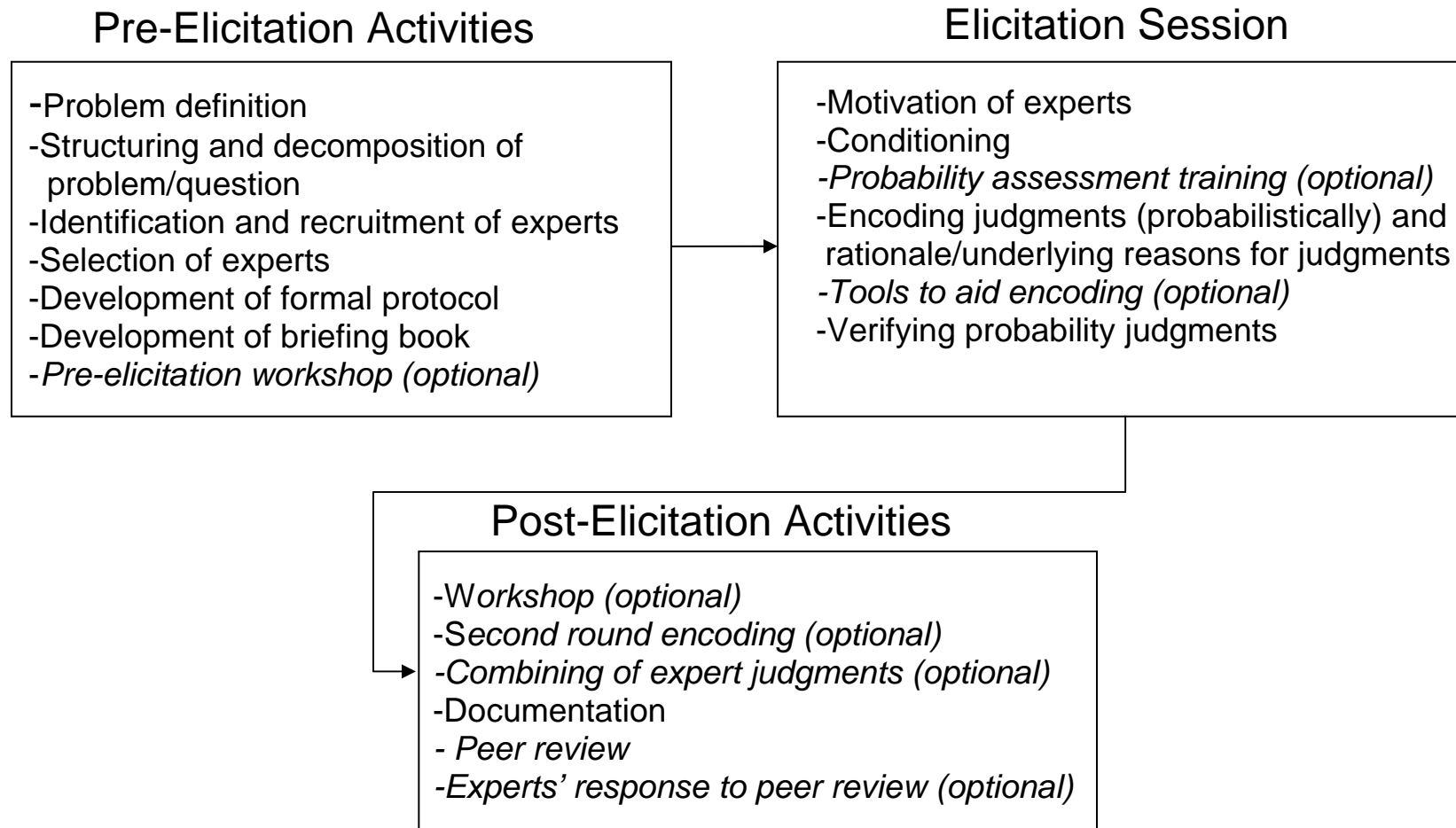
This chapter summarizes the major steps and important factors to consider when conducting an EE. It also describes good practices for conducting an EE based on a review of the literature and actual experience within EPA and other federal agencies. These practices may apply to EE conducted by EPA or by outside parties for submission to EPA. In general, these good practices consist of the following elements: 1) clear problem definition, 2) appropriate structuring of the problem, 3) appropriate staffing to conduct EE and select experts, 4) protocol development and training, including the consideration of group processes and methods to combine judgment when appropriate, 5) procedures to check expert judgments for internal consistency (verify expert judgments), 6) clear and transparent documentation, and 7) adequate peer review.

In practice, a range of approaches may be used. Hence, the protocol design for any particular EE involves considerable professional judgment. As stated by Morgan and Henrion (1990), “the process of expert elicitation must never be approached as a routine procedure amenable to cookbook solutions ... Each elicitation problem should be considered a special case and be dealt with carefully on its own terms.” As noted in Chapter 3, an EE may be used to estimate an uncertain quantity or parameter, an uncertain relationship, or an uncertain event. In the discussion that follows, the phrase “uncertain quantity” is used to represent any of these circumstances.

5.1 WHAT ARE THE STEPS IN AN EXPERT ELICITATION?

There are a number of different approaches or types of EE processes. This White Paper focuses on EEs that involve individual interviews of the experts. However, some EEs entail group processes (e.g., Delphi, group survey, and nominal group technique). In the sections below, the primary discussion concerns EE processes that focus on individuals. In addition, many, but not all, of the elements discussed below are relevant for both individual and group process EEs. Where there are significant differences with respect to the application of individual or group processes, this is discussed in sections 5.2 to 5.4.

Based on decision analysis texts and articles (Morgan and Henrion, 1990; Spetzler and von Holstein), Figure 5.1 provides an overview of the various steps included in a full EE. The overall process includes pre-elicitation activities (see section 5.2), conducting the

Figure 5.1. Overview of Expert Elicitation Process

elicitation (see section 5.3), and post-elicitation activities (see section 5.4). There are some steps that should be followed in all EEs, while other steps are optional (see Figure 5.1). The value added for these optional steps is discussed in the following sections. One should recognize that these optional EE steps also have additional costs of time and expense. Additionally, the strengths and weaknesses of EE identified and discussed in Chapter 3 should be fully considered when an EE protocol is being designed.

5.2 WHAT ARE THE PRE-ELICITATION ACTIVITIES?

Pre-elicitation activities are shown in Figure 5-1 and presented in the sections below. Some of these steps can be carried out in parallel; but, others have prerequisites that require proper sequencing.

5.2.1 What is a Problem Definition?

The initial step of an EE is to craft a problem definition that describes the objectives precisely and explicitly. What is the purpose of the EE and what are the questions that need to be addressed? Is the purpose of the elicitation to inform a regulatory decision, to guide/prioritize research needs, or to help characterize uncertainty in a regulatory impact analysis? Laying out the objectives of the elicitation is critical to the design of the EE in guiding the choice of experts, determining how information is presented to them, and determining the form of the judgments that will be required” (USNRC, 1996).

Regardless of the purpose of the EE, it is critical to its success that the uncertain quantity of interest is clearly and unambiguously defined. The quantitative question must pass the, so called, clarity or clairvoyant test (Morgan and Henrion, 1990). This “test” is fulfilled if a clairvoyant could, in theory, reveal the value of the uncertain quantity by specifying a single number or distribution without requesting any clarification. This demands that all of the significant assumptions and conditions that could impact the expert’s response are well-specified.

One should also define the uncertain quantity in such a way that the expert is able to apply his knowledge as directly and fully as possible without necessitating mental gymnastics (Morgan and Henrion., 1990). In addition, it is important that the uncertain quantity be specified in such a way that it adequately addresses the policy or analytical question(s) of interest. For example, in the chronic ozone lung injury elicitation, as briefly described in section 2.5.2.1.3 (Winkler et al., 1995), the definition of mild lesions in the centriacinar region of the lung were those that could be detected by sophisticated measurement methods such as an electron microscope; while moderate lesions were those that could be detected with the naked eye. These definitions were readily understood by the toxicology experts that participated in this elicitation.

5.2.2 How is a Problem Structured and Decomposed?

The problem definition can be designed with either an “aggregated” or “disaggregated” approach. The choice of one of these approaches will be influenced by the type of experts available, their perceived breadth of knowledge, and their ability to use integrative analysis. For the two approaches, questions will be structured and posed differently. For the aggregated approach the uncertain relationship of interest will be obtained through a single complex question. For example, if the quantity of interest is the probability that exposure to chemical x at concentration y will lead to a 10% increase in mortality, one could ask this question directly to experts. Alternatively, if following a disaggregated approach there will be a series of simpler, more granular, questions to the experts. Multiple types of experts may be used so that each can be asked a specialized question for their expertise. For example, dosimetry experts would be asked the probability that exposure to chemical x at concentration y will result in a given internal dose. Then, health scientists would be asked to provide the probability that a specified internal dose of chemical x will result in a 10% increase in mortality. In the first approach the question integrates multiple processes. In the latter approach, the larger question is broken down to more elemental questions.

In general, analysts try to present the questions at a level of aggregation (or disaggregation) for which the experts are familiar and comfortable. One might expect that decomposing complex processes to obtain uncertain quantities would aid experts in making judgments on more familiar quantities. However, Morgan and Henrion (1990) report that results are mixed as to whether decomposition actually improves outcomes. Thus, the extent to which a particular elicitation uses a more or less aggregated approach is a matter of professional judgment. Early interaction between the analysts structuring the assessment and substantive experts is encouraged. This can help to guide decisions on the extent of aggregation that is most appropriate for a given elicitation. Additional research in this area could help inform and guide analysts on the appropriate level of aggregation for EE projects.

5.2.3 What are the Staffing Requirements?

An EE project requires three types of staffing. First, it needs *generalists* who are familiar with the overall problem(s) or question(s) of interest and who are responsible for the general management of the EE. The second staffing need is for *analysts or facilitators* (often called “*normative experts*”) who are proficient in the design and conduct of EEs. These normative experts have training in probability theory, psychology, and decision analysis and are knowledgeable about the cognitive and motivational biases discussed in Chapter 3. They are proficient with methods to minimize these biases in an effort to obtain accurate expert

judgments. The third staffing need is for substantive *domain* or *subject matter* experts who are knowledgeable about the uncertain quantity of interest and any relevant theories and processes.

5.2.4 How are Experts Selected?

A critical element of an EE is the selection and invitation of the subject experts who will provide the required probabilistic and other judgments. The process for selecting these experts should ensure that the panel of experts will have all appropriate expertise to address the questions and will represent a balanced range of valid scientific opinions. Previous EE studies (e.g., Hawkins and Graham, 1988) have identified several additional criteria for the expert selection process, including: explicit and reproducible, reasonably cost-effective, and straightforward to execute. Transparency in the selection process is essential for an EE that will be used as part of a government science or regulatory activity.

5.2.4.1 What criteria can be used to select experts?

The selection of experts for an EE should use criteria that help to identify and choose experts who span the range of credible views. Because multiple disciplines may bring credible expertise to the EE, the selection criteria should seek to ensure that those disciplines are represented equitably (i.e., technical balance). Keeney and von Winterfeldt (1991) indicate that when an EE is intended to be an input to a highly influential assessment, the substantive experts “should be at the forefront of knowledge in their field, and they should be recognized as leaders by their peers.” As cautioned in Chapter 3 by Shackle (1972b?), EE study designers should be mindful that, while it is ideal to obtain experts representing the entire range of opinions in their field, this might be difficult to achieve in practice.

The selection criteria for EE experts may also vary depending on the objectives of the assessment. For example, is the goal of the assessment to characterize the range of credible views or to obtain a central tendency estimate?

Additionally, for some EEs it may be important that the selection criteria include institutional or stakeholder balance. The disclosure of institutional affiliation and attempts to achieve balance can help to avoid (or balance) conflicts of interest. There is relatively widespread agreement among EE practitioners that EE experts should fully disclose any real or potential conflicts of interest. However, if experts were always excluded from participation due to potential conflicts of interest, for some disciplines there might be few or no few experts available to participate. Since the goal of an EE assessment is to characterize uncertainties based on the most knowledgeable experts, it would be unwise to have a blanket ban on the participation of experts who might have a potential conflict of interest. It might be better to full disclosure of

any potential conflicts of interest and make a case-by-case determination about whether the nature of the conflict precludes an expert from participation.

The U.S. Nuclear Regulatory Commission (USNRC) in its Branch Technical Position (NUREG-1563) has set forth guidance for selecting an appropriate set of experts that may provide a good starting point for any future EPA guidance. It states that a panel selected to participate in an EE should include individuals who: “(a) possess the necessary knowledge and expertise; (b) have demonstrated their ability to apply their knowledge and expertise; (c) represent a broad diversity of independent opinion and approaches for addressing the topic; (d) are willing to be identified publicly with their judgments; and (e) are willing to identify, for the record, any potential conflicts of interest” (USNRC, 1996). Much of this guidance resembles that contained in EPA’s Peer Review Handbook (USEPA, 2006). While the names and institutional affiliation of the experts participating in an EE should be identified, it is not always desirable, or necessary, to attribute the particular elicited judgments to their respective experts. This follows the pattern of many peer review activities (e.g., SAB or NAS) that identify the members of the review panel; but, do not attribute particular views to their respective experts. The issue of anonymity with respect to individual judgments in the context of EE is discussed further in section 5.4.3 of this Chapter.

5.2.4.2 What approaches are available for nominating and selecting experts?

A number of approaches have been cited in the literature for nominating and selecting the substantive experts. Some approaches use literature counts as a rough measure of expertise. Others use participation on relevant NAS or SAB committees as a proxy for expertise. Another approach is to ask scientists who have published in the area of interest to recommend experts for participation. It may also be helpful to ask professional and academic societies, or other institutions that do not have a direct stake in the EE’s outcome, to submit nominations. In practice, it is possible to use a combination of these approaches. For example, in the pilot PM mortality EE, literature counts were used to identify which individuals should be asked to nominate experts for potential participation. In the full PM mortality EE, this process was modified further to include nominations from the non-profit Health Effects Institute. This allowed the pool of experts to include additional expertise in toxicology and human clinical studies that were not represented adequately by the initial approach. Having a carefully designed approach for nominating and selecting experts is advantageous because the entire process can be reproduced. This is desirable when there is a need to augment the number of experts in the assessment or to replicate the process for another study.

A currently unresolved issue is whether the sponsor of the EE (e.g., EPA) should be directly involved in the nomination and selection of experts or should allow a third party to

conduct this critical element of the EE process. This question presents a similar challenge for many EPA peer reviews. In cases where the assessment is influential or highly influential, peer reviewers are often picked by an outside independent party such as a contractor, SAB, or NAS. On one hand, the goal of minimizing involvement by the sponsor has the benefit of greater objectivity. Hawkins and Graham (1990) advocate that the selection process should minimize the level of control of the researcher who is conducting the elicitation. However, on the other hand, EPA may want to have more active control on the selection process because it is ultimately responsible for assuring the quality and credibility of its EEs. As a default practice, the EE project team is encouraged to give this issue careful consideration.

For highly influential EEs that are likely to attract controversy, it may be worthwhile to adopt a process or taking additional steps to help establish that the selection of experts was done carefully to represent the range of credible viewpoints. One possible approach is to use a transparent process with public input similar to the nomination process used to form new SAB panels. The SAB process allows outside groups to participate in the expert selection process by accepting nominations from the public for consideration. This approach has the disadvantage of the extra time involved and may raise concerns that affected stakeholders will try to influence an outcome by skewing the composition of the expert panel towards their viewpoint. Information about the design of peer review panels and selection of participants is presented EPA's Peer Review Handbook (3rd Edition, USEPA, 2006).

5.2.4.3 How Many Experts are Needed?

The number of experts involved in an EE is determined primarily by time and financial constraints, availability of credible experts, and the number of institutional affiliations or perspectives that one wishes to represent. There have been only limited efforts to develop mathematical theory for optimizing the number of experts used in studies (Hogarth, 1978; Clemen and Winkler, 1985; and Hora, 2004). Such theoretical approaches are modified frequently by other considerations including financial constraints.

It is possible to get a rough idea of the number of experts that are needed by looking at the number of experts that have been used in past EEs. A recent informal survey (Walker, 2004), based on 38 studies, found that almost 90% of the studies employed 11 or fewer experts. Nearly 60% of the studies relied on 6-8 experts and the largest number of experts used in any of these studies was 24. This survey is not intended to be representative of all EE studies; but, can provide some insight. Of the 38 studies in this survey, 27 were from a database provided by Roger Cooke from his work while at the University of Delft, Netherlands. The remaining 11 studies were obtained from a literature search. All of the studies elicited probability distributions or confidence intervals to describe uncertainty.

Clemen and Winkler (1985) argue that there can be diminishing marginal returns for including additional experts in an EE assessment. Their observations are based on a number of theoretical examples. The simplest example evaluated the impact of dependence between experts on the equivalent numbers of experts to achieve a particular level of precision in an estimate. Their findings show that when the experts are completely independent, $\rho=0$, the number of equivalent and actual experts are the same. The practical implication is that, the more different the experts are, the more experts are needed.¹² As dependence between the experts increases, the value of additional experts drops off markedly.

Clemen (1989) describes Hogarth (1978) as using test theory as a basis for discussing the selection of experts. His conclusions were that between six and twenty different forecasters should be consulted. Furthermore, the more the forecasters differed, the more experts that should be included in the combination. Libby and Blashfield (1978), though, reported that the majority of the improvement in accuracy was achieved with the combination of the first two or three forecasts. Steve Hora has argued often that “three and seldom more than six” experts are sufficient. Clemen and Winkler (1985) suggests that five experts are usually sufficient to cover most of the expertise and breadth of opinion on a given issue.

If an EE seeks to not only characterize the range of judgments; but, also to provide an estimate of the central tendency among the overall scientific community, then it may be necessary to include more experts in the process. In addition, it may be necessary to develop additional procedures to address questions about the representativeness of the group of experts. One suggestion has been to follow an EE with a survey of the broader scientific community that is knowledgeable about the issue of concern, combined with appropriate statistical techniques such as factor analysis. This will allow the analyst to compare the judgments from the EE experts with the views of this broader scientific community. To more fully develop and demonstrate this approach requires further research, development, and evaluation.

The requirements of the Paperwork Reduction Act (PRA) are an additional consideration for EEs that are sponsored by EPA or another federal government agency. OMB has indicated that EE activities falls under the requirements of the PRA. Therefore, the PRA stipulates that if more than nine individuals participate in a survey, the sponsoring agency must submit an information collection request to OMB under the PRA. This effort can add substantial amounts of time and cost to the completion of the EE. The administrative requirements of the PRA may by themselves be reason enough for EPA to limit the number of experts involved in an EE. EPA

¹² The value of ρ is the correlation and can range from 0, indicating there is no correlation and, thus in this case the experts are completely independent to 1, indicating that the experts are completely correlated or dependent.

may wish to pursue discussions with OMB about the PRA to better clarify how its requirements apply to EEs and to potentially minimize the impacts of the PRA on the EE process.

5.2.5 What is an EE Protocol?

The development of an EE protocol is one of the most resource intensive steps in the conduct of an EE. This step is particularly demanding when dealing with a complicated issue that is informed by different perspectives and disciplines. An EE protocol serves several purposes and contains the following:

1. Overall issue of interest and any relevant background information;
2. Motivation or purpose for the assessment and, at least in a general sense, the role of the EE in any larger modeling or decision process;
3. Description of the quantitative question of interest and definition of any conditions or assumptions that the experts should keep in mind; and
4. Information about heuristics and biases that are characteristic of EEs (see section 3.5.5) and provide guidance on how to minimize these problems.¹³

As an example, the protocol that was developed recently for EPA's pilot PM_{2.5} mortality EE project is available (IEC, 2004).

5.2.6 Why Should Workshops Be Considered?

Pre-elicitation workshops that bring the project team and experts together are a helpful, though not required, element of many EEs. There are three major reasons why holding one or more pre-elicitation workshops is advisable. First, these pre-elicitation workshops can be used to share information on the technical topics of the EE. This will help to assure that the experts are all familiar with the relevant literature, different perspectives about the research, and how these views might relate to the EE's questions. This can alleviate some criticism regarding the importance of a common body of knowledge (discussed in Chapter 3) and could ultimately reduce the need for more experts in future EEs on similar topics.

Second, feedback on the draft protocol that is obtained at a pre-elicitation workshop can be used to refine or restructure the EE's question. Finally, the pre-elicitation workshop can be used to introduce the concepts of judgmental probability, heuristics, and biases to the experts. The workshop provides an opportunity to conduct training so that the substantive experts will be familiar with the techniques used to elicit probability judgments. During the training, the project team and experts can practice the use of techniques to reduce bias (USNRC, 1996). This can

¹³ Appendix 5A-1 to this White Paper, "Factors to Keep in Mind When Making Probability Judgments" is an example of the type of material that should be included either in the protocol or briefing book.

allay the criticism (as discussed in Chapter 3) that experts who do not understand these techniques tend to give more haphazard responses (Hamm, 1991). Training is further discussed below in section 5.2.8.

5.2.7 What Is a Briefing Book and What Is Its Role in an Expert Elicitation?

Another important pre-elicitation step is the development of a “briefing book.” This briefing book is a binder that includes the journal articles and other technical information relevant to the topic of the EE. If this background information is particularly extensive, it may be more practical to provide the information on a CD. When gathering this information, it is important to avoid bias by selecting representative papers from all valid perspectives. According to NUREG-1563, this background material should be selected so that “a full range of views is represented and the necessary data and information are provided in a uniform, balanced, and timely fashion to all subject-matter experts” (USNRC, 1996). The experts should be given the opportunity to add pertinent information to the briefing book, including unpublished data that they are willing to share.

5.2.8 What Type of Training Should Be Conducted?

Pre-elicitation training can facilitate EEs in several ways. According to USNRC recommendations, training subject matter experts prior to elicitation has the following benefits: “(a) familiarize them with the subject matter (including the necessary background information on why the elicitation is being performed and how the results will be used); (b) familiarize them with the elicitation process; (c) educate them in both uncertainty and probability encoding and the expression of their judgments, using subjective probability; (d) provide them practice in formally articulating their judgments as well as explicitly identifying their associated judgments and rationale; and (e) educate them with regard to possible biases that could be present and influence their judgments” (USNRC, 1996). Training helps to ensure that the expert judgments accurately represent the experts’ states of knowledge about the problem of interest. In addition, this training provides an opportunity to level the knowledge base among the experts and can help to clarify the problem definition.

5.2.9 What Is the Value of Pilot Testing?

After the draft protocol and briefing book are complete, a pilot test can provide valuable feedback on the quality of the protocol and help identify any obstacles. The objective of this step is to improve the clarity of the protocol and to determine whether the questions are framed appropriately. Ideally, pilot testing should be conducted with substantive experts who are not among the pool of experts that will participate in the actual EE. Pilot testing can include several experts; but, it is essential to pilot test the draft protocol with at least one person.

5.2.10 How Should an EE be Documented?

It is absolutely critical that all significant aspects of the steps listed above be documented clearly. This documentation is required to assure a transparent process as required under EPA's Information Quality Guidelines. In addition, clear documentation is essential to establishing the credibility of the results from an EE.

5.3 WHAT APPROACHES ARE USED TO CONDUCT EXPERT ELICITATIONS?

Three different approaches for conducting EEs have been demonstrated and documented (see pp.141-154 of Morgan and Henrion (1990) for a more detailed description). These approaches for eliciting expert judgment, often referred to as "probability encoding," include: (a) the approach used by Wallsten and Whitfield in EEs carried out for EPA's OAQPS (see Wallsten and Whitfield (1986) for an example); (b) the approach used by Stanford/SRI, pioneered by Howard, North, and Merkhoffer, and described in Spetzler and von Holstein (1975); and (c) the approach used by Morgan and his colleagues at Carnegie Mellon University (Morgan et al., 1984). While these approaches have some case-specific characteristics and other features that differ based on the tools chosen by the analyst, most EE practitioners agree about general principles that constitute good practice for the encoding process.

The encoding process is typically divided into five phases:

- **Motivating:** Rapport with the subject is established and possible motivational biases are explored.
- **Structuring:** The structure of the uncertain quantity is defined.
- **Conditioning:** The expert is conditioned to think fundamentally about judgments and to avoid cognitive bias.
- **Encoding:** This judgment is quantified probabilistically.
- **Verifying:** The responses obtained from the encoding session are checked for internal consistency.

The encoding session is conducted in a private setting (e.g., typically the expert's office) so that the subject is comfortable and the discussion can be uninterrupted and candid. As discussed previously (section 5.2.5), the EE protocol is used to guide the encoding session so that the topics covered and responses to experts' questions asked are treated consistently among the several experts. Responses and other feedback from the subject matter experts are documented thoroughly with one or more of the following: written notes, transcripts, and audio or video tape.

5.3.1 What Are the Staffing Requirements for the Encoding Session?

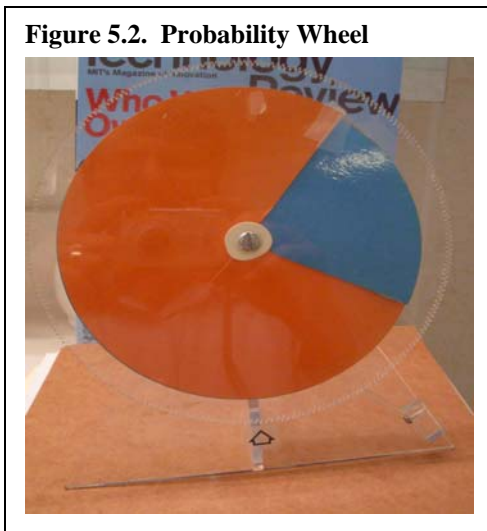
In general, a minimum of two individuals are required to conduct the encoding session. These usually include at least one subject matter expert and one analyst (see Section 5.2.3 for a description of the roles of these individuals). In addition, the project team requires a generalist who will be responsible for the general management of the EE.

5.3.2 What Methods And Tools Are Available To Aid An Elicitation?

A variety of innovative methods and tools to aid the elicitation of probabilities during the encoding session have been developed and used by Morgan and Henrion (1990), Spetzler and Staehl von Holstein (1975), Wallsten and Whitfield (1986), and Hamm (1991). These methods include: 1) fixed values, 2) fixed probabilities, 3) probability wheel, and 4) specialized software or models for incorporating and/or displaying judgments for feedback to experts (Morgan and Henrion, 1990; Hamm, 1991).

- **Fixed value methods:** In this method the probability that the quantity of interest lies within a specified range of values is assessed. Generally, this involves dividing up the range of the variable into equal intervals.
 - **Fixed probability methods:** In this method the values of the quantity that bound specified fractiles or confidence intervals are assessed. Typically, the fractiles that are assessed include the median (0.5), quartiles (0.25, 0.75), octiles (0.125, 0.875) and extremes such as (0.1, 0.99).
 - **Specialized software or models:** Software programs (e.g., Analytica®) that have been developed to facilitate the elicitation of probability judgments from experts. These software tools use graphics to illustrate the implications of an expert's judgments for the shape of the probability distribution provided. In addition, software (e.g., Excalibur) has been developed to address the calibration or performance assessment of experts.
 - **Probability wheel and other physical aids:** The probability wheel (Figure 5-2) is among the most popular physical aids used for EEs. It supports the visualization of probabilities by using a colored wheel (e.g., blue) with a spinner and pointer that has an adjustable pie shaped wedge of a different color (e.g., orange). To use the probability wheel, the expert varies the pie shaped wedge, changing the blue/orange proportions of the wheel until the probability that the spinner will end up on blue equals the probability of occurrence for the event of interest. The back of the probability wheel shows the fractional proportions of the two colors. The analyst
-

uses these numbers for the expert's response. See Morgan and Henrion (1990, pp. 126-127).



5.3.3 What Group Processes Can be Used for Expert Elicitation?

Many methods exist for assembling groups and obtaining their judgments. These methods vary in how they encourage interaction, allow iteration, and seek consensus. In addition, the application of each method depends on the number of group members, the time available, the financial resources, and the context of the participation. The Delphi method and Q-Methodology are described herein. Detailed discussions of many other methods are available elsewhere.¹⁴

5.3.3.1 Delphi

The Delphi method is a common technique for assembling a group, in a manner similar to focus groups, to obtain their judgments. One distinguishing feature of the Delphi method is that its members generally do not meet as a group. In a study by the Delphi method, participants are selected because they have expertise about a specific domain of interest. In most cases, correspondence is remote (i.e., mail, e-mail, fax, or telephone); however, face-to-face interviews can also be used. The initial interview session with each participant is conducted by a facilitator.

¹⁴For example, see SMARTe's (Sustainable Management Approaches and Revitalization -electronic) Tool for Public Participation (Go to: www.smarte.org → select: "Tools" → select: "Public Participation"), or: <http://www.smarte.org/smarte/tools/PublicParticipation/index.xml?mode=ui&topic=publicinvolvementaction>

The facilitator serves as a clearinghouse of the panelists responses. For subsequent iterations of the interview session, each participant sees and reacts to the views expressed by the other participants. Through a series of iterations of the interview, the panelists share and generate new ideas with the objective that consensus will emerge. Through this process, the Delphi method generates ideas and facilitates consensus among participants even though they are not in direct contact with each other.

The advantages of the Delphi method are that it encourages the sharing of ideas and promotes consensus with a large number of stakeholders who may be geographically distant from each other. It is a transparent and democratic technique that may be appropriate for highly technical issues where the goal is to obtain a consensus judgment. Its shortcomings are: 1) it may be resource intensive (time and money), 2) it can require large amounts of data to be assessed and distributed, 3) by emphasizing consensus, its final judgment may not characterize the full range of uncertainty, and 4) if participants are too rigid or their commitment wanes it may be impossible to obtain consensus judgments.

5.3.3.2 Q-Methodology

The Q-Methodology is a social science technique that was invented in 1935 by William Stephenson, a British physicist-psychologist (Stephenson 1935a; 1935b; 1953). It provides a quantitative basis for determining the subjective framing of an issue and identifies the statements that are most important to each discourse. This method provides a quantitative analysis of subjectivity, i.e., “subjective structures, attitudes, and perspectives from the standpoint of the person or persons being observed” (Brown, 1980; 1996; McKeown and Thomas, 1988).

In Q-Methodology, participants map their individual subjective preferences by rank-ordering a set of statements (typically on a Likert-like scale with endpoints usually representing “most agree” to “most disagree” and zero indicating neutrality). This is accomplished by sorting statements that are printed on cards using a traditional survey instrument and numerical scale; or, more recently, with an internet software application. The sorting is often performed according to a predetermined “quasi” normal distribution (i.e., the number of allowable responses for each value in the scale is predetermined with the greatest number of responses in the middle of the scale and fewer responses at either end). The collection of sortings from the participants form a kind of cognitive map, or mental model, of subjective preferences about the particular issue.

Following the sorting phase, participants are generally asked to reflect on their responses. Participants’ individual preferences are then correlated against each other and are factor analyzed. Factor rotation is commonly conducted either judgmentally, based on theoretical considerations, or by using varimax rotation. The factor outputs, as indicated by individual

loadings, represent the degree to which the study participants are similar in their responses. Factor scores represent the degree to which each statement characterizes the factor and can be used to construct a narrative that describes (subjectively) the specific discourse associated with the factor.

5.3.4 Use of Other Media to Elicit Judgments from Remote Locations

Though face-to-face interviews are often the preferred method for EEs, constraints on time and money may necessitate conducting the interviews via another medium. Whether questionnaires (by mail or e-mail) (Arnell et al, 2005), telephone, video-conference, or some combination (Stiber et al., 2004; Morgan, 2005) are used, the essential elements are consistency and reproducibility. To the degree possible, each expert should be presented with identical information in a standardized order. Answers to questions from the experts should be uniform and, ideally, should be prepared in advance for anticipated queries. It is important that the elicitation process produces a common experience for the experts so that their responses reflect a shared understanding of the questions.

Although face-to-face interviews are the standard and preferred approach, there are important advantages to other media. In addition to being less expensive and permitting greater flexibility with schedules, eliciting judgments remotely can engender greater consistency. Because there is less body language (e.g., hand movements and facial expressions), the emphasis is on the content of the interview and this can be more easily standardized for all of the interviews. In addition, if experts are responding to written surveys they may feel less pressured, and as a result, may be more thoughtful in their responses. Also, if follow-up questions become necessary, they can be handled with the same format as the original questions (e.g., telephone or written).

5.4 WHAT POST-ELICITATION ACTIVITIES SHOULD BE PERFORMED?

5.4.1 When are Post-Elicitation Workshops and/or Follow-up Encoding Appropriate?

Conducting a post-elicitation workshop is not required for every EE; however, the EE project team should consider and weigh its potential value added against the additional cost (i.e., resources and time). A post-elicitation workshop provides an opportunity for all of the subject matter experts to see the judgments from their peers. Typically, both the probabilistic judgments and the reasoning behind those judgments would be shared. At a workshop the experts and the project team can probe reasons for differences in judgments. The exchange of views may unearth new insights (i.e., new data, theory, or perspectives) that could influence experts to modify their judgments. This reconsideration and modification of judgments is consistent with the goal of obtaining the most accurate representation of the experts' beliefs based on their

understanding of the state of information. If resources or timing preclude an in-person post-elicitation workshop, an alternative is to meet via audio or video conference.

Holding a post-elicitation workshop for the experts to reflect and potentially change their views has two potential additional benefits. First, the experts can change and refine their responses so that the elicitation results more accurately represent their judgments. Second, where movement toward consensus is possible, it becomes more likely that uncertainty can be reduced by more closely representing the collective judgment of the relevant experts for that discipline.

5.4.2 Verify Final Judgments

As soon as practical after the elicitation, the subject matter experts should be provided with their elicitation results (USNRC, 1996). Then, the analysts can query the experts to ensure that the experts' responses have been represented accurately and even-handedly. It is the job of the project team to determine if any revision or clarification of the experts' judgments and rationale is needed. Any revisions should be obtain in a manner consistent with the original elicitation and documented carefully. Finally, the experts can confirm concurrence with their final judgments and associated qualitative discussions of rationale.

5.4.3 Should Individual Judgments be Combined, and if so, How?

In many decision-making contexts, decision makers want a single unambiguous result, not a complex spectrum of findings. When an EE uses multiple experts, which is most often the case, the result is many independent sets of judgments, each representing the beliefs of a single expert. These sets of judgments are the "experimental" results of the EE exercise. However, these data are very different from traditional experimental results. With traditional scientific experiments, if the process that produced the results, including measurement errors and uncertainties, is known, it may be possible to arithmetically combine the results into a single aggregate finding. Handling the results of an EE is more complex. Although each expert was prepared similarly for the elicitation (pre-elicitation workshop), was presented with the same data, and was elicited by essentially the same process (elicitation protocol), the experts differ in their training, experiences, and the manner of considering the relevant information to produce beliefs. Consequently, whether and how to combine multiple expert beliefs requires consideration of both theoretical and practical constraints and needs. In many cases, combining expert judgments may not be theoretically defensible or practical. See sections 3.4.3 and 3.4.4 for a more detailed discussion of the theoretical advantages and disadvantages of combining expert judgments.

Nevertheless, the practical nature of decision making sometimes motivates analysts to produce a single aggregate result. This section examines when it is appropriate to aggregate the judgments of multiple experts and describes how this can be done while preserving the individual nature of the EE data and maintaining the richness of the original findings. The decision to combine expert judgments and the selection of a method for doing so must consider the attributes of the particular elicitation and how the findings will be used. The application of a method to combine the judgments of multiple experts can be project-specific. This section focuses on the mechanics of expert aggregation. Section 6.3 provides a policy-related discussion of these issues.

Whether or not judgments are aggregated, there is a question about the desirability of associating each individual's judgments by name; or, whether it is sufficient to list the experts who participated and identify the individual judgments via an anonymous lettering system. This issue was raised in the USNRC's Branch Technical Position guidance (USNRC, 1996) and is still debated in the decision analysis community. Cooke (1991) takes the perspective that EE should be consistent with scientific principles and has argued that the goal of accountability requires each judgment to be explicitly associated with a named expert. Others consider EE to be a trans-scientific exercise. From this perspective, preserving anonymity for the specific judgments made by an expert best serves the overarching goal of obtaining the best possible representation of expert judgment. Given the current norms within the scientific community, experts may be unwilling to participate and share their judgments honestly if they fear a need to defend any judgments that divert from the mainstream or conflict with positions taken by their institutions.

5.4.3.1 Why are judgments different?

When considering the possible combination of expert judgments, the first step is to ask the question: "Why are judgments different?" Unlike a traditional scientific experiment in which the selection of a technique for combining results can (and should) be made before any data is collected, with an EE, it is necessary to see the results before determining if aggregation is appropriate.

Understanding the source of differences between experts can lead to insights, consensus, and/or revision of the elicitation protocol (Morgan and Henrion, 1990). Indeed, this understanding about the source of the differences can be more valuable than any aggregate finding. Furthermore, for many situations, variability among experts is not a problem; but rather, the objective of the elicitation. Many scientific questions are unanswerable; but, have a spectrum of legitimate approaches providing different answers. Expert elicitation is often undertaken to

obtain a feel for the range of potential answers. Hence, diversity of judgments may be a good thing and it would be inappropriate to replace this richness outcome with a crude average.

The experts' judgments may be different for a number of reasons, including: unused information, misunderstanding about the question, different paradigms among experts, and motivational biases. It is possible that some of the experts failed to use information that others found to be influential, or weighed evidence differently. Alternatively, it may be clear from the experts' responses that one of them misunderstood the question (or at least, understood it differently). In these cases, it may be possible to re-elicite the mistaken expert to rectify the irregularity. If the elicitation process is uniform, there may be less variance in responses.

In other cases, the differences in response may result from different paradigms by which the experts view the world and the data. This will be particularly true when the experts come from different disciplinary backgrounds. Experts tend to trust data obtained through methods with which they have direct experience. For example, when one is trying to estimate the relationship between exposure to a substance and increased morbidity or mortality, epidemiologists may tend to find epidemiological data compelling, while being more suspect of toxicological studies on animals. Toxicologists may have the opposite preference. In this situation, the variability among the findings represents a spectrum of beliefs and weights that experts from different fields place on the various types of evidence.

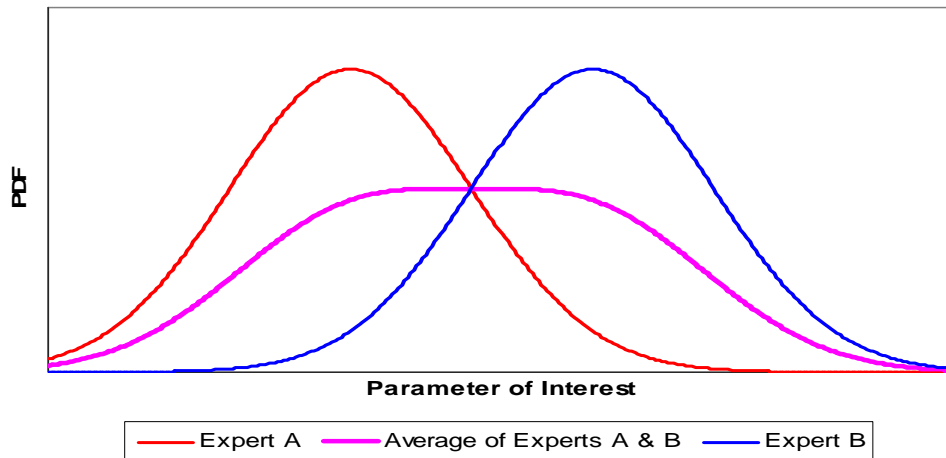
Although one of the goals in the selection of expert is to obtain an impartial panel of experts, there may be cases of motivational bias. Identifying such experts and determining how to use their beliefs is best handled on a case-specific basis.

Evaluating the source of differences among experts is intended to produce insights that may obviate any needs to aggregate. These insights may lead to improvements in the elicitation protocol, understanding about varying disciplinary perspectives, or ideas for future research that can (ideally) reduce the inter-expert differences. In any case, the knowledge gained from insights may be more valuable than the benefits of a single aggregate finding.

5.4.3.2 Should judgments be aggregated?

The next step is to determine if judgments should be aggregated; or, in less normative terms, if it is appropriate to aggregate. In many situations, part of the answer to this question depends on the relative value of the uncertainty of each individual's judgments with respect to the difference between the individual judgments. If the inter-individual variability is less than each individual's uncertainty (see Figure 5.3), then aggregation may be appropriate. In this case, knowledge sharing may result in convergence because aggregation may average out the "noise" in the characterization of the different judgments.

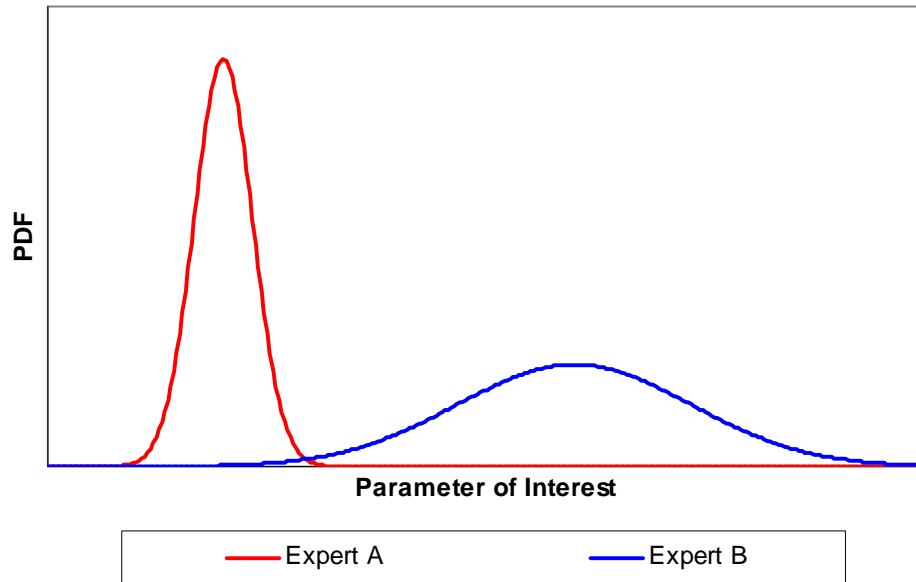
Figure 5.3. Experts with Similar Paradigm But Different Central Tendency



When the inverse is true, the inter-individual variability is greater than each individual's uncertainty (see Figure 5.4), then it may be inappropriate to aggregate. If it appears that the experts' judgments are based on fundamentally different paradigms, then aggregation may create an average set of judgments that lacks any phenomenological meaning. The dissonance between the experts, itself may be the insight that the EE provides. The existence of heterogeneity is itself important and should encourage alternative modes of policy analysis (Keith, 1996). If the decision maker expects to receive a single value outcome, and it does not appear appropriate based on the considerations of individual uncertainty and inter-individual variability, then the analyst should explain these limitations to the decision maker. If it is possible to combine the results to provide a meaningful characterization, these judgments may be aggregated. However, the analyst should exercise caution and try to share sensitivity analyses with the decision maker

In any situation when multiple judgments are combined, sensitivity analyses should be used to examine the effects of each expert on the aggregate outcome. Moreover, the presentation of results should include the individual responses along with the combined response. In general, the decision to aggregate is case-specific and "depends on the individual circumstances and what is meant to be accomplished" (Kraymer von Krauss et al., 2004).

Figure 5.4. Experts with Different Paradigms



5.4.3.3 How can beliefs be aggregated?

Once the decision has been made to aggregate expert judgments, several methods are available. There are mathematical methods of combination including simple averages, weighted averages, the Classical Model, and Bayesian approaches (Clemen and Winkler, 1998; Ayyub, 2000; Ayyub, 2001). Simple averages are the easiest to implement and often have good and equitable performance. Weighted averages can vary in sophistication and scheme for developing weighting factors. In the Classical Model *calibration* and *information* are used for a statistical test to create weights for averaging (Bedford and Cooke, 2001). Bayesian approaches are discussed in numerous references (Genest and Zidek, 1986; Jouini and Clemen, 1996; Stiber et al., 2004).

Behavioral methods offer a means by which group interaction can produce consensus, or at least generate a better understanding of differences between experts and their sources. Included in these behavioral methods are the Delphi method (Dalkey, 1969; Linstone & Turoff, 1975; Parenté and Anderson-Parenté, 1987) and the Nominal Group Technique (Delbecq et al.,

1975; and Kaplan, 1990). Although a group can lead to a consensus among experts, this is not always possible or desirable. Inter-expert variability is a source of richness in EE studies and often leads to insights. The analyst should not seek to remove these differences without first understanding them.

5.5 WHEN AND WHAT TYPE OF PEER REVIEW IS NEEDED FOR REVIEW OF EXPERT ELICITATION?

The purpose of conducting a peer review of an EE project is to evaluate whether the elicitation and other process elements were conducted in a professional and objective manner. The judgments provided by the experts are not subject to evaluation by peer review. The mechanism for peer review should be selected in consultation with EPA's *Peer Review Handbook* (2006) and in consideration of the intended use of the EE results. In some circumstances, it may also be appropriate to conduct a peer review to provide advice on how the results of an EE should be considered relative to other analyses or scientific assessments in a given regulatory context.

5.6 SUMMARY

This chapter discussed important factors to consider when conducting an "acceptable" EE. It presented "good" practices for conducting an EE whether it is conducted (or sponsored) by EPA or conducted by outside parties for submission to EPA. The discussion of "good" or "acceptable" practice was based on a review of the literature and actual experience within EPA and other federal agencies. In general, the degree to which practices are "good" or "acceptable" depends substantively on the following: 1) clear problem definition, 2) appropriate structuring of the problem, 3) appropriate staffing to conduct EE and selection of experts, 4) protocol development and training, including the consideration of group processes and methods to combine judgment, if appropriate, 5) procedures to verify expert judgments, 6) clear and transparent documentation, and 7) appropriate peer review for the situation. While this White Paper presents what the EE Task Force believes can constitute "good" practice for EPA EEs, we recognize that a range of approaches are currently used among EE practitioners. Hence, the design of a particular EE involves considerable professional judgment.

6.0 HOW SHOULD RESULTS BE PRESENTED AND USED?

6.1 DOES THE PRESENTATION OF RESULTS MATTER?

The presentation of results not only provides the findings of an EE, it also serves as a window to the methods, assumptions, and context of the EE itself. For many stakeholders, viewing the results will be their first impression of the EE. Hence, it is important that results are presented thoughtfully and with attention to the particular needs of the intended audience. At the EPA, it is expected that EE results will be used to support regulatory decision making. In this context, the decision maker will be one of the primary recipients of EE findings. The presentation of results will inform decision makers about EE findings and help them to consider the EE results along with many other sources of information. To this end, results should be presented in ways that will enable their understanding and promote their appropriate use. In addition, regulatory decisions require a transparent process and therefore EE results should also be presented in a manner that will enhance their understanding by members of the public.

6.2 WHAT IS THE STAKEHOLDER AND PARTNER COMMUNICATION PROCESS?

Communicating about an EE is a two-way process involving the transfer of information between parties with different professional training, data needs, motivations, and paradigms for interpreting results. Hence EE analysts must consider how to present results for each intended audience. Table 6.1 summarizes potential stakeholders for EPA EEs. Presentation of EE results should consider the stakeholder perspectives, knowledge of the EE subject matter, and the context in which the EE is being developed and used.

Table 6.1. List of Stakeholders and Partners for EPA Expert Elicitations

EPA risk managers/assessors

Members of the public

State and local environmental and health agencies

Federal agencies (e.g., HHS, USGS, DOI, etc.)

The Office of Management and Budget

Tribal governments

Regulated community

Scientific community

Managers of federal facilities (e.g., DOD, DOE, etc.)

(After: USEPA, 2001a)

Communication of technical information to the public involves the establishment of “trust and credibility” between the Agency and the Stakeholders. Approaches to establishing trust and credibility include continued community involvement (USEPA, 1994) and providing information at the appropriate level to aide an understanding of results (USEPA, 2002). A wide range of literature is available regarding stakeholder perceptions of risk (Slovic et al., 1979) and the need to consider these perceptions in developing communication strategies for EE results. For example, it may be helpful to discuss the EE results within the context of other studies or regulation decisions.

EPA has developed several guidance documents on public involvement and communication that may be considered when developing EE communication strategies (USEPA, 1995, 1997, 2000a) along with the large research literature (Covello, 1987; Deisler, 1988; Fischhoff, 1995, 1997 and 1998; Hora, 1992; Ibrekk and Morgan 1987; Johnson and Slovic 1995; Kaplan 1992; Morgan et al., 1992; Ohanian et al., 1997; and Thompson and Bloom, 2000). Depending on the complexity of information presented in the EE, it may be necessary to test and review communication materials to improve clarity of presentation.

6.3 HOW CAN COMMUNICATIONS BE STAKEHOLDER-SPECIFIC?

Each stakeholder has different needs that should influence the content and form of communications products. The types of communication and information products that would be appropriate for the major types of stakeholders are listed below.

Risk managers and state/federal officials may have technical and/or policy training; but, probably have a limited understanding of EE. Useful products for this audience include: executive summaries, bulleted slides, and briefing packages with a short description of EE process. When developing these communication materials, one may consider placing the results of the EE in the broader context of the decision. This may include an overview of the reasons and importance of the decision, previous decisions, and the positions of major stakeholders. The presentation should also include an appropriate discussion of the uncertainties and their potential impact on the decision (Thompson and Bloom, 2000).

Researchers with technical domain knowledge and/or expertise in expert elicitation will be the most literate about EEs. This audience typically has a more in-depth scientific knowledge of the EE process and related issues. Useful products include technical reports, peer-reviewed scientific papers, and presentations at professional meetings.

Community members generally have a limited knowledge of EE and may require more background and a concise discussion of the EE process. Documents may include fact sheets, press releases, slide shows, and speeches that summarizes the key issues and conclusions of the

EE in a lay person's context. Synopsis and simplification do not mean simplistic products. Consideration of the user's knowledge base is important.

6.4 WHAT IS IN A TECHNICAL SUPPORT DOCUMENT?

The Technical Support Document (TSD) provides the basis for development of all information products. The TSD contains all relevant information for the EE, including background, methods, data, analysis, and conclusions. Appendices to the TSD may include the EE protocol, prepared questions and answers for the interviews, list of experts with their affiliation, and other information documenting the process.

The following sections provide a checklist of suggested topics for the TSD of an EE. It covers the introductory information that should be included in documentation (Section 6.4.1), the technical details of the EE (Section 6.4.2), and finally, examples of means to summarize the results of the EE (Section 6.5). This template for information that should be covered does not preclude other relevant requirements for data quality or peer review. Most of the items in this checklist are self-explanatory and do not require a detailed description. Descriptions are only included when they may be helpful.

6.4.1 What is in the Introduction of a TSD?

The Introduction of a TSD should include:

- The data or information gap(s) that this EE addresses. The quantities that are the subject of the elicitation may be poorly defined or variously conceptualized (Hora, 2002) and should be clarified. A well-crafted objective statement should specify the problem that the EE addresses, define the meaning of elicited values, and identify the intended audiences of the findings.
- A brief summary of what EE is, and what it is not (especially compared to “expert judgment” and “peer review”).
- The rationale for using EE in the context of the larger policy or research question. Describe the policy/research question to which these elicitation results will be applied.
- How the methods used in this instance compare to “gold standard” methods.
- What the results of EE mean.
- What are the limitations/cautions of the elicitation?

When the documentation uses words to describe uncertainty, the audience may have varying understandings of what is meant. Hence, the authors need to be sensitive that even

quantitatively presented “probabilities” or “likelihoods” are often misunderstood. Research has shown (Anderson, 1998; Edwards, 2005; Morgan, 1990) that these ambiguous terms can create confusion or misunderstanding for the decision-maker about what presented results may mean. The Technical Support Document should describe what is meant in the EE study, by “probability” or “likelihood.” People tend to reason differently about established frequencies and probabilities of unique future events.

Anderson (1998) provides a useful summary of cognitive research that demonstrates that different people may interpret the word “probability” in very different ways. She abbreviates the concepts as in Table 6.2, and asserts that people use different solution algorithms or heuristics to take meaning from a provided probability, depending on which definition of “probability” they are using. Likewise, she stresses that it is important for an audience to know how probability was understood by the experts.

The research on how this affects elicited information implies that the results of the EE must be presented with attention to the format of results and to encourage the audience to use the results correctly. In light of these challenges, Anderson recommends formatting and presentation that will be discussed below. Research by cognitive scientists indicates that a presenter must take care in the introduction of a presentation to define what is meant by “probability” or “likelihood,” and that this should correspond to the concept held by the experts in the elicitation.

Table 6.2. Classification of Subjective Concepts Associated with Probability

For most ecologists and statisticians, the word “probability” seems to have a clear meaning. However, cognitive scientists recognize that subjective meanings vary depending on context. Teigen (1994) classified several ideas associated with probability and uncertainty. Each of the subjective concepts implies its own calculus of “probability” and each seems to be processed by a different cognitive mechanism (Teigen, 1994). It is important for Bayesian analysts to realize which idea they are activating when they refer to “probability” in a paper or ask an expert for a probability estimate.

Concept	Definition of “Probability”
<i>Chance</i>	The frequency of a particular outcome among all outcomes of a truly random process.
<i>Tendency</i>	The tendency of a particular outcome to occur, or how “close” it is to occurring.
<i>Knowledge</i>	It is allocated among the set of known hypotheses.
<i>Confidence</i>	The degree of belief in a particular hypothesis.
<i>Control</i>	The degree of control over particular outcomes.
<i>Plausibility</i>	The believability, quantity, and quality of detail in a narrative or model.

Adapted by Anderson (1998) from Teigen (1994)

6.4.2 What Technical Details of the Expert Elicitation Methods are in the TSD?

The previous section covered what might be included in the introduction for an EE’s TSD or other EE communication products. This section addresses what technical details of the EE methods should be included in the body of the document. The documentation should cover:

- The process (protocol) used for the EE and reasons for selection of key elements of the process.
- What criteria were used in selecting the experts (both criteria for individuals such as type of expertise and overall group criteria such as representing the range of credible viewpoints or range of disciplines. It should also identify who selected the experts.
- How well the set of experts selected meets the criteria set forth for the elicitation:
 - Identification of the list of experts, with affiliation and discipline/field that were selected and who agreed to participate and who, if any, did not.
 - Any potential conflicts of interest concerns (or appearances of conflict) and, if any, how they were addressed.
- Clear characterization of what experts were asked:

- Which uncertainties, parameters, relationships, etc. the experts addressed.
 - The information elicited from the experts.
 - What data experts used on which they may have based their judgments, as well as identified key data gaps. Presenting these alongside the EE results, however, might misleadingly imply an “apples to apples” comparison. The EE may address a broader question, or may introduce complexities that cannot be analyzed with available data. Or, if the EE includes a wider range of sources of uncertainty, one would expect the uncertainty bounds to be wider.
 - The degree to which the elicited results conform to axioms of probability theory and to the available empirical data.
 - Where biases may have been introduced and, if possible, insights into the likely direction of any biases (individually and overall).
 - How well the extreme values are likely to be represented (especially if there are potential catastrophic effects).
 - Possible correlations with non-elicited components of the overall analysis or policy question.
 - Text or graphics (e.g., influence diagrams or frequency-solution diagrams) that describe the mental models of the experts.
 - Presentation of results.
 - Findings of uncertainty and sensitivity analysis, including sensitivity of results to different methods of aggregating expert judgments from the elicitation.
 - Insights/explanations for differences in judgments among experts:
 - Degree of consensus or disagreement.
 - Degree to which views changed from initial judgments to final judgments – how much exchange and clarification of definitions and issues helped to resolve differences.
 - Principle technical reasons for difference in views, especially for outlying judgments. These may reflect different conceptual models, functional relationships, or beliefs about the appropriateness of evidence or parameters. This qualitative explanation is an important complement to the quantitative presentation of results (Morgan 2005).
 - Whether this is the first EE on this parameter, or whether there is a history or evolution of judgments that would be helpful to understand.
-

- Remaining uncertainties and weaknesses – Possible future strategies (e.g. data development) to reduce important uncertainties or to eliminate possible biases.
- Summarize the any peer review comments and what was done (and not done) to address them, including preliminary peer review conducted on methods.

Table 6.3 summarizes the technical details that should be included in the TSD.

Table 6.3. Summary of Technical Details of Expert Elicitation Methods

EE Process	Key Elements	Additional Data
Process	Description of EE process Description of reasons for elicitation and elements	Appendix to Technical Report
Expert Selection	Expertise requirements Range of affiliations and disciplines/fields Comparison of experts and how they met criteria	Appendix with criteria and basis List of Experts in Appendix Criteria for determining potential conflicts of interest
EE Questions	Charge Questions summarized Definitions of uncertainties, parameters, relationships, etc. experts addressed Definition of information elicited from experts Definition of data gaps Discussion of why data was aggregated or not aggregated	Appendix with detailed questions and supporting information
EE Results	Raw data tables and post-processed results that apply elicited values to calculate relevant quantities.	Appendix with all elicited probabilities and qualitative responses.
EE Analysis	Comparison of how elicited results conform to axioms of probability theory and to empirical data Biases introduced? Direction of Biases Extreme value presentation Correlations with non-elicited components of overall analysis or policy questions Mental models	Appendix with detailed analysis of calculations, graphics, etc. Calculations and detailed analyses in Appendix Appendix with influence diagrams and frequency-solution diagrams
EE Conclusions	Insights/explanations for differences in judgments among experts Degree of consensus and disagreement Analysis of changes in views from initial judgments to final judgments (how exchange and clarification of definitions and issues resolved differences). Technical reasons for differences in views and outlying judgments Results in context – history or evolution that would be helpful to understand	Appendix may include dissenting opinions if appropriate
Uncertainties and Weaknesses	Future strategies (e.g., develop data to reduce uncertainties and eliminate possible biases).	

6.5 WHAT ARE EXAMPLES OF EFFECTIVE EXPERT ELICITATION COMMUNICATIONS?

Many alternatives are available for conveying results to the users of EE findings. The following section provides examples of how results of EEs may be displayed qualitatively and quantitatively in text, figures, and tables. These examples are intended to demonstrate effective communication; but, the suitability of any presentation depends on the particular results, context, and audience. These different presentations of results require a variety of software and range of expertise to create and interpret. When considering among different displays, the analysts should consider the technical level of the audience and the aspect of the results that are to be highlighted. This section also identifies where research suggests that particular means of communicating results may be more effective.

6.5.1 How can Probabilistic Descriptions be Used?

If qualitative terms (e.g., “likely” and “probably”) are used, they should be associated with their quantitative meaning in the EE. Wallsten et al. (1986) and other researchers have demonstrated that the quantitative probability associated with a term of likelihood varies substantially from person to person. To overcome that inter-individual variability, some researchers have proposed systematizing the uses of specific terms (Moss and Schneider, 2000; Karelitz et al., 2002). For example, the terminology system of the Intergovernmental Panel on Climate Change (IPCC, 2005) is shown in Tables 6.4 and 6.5. Table 6.4 shows quantitatively calibrated levels of confidence. These can be used to characterize uncertainty that is based on expert judgment as to the correctness of a model, an analysis, or a statement.

Table 6.4. Quantitatively Calibrated Levels of Confidence

Terminology	Degree of confidence in being correct
Very High confidence	At least 9 out of 10 chance of being correct
High confidence	About 8 out of 10 chance
Medium confidence	About 5 out of 10 chance
Low confidence	About 2 out of 10 chance
Very Low confidence	Less than 1 out of 10

Table 6.5 shows a likelihood scale. This refers to a probabilistic assessment of some well defined outcome having occurred or occurring in the future – fuzzy boundaries.

Table 6.5. Likelihood Scale

Terminology	Probability of Occurrence/Outcome
Virtually certain	> 99%
Very likely	> 90%
Likely	> 66%
About as likely as not	33 to 66%
Unlikely	< 33%
Very unlikely	< 10%
Exceptionally unlikely	< 1%

6.5.2 How can Text Summaries be Used?

To improve the audience’s understanding of results, Anderson et al. (1998) recommend that results should be presented:

- As a frequency (as 43 out of 10,000, rather than 0.0043);
- Within a well-defined “reference class,” such as the general population to which the frequency might apply;¹⁵
- Keeping constant the denominator of the frequency statement (i.e., the size of the population) constant across comparisons (such as 43 out of 10,000 and 2082 out of 10,000, rather than 43 out of 10,000 and 2 out of 10).

Anderson and other researchers have concerned themselves with the ways in which humans receive and process information, specifically with respect to uncertainties, and conclude that humans have difficulty in interpreting probabilities expressed as decimals between 0.0 and 1.0. They note that frequencies are taught early in elementary school mathematics, being part of set theory, classification and counting, while probabilities generally have not been taught until advanced math in high schools or universities. Consequently, the heuristics needed for an audience to correctly interpret results from an EE are more easily available, to most people, when presented as frequencies.

From a survey of summaries of EEs, it appears that few researchers provide simple, quantitative summaries of results. Instead the summaries rely on graphical or tabular presentations. Policy-makers and others may have difficulty reading and understand these technical presentations. Three examples of effective textual summaries of EE results are provided in the following box:

¹⁵ Modified slightly from Anderson’s example: “If there were 100 similar populations of Spectacled Eiders nesting in eastern arctic Russia, how many would you expect to exhibit a rate of population increase of less than -0.05?”

IEI (2004): “...the experts exhibited considerable variation in both the median values they reported and in the spread of uncertainty about the median. In response to the question concerning the effects of changes in long-term exposures to PM_{2.5}, the median value ranged from values at or near zero to a 0.7 percent increase in annual non-accidental mortality per 1 µg/m³ increase in annual average PM_{2.5} concentration. The variation in the experts’ responses regarding the effects of long-term exposures largely reflects differences in their views about the degree of uncertainty inherent in key epidemiological results from long-term cohort studies, the likelihood of a causal relationship, and the shape of the concentration-response (C-R) function.”

Morgan and Keith (1995): “Of the 13 responses received, 4 of the means lie in the interval -1 to ≤ 0 and 9 lie in the interval -2 to ≤ -1 .”

Titus and Narayanan (1995): “Global warming is most likely to raise sea level 15 cm by the year 2050 and 34 cm by the year 2100. There is also a 10 percent chance that climate change will contribute 30 cm by 2050 and 65 cm by 2100. “

6.5.3 How can Figures be Used?

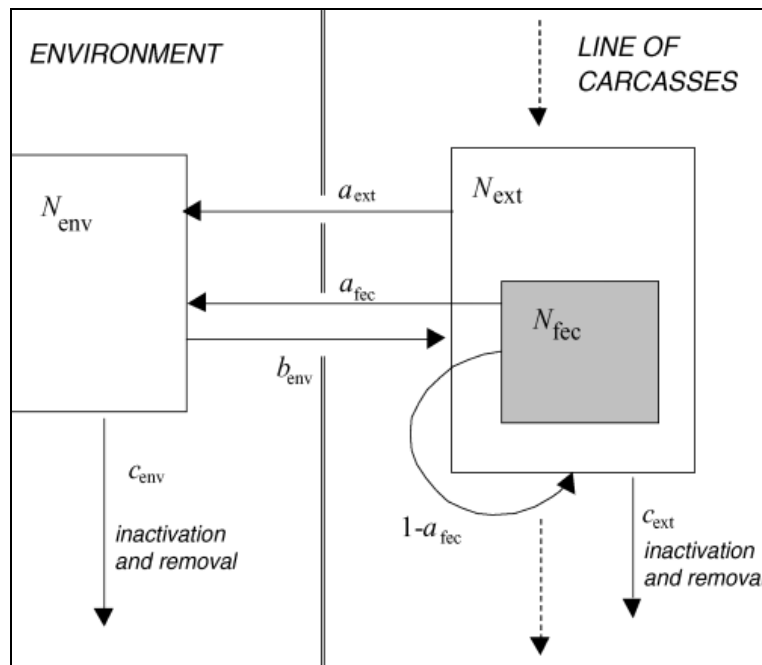
The audience for EE results requires a clear and unambiguous understanding of the elicitation questions. Hora and Jensen (2002) note that the attribute, parameter, or relationship that is the subject of an elicitation may itself generate debate among experts. “In a sense, this is logical: if there were nothing unclear about a quantity, it would probably not be selected for elicitation. The mere fact that it was chosen in the first place implies that it is critical in some sense, and perhaps the difficulties extend to its definition.” Because it is difficult to present uncertain quantities with text alone, diagrams and figures can lead to effective communication.

6.5.3.1 Influence Diagram

The use of “influence diagrams” is frequently used to illustrate the question, the important parameters and relationships that are understood to compose the elicited question. Providing the influence diagram used with the experts can be an instructive introduction for the audience and prepare them to put various results in proper context.

In Figure 6.1, Nauta et al. (2005) provide a useful influence diagram that captures the uncertain model parameters to be elicited within the graphic presentation. This type of presentation improves the clarity of the model and parameters by helping to facilitate both common conceptual models (or identification of differences) and the meaning of the parameters elicited.

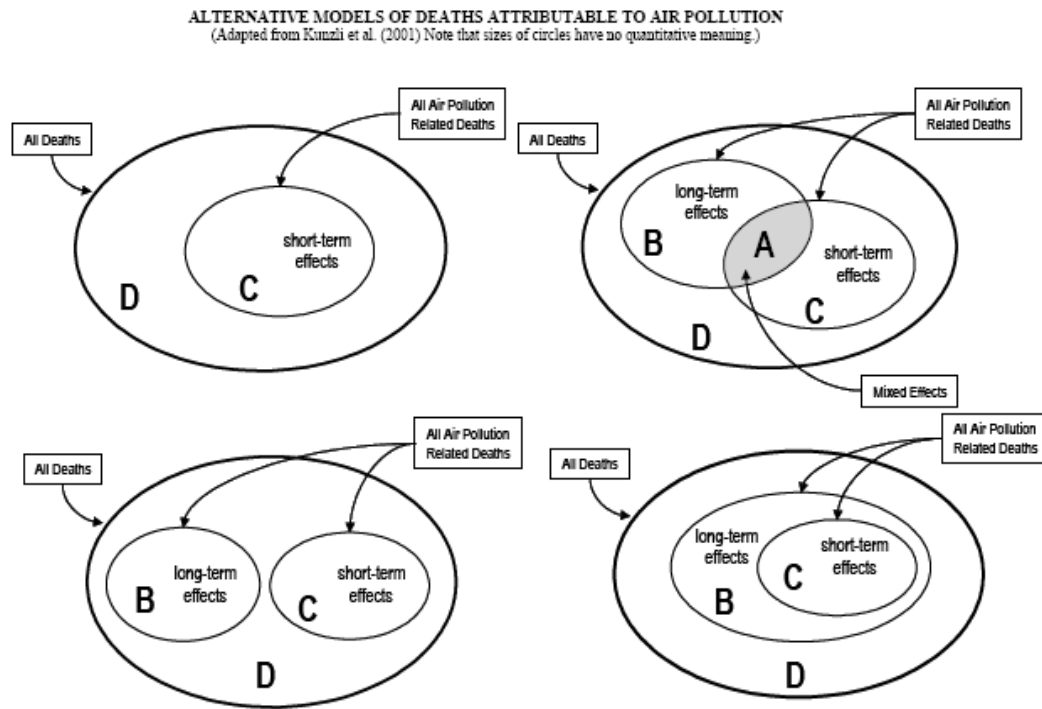
Figure 6.1. Example Influence Diagram



Source : Nauta et al. (2005)

In Figure 6.2, IEC (2004) presents a set of influence diagrams that were adapted from Kunzli et al. (2001). The original model, in the upper right of Figure 6.2, was adapted to illustrate variant conceptual frameworks (mental models) for describing the relationship among different causes of mortality.

Figure 6.2. Example of Alternative Mental Models Held of Different Experts

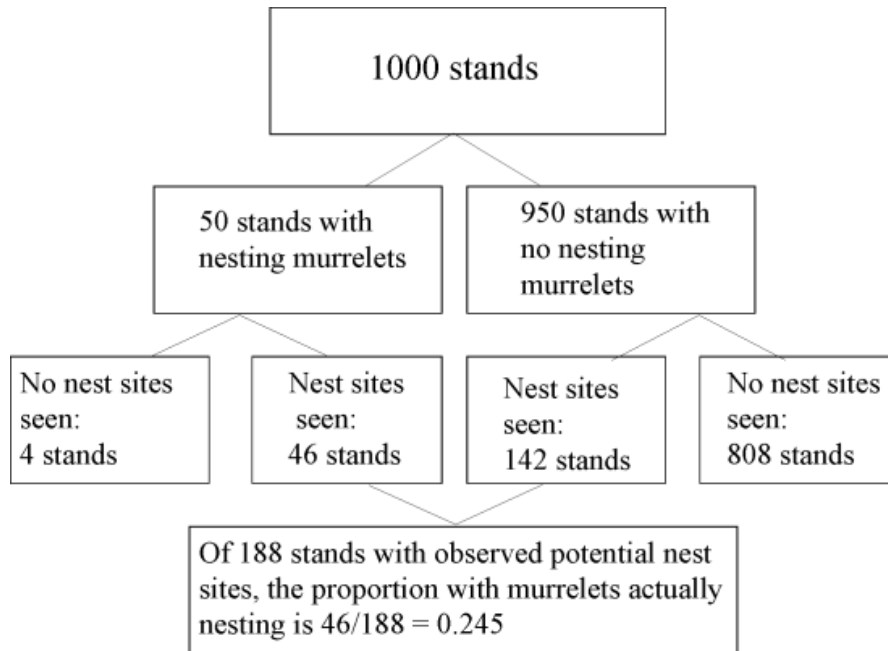


Source: Adapted from Kunzli et al. (2001) by IEC (2004)

6.5.3.2 Frequency-Solution Diagram

An example of a frequency-solution diagram, including hypothetical statements of frequencies and the reference population, taken from Anderson (1998), is provided in Figure 6.3. Anderson (1998) found that “presentation of the data in frequency format seems to encourage mental imagery and facilitate estimation of the correct answer.”

Figure 6.3. Example of a Frequency-Solution Diagram

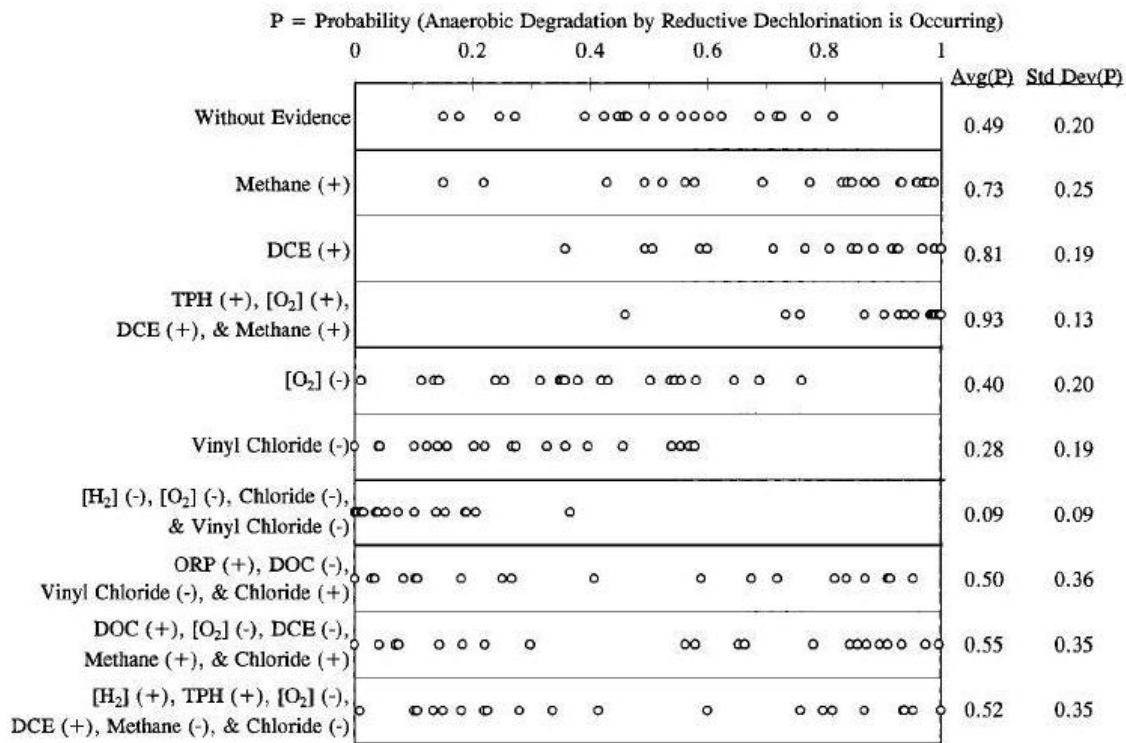


Source: Anderson (1998)

6.5.3.3 Results of Multiple Models Based on Experts' Probabilities

The probabilities obtained from experts may be used as the quantitative parameters of models. These multiple models (one for each expert) can be used to provide results under different scenarios of evidence. Figure 6.4 shows an example from Stiber et al. (1999) where the outputs of 22 models, built from the probabilities obtained from 22 experts, are compared for different scenarios (or cases of evidence).

Figure 6.4. Distribution of Expert Models' Predictions for Different Cases of Evidence



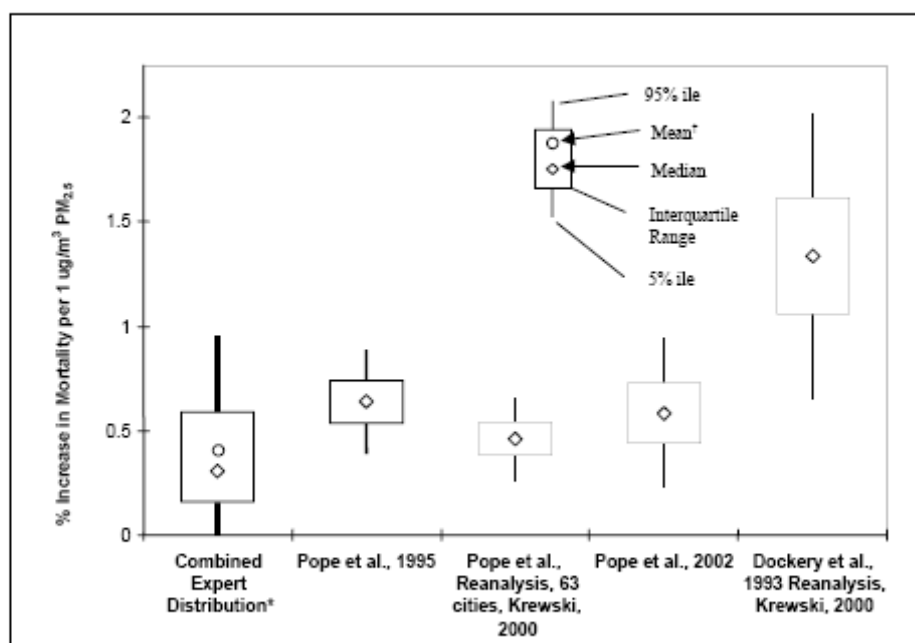
Source: (Stiber et al., 1999)

6.5.3.4 Box and Whisker Graphs

Sometimes experts are asked to specify a probability density function or a cumulative probability function. Below are several example ways in which such results may be presented. The “box and whisker” diagrams are perhaps the most intuitively understood of the formats. Figure 6.5 is a complex, but effective, box-and-whisker graphic showing multiple responses from the multiple experts and providing combined expert distribution.

Figure 6.5. Example Presentation of a Simple Box and Whisker Diagram Comparing Expert Elicitation Results with Other Studies

Comparison of Combined Expert Judgment Distribution to Results from Selected Studies: Percent Increase in Annual Non-Accidental Mortality Associated with a $1 \mu\text{g}/\text{m}^3$ Increase in Annual Average $\text{PM}_{2.5}$



*The experts' judgments were combined assuming equal weight to each expert and an underlying distribution of population weighted annual average $\text{PM}_{2.5}$ generated from the BENMAP model (see Analytical Methods section for details).

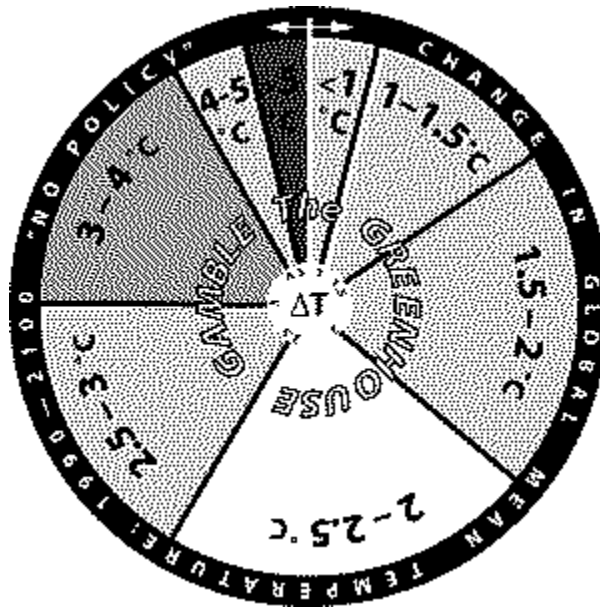
†No mean values were elicited from experts. All means were estimated using Monte Carlo sampling methods.

Source: IEC (2004)

6.5.3.5 Roulette or Probability Wheel

Because it is based on a familiar gambling paradigm, the roulette or probability wheel is a useful and intuitive display (Figure 6.6). People intuitively understand that the likelihood of any particular outcome is proportional to its area on probability wheel. This display provides an effective presentation of a future uncertain event because one of these states of the world will occur; but, we do not know which.

Figure 6.6. The Roulette or Probability Wheel to Express EE Results



Source: <http://web.mit.edu/globalchange/www/wheel.degC.html>

6.5.3.6 CDFs and PDFs or Both

In the Thompson and Bloom study (2000) of EPA decision-makers, the focus group liked “the format of showing risk as a distribution, although several members indicated a preference for seeing a cumulative distribution function instead of, or in addition to, a probability density function. They expressed some confusion about the level of aggregation of the distribution (i.e. whether it was representing variability in the distribution of all the maximum individual risks for the source category, or uncertainty for the maximum individual risks for one source). Most said that they would need more information about what the distribution represents and the underlying assumptions.”

The results of lead induced health effect from the elicitations of multiple experts can be easily compared with CDFs as in Figure 6.7.

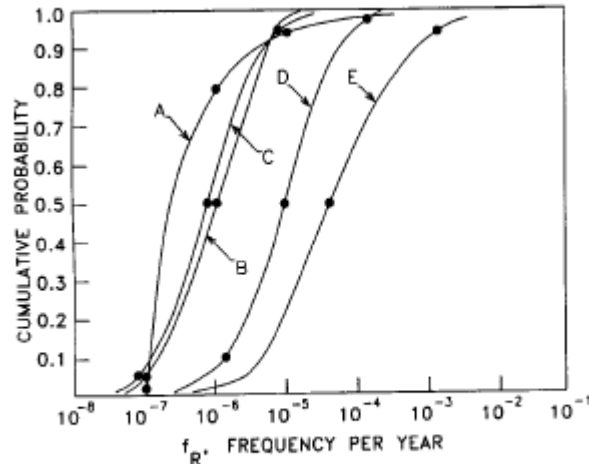
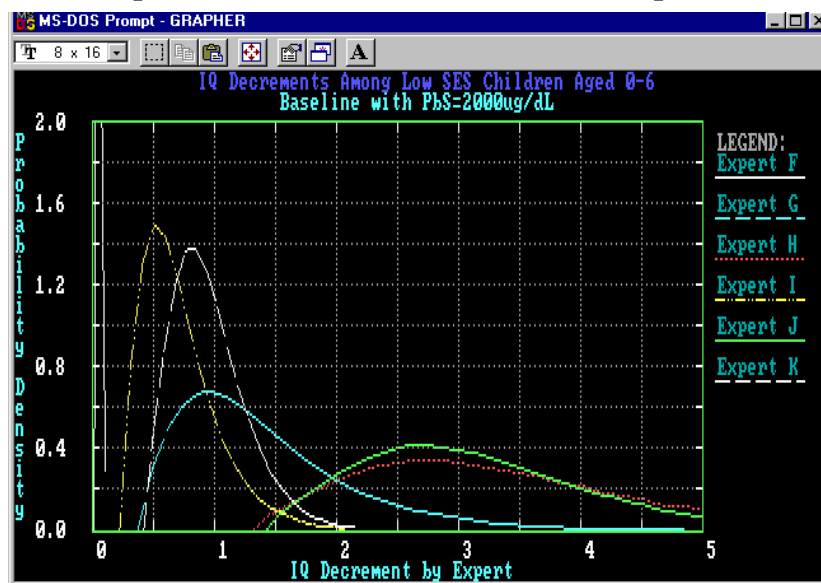
Figure 6.7. Example Presentation of CDFs of Multiple Experts in an Elicitation

Fig. 2. Cumulative probability distributions for five experts over annual RCS rupture frequency.

Source: Keeney and von Winterfelt (1991)

For the same study as in the figure just above, the Argonne National Laboratory ran the elicited probability distributions through their model, and presented the alternative results (in this case, estimated IQ decrement) that would result from application of each expert's estimates (Figure 6.8).

Figure 6.8. Sample Presentation of Results of Each Expert from an Elicitation

Source: Whitfield and Wallsten (1989)

6.5.4 How can Tables be Used?

Some elicitations, or parts of an elicitation, may not seek quantitative data, and even when quantitative results are elicited, a table may be used to summarize the differences in paradigms among experts, helping to explain differences among elicited values. Table 6.5 provides a concise but excellent non-quantitative summary of expert responses.

Table 6.6. Sample Compilation of Qualitative Expert Responses

A	Theoretical Construct for Long-and Short-term Exposure Effects on Premature Mortality
	See discussion of F1 and F2 in text
B1.	What do you believe to be the major causes of death associated with long-term exposure to PM_{2.5}? (In order of importance)
A	<ul style="list-style-type: none"> • Cardiovascular disease • lung cancer • Respiratory disease
B	<ul style="list-style-type: none"> • Cardiovascular disease • Respiratory disease • Not cancer – does not believe PM is likely to be a significant contributor to cancer risk
C	<ul style="list-style-type: none"> • Cardio-respiratory diseases probably constitute the bulk of the effects of PM but because cardiac deaths represent a very substantial portion of all deaths in the U.S. “But then our air pollution related effects are a very small part of that total.” • “I think our data is highly uncertain with regard to the issue of lung cancer associated with contemporary levels of airborne particulate material and I think even more uncertain with regard to other cancers.” • He thinks that PM exposure does not create a unique disease related to PM exposure. Instead, it “adds to the wear and tear of life.”
D	<ul style="list-style-type: none"> • Broad Category of Effects: Cardio-respiratory deaths <ul style="list-style-type: none"> • Heart Disease (CHD) • COPD (particles likely contribute, but are not a major contributor) • stroke, possibly • Cancer
E	<ul style="list-style-type: none"> • cardiovascular deaths • respiratory deaths (COPD, pneumonia, flu, infectious disease) • lung cancer

Source: USEPA, 2004, technical appendix

Alternatively, a complex but clear presentation of quantitative expert elicitation results in tabular format is provided in Table 6.6.

Table 6.6. Example Table for Responses of Multiple Experts

Summary of expected extinctions relative to the number of extinctions experienced during historical times, estimated by the 11 forest ecosystem experts. The notations 'w/o CC' and 'w/CC' refer to extinctions without and with a $2 \times [\text{CO}_2]$ climate change

Expert	Condition	North America above 45° N			Tropical forests between 20° N and 20° S		
		Mammals	Birds	Vascular plants	Mammals	Birds	Vascular plants
1	w/o CC	2	1.3	4	3.3	8	2
	w/ CC	5	1.4	8	6	10	4
2	w/o CC	0.5	3.3	1.1	2	10	10
	W/ CC	$0.5 + e$	$3.3 + e$	3.9	$2 + e$	12	12
3	w/o CC	1	1	0.2	10	2	3
	w/ CC	2	1–2.5	0.2	10	2–30	3–20
4	w/o CC	1	1.1	1.1	2	1.2	1.2
	w/ CC	1.3	1.2	1.2	2	1.2	1.3
5	w/o CC	10	3.3	1.05	250	10	2
	w/ CC	10	5	1.05	250	12	2
6	w/o CC	dk	dk	1	dk	dk	2.5
	w/ CC	dk	dk	1	dk	dk	6
7	w/o CC	2	6.7	3	10	>10	dk
	w/ CC	$2 + e$	10	4.5	$10 + e$	$>10 + e$	$1.5-2(\text{ma})$
8	w/o CC	1–1.5	2.5	2.5–3	4	4	6
	w/ CC	1–1.5	4.5	2.5–3	$4 + e$	$4 + e$	$6 + e$
9	w/o CC	1	1	'	10	10	10
	w/ CC	1	1.2	1.3	11	$10 + e$	11
10	w/o CC	~0	dk	0.5	1	dk	1
	w/ CC	e to 10	dk	2	$1 + e$	dk	2
11	w/o CC	~0	0.33	~0	5	2	~2
	w/ CC	0.5	$0.33 + e$	2	$5 + e$	$2 + e$	$\sim 2 + e$

dk = don't know.
 Note 1: Some North American vascular plants may go extinct in the wild, but people will continue to propagate them.
 Note 2: If there is a climate state change that results in the loss of the tropical forests tropical species loss could be much greater, but the climate scenario considered is too modest to produce this change.
 Note 3: Expert considered all of North America including mountains of the SW. Without human intervention mammalian extinctions could be 10 times historical level. Human interventions could largely eliminate this.
 Note 4: There would be seed bank efforts.

Source: Morgan et al. (2001)

As will be discussed in Section 6.4, the technical documentation and other presentations should also present important sensitivity analyses. One presentation in particular that would be useful is analysis of the importance overall of the uncertainty in the elicited results to the research question or decision at hand. Given the controversy potentially surrounding aggregation of experts' beliefs, a very useful presentation is the sensitivity analysis to different methods of combining expert judgments, or of using them individually. Table 6.7 illustrates such a sensitivity analysis.

Table 6.7. Example Table Presenting a Sensitivity Analysis For Combining Experts

Table D-5

Sensitivity Analysis of Combined Results for Effects of Long-term PM_{2.5} Exposure to Individual Expert Results

Percentiles	Combined Results All Experts ^a	Percent Change in Combined Results				
		Minus A	Minus B	Minus C	Minus D	Minus E
95 th	0.94	-2%	-0.4%	15%	-4%	-20%
75 th	0.59	-5%	13%	18%	-0.5%	-26%
50 th	0.30	-19%	21%	21%	-3%	-35%
25 th	0.15	-22%	21%	21%	5%	-41%
5 th	0.00	0%	0%	0%	0%	0%
Mean (Estimated) ^b	0.40	-8%	10%	20%	-3%	-28%

- a. The combined values are averages across experts at each percentile. The method gives equal weight to each expert's distribution. Combination method uses population-weighted distribution of annual mean PM_{2.5} concentrations in U.S. (from BENMAP model).
- b. The mean is estimated by combining the distributions using Monte Carlo simulation (see text for full discussion of methodology).

Source: USEPA (2004)

The results of sensitivity analyses can be summarized in tables and in graphics, and are often a critical component of high-level and public presentations, not only technical documentation. The decision-maker will want to understand the implications – or the lack of influence – of such methodological choices on the support for decision choices.

6.6 HOW CAN EES BE TRANSPARENT, DEFENSIBLE, AND REPRODUCIBLE?

EPA's Information Quality Guidelines (USEPA, 2002) requires that all information disseminated by the Agency meet a high standard of quality. This rigorous attention to quality is particularly relevant for EEs that were conducted as part of a regulatory process. When an EE is a component in a regulatory analysis, it may receive significant public attention.

EPA places great value in a regulatory process that is transparent, deliberate, and reviewable. In accordance with EPA's Peer Review Handbook (2006) influential scientific and technical work products used in decision making will be peer reviewed. The mechanism of review for a work product depends on its significance, the decision-making timeframe, level of public interest, and other factors. Regardless of the peer review mechanism selected, it is important that the reviewers (whether they are other EPA employees, independent external experts, or members of the public) are able to follow and understand the process of an EE.

The methods selected for analyzing EE data are of interest to peer reviewers. Given that many methods are available and the choice of a particular method could influence the outcome, peer reviewers are certain to examine this process. If a method is selected for arbitrary, subjective reasons, this is sure to attract criticism.

6.7 SHOULD EXPERT JUDGMENTS BE AGGREGATED FOR POLICY DECISIONS?

In section 5.4, the appropriateness and conduct of multi-expert aggregation was discussed within the context of the EE process. As was noted, not all EEs are amenable to meaningful aggregation (Morgan and Henrion, 1990). Even when a particular project and its results are well suited for aggregation, the analysts should show the aggregated results while preserving and presenting the richness of each individual expert's beliefs. This section provides additional thoughts on issues concerning the aggregation of multiple experts.

Given potential impact that aggregation can have on the interpretation of expert results, the use of aggregation is likely to be decided on a case-by-case basis. Such decisions may be addressed by the peer review. To support this peer review, any decision to combine experts should be well documented to illustrate why and how the experts were aggregated. This documentation should include a well-developed explanation of the rationale for the decision, the method selected, the influence of aggregation on findings, and sensitivity of the results to aggregation by different methods (or not aggregating). Meta-analytic techniques can be used to estimate relative importance of differences among expert views or of combining elicited judgments.

Selecting an approach for aggregation and making the decisions to implement that method may reflect the preferences of an analyst. In a public policy process that needs to be transparent and defensible. Simple averaging may be the method of aggregation most likely to receive broad acceptance and survive the scrutiny of peer review. Aggregation by other methods may require the analysts to make choices that could be perceived as arbitrary. Simple averaging is easy to put into practice and avoids the numerous choices that must be made for other

methods. In addition, simple averaging is equitable and has been shown to be robust and perform well in a variety of situations (Clemen and Winkler, 1999).

6.8 HOW CAN EXPERT ELICITATION RESULTS AND OTHER PROBABILITY DISTRIBUTIONS BE INTEGRATED?

EEs are often conducted to substitute for missing empirical data (see Chapter 3). After an EE has been conducted, the EE's results may be presented independently and alongside the results of other analyses. However, in general, it is useful to integrate the EE results with the other empirical data.

Risk assessments and policy analyses frequently require the combination of multiple types of disparate data. Because these data, including EE results, come from multiple disciplines and sources, it is important to integrate the data with caution. After integrating the elicited results with other parameters in a larger analysis, the analyst should critically evaluate the outcome to consider whether it is physically, biologically, and logically plausible.

Bayesian updating is a useful method for integrating the results of an EE with a prior distribution. In simple terms, Bayesian updating may be described as having one set of data or one distribution – or “prior,” then updating that prior as new data become available. The state of knowledge before the EE may be used as a prior. The findings of the EE can then provide an update on the prior. As better observations become available, they should be used.

6.9 HOW CAN AN EXPERT ELICITATION BE EVALUATED POST HOC?

The importance of presenting sensitivity and uncertainty analyses to provide insight into the strengths and weaknesses of the EE has already been discussed. In addition, there is additional evaluation of the EE that serves both future research design and the influence of the EE on the analytic support for various decision options.

A posteriori (post hoc) analyses should consider choice of model, distributions, bounding of the parameters, and method of combination (if any), as well as the parameters themselves. Whether a “back of the envelope” analysis or a more formal approach is used will depend on the importance of the findings. A post hoc analysis of the EE in the context of an integrated result is standard practice program evaluation and considering possible needs for future work. Do the results of the EE sufficiently answer the question that was at issue? Did the EE reduce or resolve controversy? Does the reduction of differences assist in determining what should (or should not) be done? If not, has the EE revealed how a different analysis might resolve the controversy, or has it exposed the need for different or new data or models? Was the original question well posed, or should the question be restated or better defined?

6.10 SUMMARY

The protocol is developed, the experts are selected, and the probabilities are elicited; but, the findings of an EE are only beneficial after they are presented. Because the presentation of results is what most readers and decision makers will see and read, it is used to judge the findings, form opinions, and make decisions. While a Technical Support Document contains all relevant details, pertinent findings must be abstracted for effective presentation to different users. The manner of presentation is critical because users have various backgrounds, preferences, and paradigms for using data. Hence, the presentation of results should ideally be part of a communication strategy that focuses on users and their needs. This chapter provided some examples for communicating EE results via probabilistic descriptions, text, figures, and tables. These examples were effective in their contexts and similar presentation can be considered by other practitioners. In addition, this chapter discussed issues concerning how to make defensible decisions about the aggregating expert judgments, combining EE results with other data, and providing peer review of findings.

7.0 FINDINGS AND RECOMMENDATIONS

The purpose of this Task Force was to initiate a dialogue within the Agency about the choice, conduct (including selection of experts), and use of EE and to facilitate future development and appropriate use of EE methods. The Task Force facilitated a series of discussions to explore the potential utility of using EE and to evaluate and address issues that may arise from using this approach. Based on those discussions, the Task Force has developed a set of findings and recommendations concerning 1) when it is appropriate to use EE (i.e., what constitutes “good practice” in deciding whether to conduct an EE), 2) how to plan and conduct such assessments, and 3) how to present and use the results of such assessments. The Task Force has also identified various recommended steps to facilitate future development and application of these methods within EPA. Section 7.1 and Section 7.2 summarizes the Task Force’s summarizes findings and recommendations, respectively.

7.1 FINDINGS

The findings of the Task Force are as follows:

7.1.1 What is Expert Elicitation?

- For the purposes of this White Paper, the Task Force has developed an operational definition of EE as the formal, systematic process of obtaining and quantifying expert judgment on the probabilities of events, relationships, or parameters. This definition also applies to expert judgment as identified in existing or proposed guidelines, including OMB’s *Circular A-4* and EPA’s revised *Cancer Risk Assessment Guidelines*.
 - EE is recognized as a powerful and legitimate tool. It can enable quantitative estimation of uncertain values and can provide uncertainty distributions where data are unavailable or inadequate. In addition, EE may be valuable for questions that are not necessarily quantitative such as model conceptualization or design of observational systems.
 - EE is one type of expert judgment activity. In general, expert judgment is an inherent and unavoidable part of many EPA assessments and decisions. Expert judgment is required in many stages of EPA analyses (e.g., problem formulation, model selection, study selection, estimation of input values, etc.) as well for the interpretation and communication of results. In addition, expert judgment is a component in the external peer review of EPA assessments.
-

- EE concerns questions of scientific information rather than societal values or preferences. In the broader set of tools for expert judgment, there are methods for capturing and incorporating values and preferences. In addition, the elicitation of preferences or economic valuation (e.g., willingness to pay for avoided risks) are related topics; but, were not the focus of the Task Force are not included in this White Paper's definition of EE.
- The results of an EE provide a characterization of the current state of knowledge for some question of interest. This is useful when traditional data are unavailable or inadequate. However, because an EE does not include measurements, observations, or experiments of the physical environment, it does not create new empirical data. Rather, it provides subjective estimates from experts that characterize the state of knowledge about some uncertain quantity, event, or relationship.

7.1.2 What is the Role of Expert Elicitation at EPA?

- Past experience with EE at EPA (e.g., in OAQPS since the late 1970s) indicates that it can provide useful, credible results. NAS has highlighted these past efforts as exemplary and recommended that EPA continue in the direction established by these precedents.
 - The use of EE is appropriate for some situations; but, not for others. Factors favoring the use of EE include: inadequate information to inform a decision, lack of scientific consensus, and the need to characterize uncertainty. Factors favoring alternatives to EE include theoretical and practical limitations. See section 4.6 for a summary of these factors. Typically, an EE requires a significant investment of resources and time to provide credible results.
 - EE can work well when a scientific problem has a body of knowledge; but, lacks a consensus interpretation. For this case, expert beliefs about the value and meaning of data can provide valuable assessments and insights. This may be the case for an emerging scientific challenge or one that depends on uncertain future events. However, when a problem has abundant relevant empirical data and relative consensus exists in the scientific community, there is probably little need to conduct an EE. At the other end of the spectrum, if data are inadequate for the experts to develop judgments, an EE may not be worthwhile.
 - Given that EPA uses other more familiar approaches to characterize uncertainty, the application and acceptance of EE at EPA will likely grow with experience. If early EE
-

efforts are well-designed and implemented, this will promote the credibility and endorsement of EE within the Agency and by external stakeholders.

- The nature of the regulatory process (i.e., legal, political, financial, technical, and procedural considerations) will influence whether and how to conduct an EE and how to communicate and use results. Within the regulatory process, EPA can use EE to encourage transparency, credibility, objectivity (unbiased and balanced), rigor (control of heuristics and biases), and relevance to the problem of concern.

7.1.3 What Factors are Considered in the Design and Conduct of an Expert Elicitation?

- Designing an EE and interpreting results are generally case-specific and context-specific. Although the conduct of an EE does not lend itself to a rigid cookbook approach, there are a number of steps to follow to promote a credible and defensible EE (Chapter 5).
- An EE includes distinct roles for the members of the project team (generalist, analyst, and subject matter expert) and the experts whose judgments are the subject of the EE.

7.2 RECOMMENDATIONS

The recommendations of the Task Force are as follows:

7.2.1 What Challenges are Well-Suited for an Expert Elicitation?

- EE is well-suited for challenges with complex technical problems, unobtainable data, conflicting conceptual models, available experts, and sufficient financial resources.
 - EE should be considered to characterize uncertainty, where it can not be addressed adequately by existing data or additional studies within the necessary timeframe. EE should also be considered to fill data gaps where traditional data are unobtainable given the importance and relevance of the data to the decision and the decision schedule.
 - EE results can provide a proxy for traditional data, but, it is not equivalent to valid empirical data. When appropriate empirical data can be obtained given the available time and resources, EE should not be used as a substitute for conventional research.
 - Before deciding to conduct an EE, managers and staff should engage in discussions about:
-

- The goals of the EE and the basis for the selection of this approach;
 - The anticipated output from the EE and how it may be used in the overall decision;
 - EEs that have been used for similar types of decisions;
 - The outcome of the EE upon completion; and
 - The time and cost of conducting a defensible EE.
- Investigators considering the use of EE may want to consider alternative methodologies to characterize uncertainty or fill data gaps (Chapter 4).

7.2.2 How Should an Expert Elicitation be Designed and Conducted?

- EPA should develop policy, guidance, training, and/or tools (drawing on this White Paper and the literature cited in herein) to support the conduct and use of EE. These resources should address the following issues:
 - Standards of quality for EEs that are a function of their intended use (e.g., to inform research needs, to inform regulatory decisions, etc.) and a minimum set of best practices.
 - How to interpret the quality of the results and the EE process.
 - How to review or interpret efforts in the context of their use (i.e., how does acceptability depend upon context).
 - The role of stakeholders early in the EE planning process to provide input on relevant questions or issues.
 - Appropriateness of secondary application of EE results (i.e., the use of results beyond the purpose intended when the study was designed).
 - Under what circumstances and how should experts' judgments be combined.
 - Comparison of quality/usefulness of various types of research findings: empirical data, external expert recommendations, and EE result.
 - Whether the judgments of individual experts should be weighted differentially to produce an aggregate judgment. If so, what criteria measures are most equitable?
 - How results should be used and communicated to decision-makers.
-

- Until this EE resource is developed, those considering and/or conducting an EE within EPA should be encouraged to carefully consider the issues, examples, concerns, and references presented in this White Paper.
 - EEs should focus on those aspects of uncertainty that cannot be adequately described by empirical data. The EE protocol should avoid overlap between what the experts are asked and what the data adequately describe.
 - For some questions that require characterization of a quantity that encompasses several aspects of uncertainty, it may be appropriate to disaggregate the problem and ask the experts to assess each aspect separately.
 - The long-term success of EE may depend heavily on whether its early applications are considered credible and helpful by decision makers and stakeholders. Hence, the utility and acceptability of EEs at EPA can be facilitated by well-designed and implemented studies. Therefore, the Task Force recommends that:
 - Early efforts by EPA program or regional offices with little EE experience should include collaboration with knowledgeable staff within the Agency (e.g., members of this Task Force) and/or external EE specialists. These efforts should also bear in mind the approaches and considerations outlined in Chapter 5 for design and conduct of an EE.
 - Given that the success of the EPA/OAQPS 1980s efforts (cited by the 2003 NAS panel as exemplary efforts) benefited significantly from early collaborations with SAB and external EE specialists, similar collaborations may be highly desirable in early efforts by other offices. This collaborative approach can help to ensure the quality, credibility, and relevance of these efforts.
 - Training materials should be developed to teach EE basics. Those involved in the design, conduct, or use of EEs should draw on these materials to promote familiarity with EEs and to obtain the advice of those with greater experience.
 - To facilitate learning among EPA staff and offices about EEs, EPA should make EE resources available, including:
 - Examples of well-conducted EEs, including protocols, criteria used for selecting experts, peer reviews, etc.
-

- Documentation for these EEs should include discussion of advantages and limitations and lessons learned.
- Internal tracking of ongoing and planned EE efforts.

This could be provided as part of the probabilistic analysis web site that is under development by the RAF's Probabilistic Work Group.

- Additional research and discussion is recommended to better determine the appropriate level of disaggregation for EE questions.

7.2.3 How Should Experts be Selected?

- For highly influential and potentially controversial EEs, additional steps in the expert selection process are recommended. These steps can help to establish that experts were selected without bias and to include the range of scientific perspectives. This is of special important if an EE project may aggregate expert judgments.
 - Any potential conflicts of interest should be disclosure fully. Any potential conflicts should be considered on a case-by-case basis to determine if the nature of the conflict precludes participation in the EE.
 - The involvement of the EE's sponsor (e.g., EPA) in the process of nominating and selecting experts should be decided on a case-by-case basis. If the sponsor does not participate, this could improve the perception of objectivity. On the other hand, EPA may want to have more active control on the selection process because it is ultimately responsible for assuring the quality and credibility of its EEs. As a default practice, the EE project team encouraged to give this issue careful consideration.
 - To comply with the Paperwork Reduction Act, if more than nine experts will participate in an EE, EPA must submit an information collection request to OMB. To avoid this time-consuming request, using a maximum of nine experts may be expedient; however, this must be balanced by the importance of the EE, the range of different believes, and the availability of experts.
-

7.2.4 How Should Expert Elicitation Results be Presented and Used?

- Experts who participate in an EE should be identified by name and institutional affiliation; but, their actual judgments may be anonymous. However, a record of all judgments should be maintained and provided for any required auditing or if needed for peer review.
- EPA should link and/or integrate its EE efforts with its ongoing efforts to promote the use of PRA. Lessons about communicating to decision makers and other stakeholders can be derived from common efforts that are related to probabilistic analysis.

7.2.5 What is the Role of Peer Review and Peer Input in Expert Elicitation Projects?

- Peer review of any EE exercise should focus on the EE process, including how the experts were selected, what they were provided, how the EE was conducted (including controlling for heuristics and biases), and how the results were analyzed. The peer review should include subject matter experts and EE specialists. The purpose of peer reviewing an EE is to review the process, not be to second-guess the expert judgments.
- Depending on the purpose of the EE, a peer review of the expert selection process, the EE methods, and the results of any pilots may be appropriate prior to conducting the actual EE. Peer input about the EE protocol (i.e., prior to conducting the elicitations) may be very useful. Receiving this ex ante consultation can improve the quality of the EE and maximize resource efficiency.

7.2.6 What Outreach, Research, and Future Steps are Recommended?

- EPA should work to develop guidance and policy on the conduct and use of EE. Internal discussions should use this White Paper as a guide to discuss how the use of this method might improve Agency decision making, including discussions of lessons learned and potential guidance development.
 - EPA should work collaboratively with other Federal agencies such as FDA and USDA (using this White Paper as a guide) to discuss how the use of EE might improve government decision making, including discussions of lessons learned and potential guidance development.
-

- EPA should identify cross-cutting scientific issues where an improved characterization of uncertainty could impact multiple Agency assessments (e.g., a parameter or relationship that affects numerous exposure analyses) and which are good candidates for EE.
 - EPA should support research on EE methods development and evaluation that are related its environmental and regulatory mission. This research should include the investigation of probabilistic and non-probabilistic EE methodologies and seek to determine their appropriate use and limitations. EPA should consider how these efforts can be implemented through building intramural expertise (the Economics and Decision Sciences multiyear plan) and extramural research (the STAR grants program).
 - Management should continue to support EPA efforts to co-sponsor and participate in workshops, colloquia, and professional society meetings (e.g., SRA, SOT, etc.). These engagements promote dialogue, encourage innovation, and the facilitate experience sharing that can ultimately improve the quality and use of EE assessments.
-

REFERENCES

- Anderson, J. L., 1998. Embracing Uncertainty: The Interface of Bayesian Statistics and Cognitive Psychology. *Conservation Ecology* online: 1(2). (Available at www.ecologyandsociety.org/vol2/iss1/art2/#AdaptingBayesianAnalysisToTheHumanMind:GuidelinesFromCognitiveScience).
- Ariely, D., W.T. Au, R.H. Bender, D.V. Budescu, C.B. Dietz, H. Gu, T.S. Wallsten, and G. Zauberman, 2000. The Effects of Averaging Subjective Probability Estimates Between and Within Judges. *Journal of Experimental Psychology: Applied*, 6(2):130-147.
- Arnell, N. W., E. L. Tompkins, and W. N. Adger, 2005. Eliciting Information from Experts on the Likelihood of Rapid Climate Change. *Risk Analysis* 25:6:1419-1431.
- Ayyub, B.M., 2000. Methods for Expert-Opinion Elicitation of Probabilities and Consequences for Corps Facilities. Prepared for U.S. Army Corps of Engineers Institute for Water Resources, IWR Report 00-R-10. (Available at: www.iwr.usace.army.mil/iwr/pdf/MethodsforEEfinal1.PDF).
- Ayyub, B.M., 2001. A Practical Guide on Conducting Expert-Opinion Elicitation of Probabilities and Consequences for Corps Facilities. Prepared for U.S. Army Corps of Engineers Institute for Water Resources, IWR Report 01-R-01, Alexandria, VA. (Available at: www.iwr.usace.army.mil/iwr/pdf/PEEfinal.PDF).
- Batz, M.B., M.P. Doyle, J.G. Morris, J. Painter, R. Singh, R.V. Tauxe, M.R. Taylor, M.A. Danilo, and L.F. Wong, 2005. Attributing Illness to Food. *Emerging Infectious Diseases*, 11(7), July 2005. (Available at: www.medscape.com/viewarticle/507913_11)
- Bedford and R.M. Cooke, 2001. Probabilistic Risk Analysis: Foundations and Methods.
- Berger, J. O. and D. A. Berry, 1988. Statistical Analysis and the Illusion of Objectivity. *American Scientist* 76: 159-165.
- Bloom, D.L. et. al., 1993. Communicating Risk to Senior EPA Policy Makers: A Focus Group Study. U.S. EPA Office of Air Quality Planning and Standards, Research Triangle Park, NC.
- Brown, S. R., 1980. Political Subjectivity: Application of Q Methodology in Political Science. Yale University Press, New Haven.
- Brown, S. R., 1996. Q methodology and qualitative research. *Qualitative Health Research* 6(4): 561-567.
-

Bruine de Bruin, W., P. S. Fischbeck, N. A. Stiber and B. Fischhoff, 2002. What Number is "Fifty-Fifty"?: Redistributing Excessive 50% Responses in Elicited Probabilities. *Risk Analysis* 22: 713-723.

Brunner, N. and M. Starkl, 2004. Decision Aid Systems For Evaluating Sustainability: A Critical Survey. *Environmental Impact Assessment Review* 24: 441-469.

Bunn, Derek W., 1984. Applied Decision Analysis. McGraw-Hill Book Company, New York.

Chechile, R. A., 1991. Probability, Utility, and Decision Trees in Environmental Analysis. In: Environmental Decision Making: A Multidisciplinary Perspective, R. A. Chechile and S. Carlisle, eds., Van Nostrand Reinhold, New York, p. 64-91.

Clemen, R., 1996. *Making Hard Decisions*, Duxbury Press, Belmont, CA

Clemen, R. and R. L. Winkler, 1999. Combining Probability Distributions from Experts in Risk Analysis. *Risk Analysis* 19 (2): 187-203.

Clemen, R.T and R.L. Winkler, 1985. Limits for the Precision and Value of Information from Dependent Sources. *Operations Research* 33(2):427-442

Clemen, R.T., 1989. Combining Forecasts: A Review and Annotated Bibliography. *International Journal of Forecasting* 5:559-583.

CNS (Center for Nonproliferation Studies), 2003. Bioterrorism Threat Assessment and Risk Management Workshop, Report for US Department of Energy, Final Report and Commentary, Center for Nonproliferation Studies, Monterey Institute of International Studies.

Cocks, D. and J. Ive, 1996. Mediation Support for Forest Land Allocation: The SIRO-MED System. *Environmental Management* 20: 41-52.

Cooke, R.M., 1991. Experts in Uncertainty: Opinion and Subjective Probability in Science. Oxford University Press, New York.

Cooke, R. M. and L. H. J. Goossens, 2000. Procedures Guide for Structured Expert Judgment in Accident Consequence Modelling. *Radiation Protection Dosimetry* 90(3): 303-309.

Cosmides, L., and J. Tooby, 1996. Are Humans Good Intuitive Statisticians After All? Rethinking Some Conclusions From The Literature On Judgment Under Uncertainty. *Cognition* 58:1-73.

_____, 1994. Why The Distinction Between Single-Event Probabilities And Frequencies Is Relevant For Psychology (And Vice Versa). In G. Wright and P. Ayton, editors. Subjective probability. Wiley, New York, pp.129-162.

Crawford-Brown, D., 2005. The Concept of Sound Science in Risk Management Decisions. *Risk Management: An International Journal* 7: 7-20.

Crawford-Brown, D., 2001. Scientific Models of Human Health Risk Analysis in Legal and Policy Decisions. *Law and Contemporary Problems* 64(4): 63-81.

Cullen, A.C. and H.C. Frey, 1999. *Probabilistic Techniques in Exposure Assessment: A handbook for dealing with variability and uncertainty in models and inputs*, Plenum, New York

Dalkey, N.C., 1969. The Delphi Method: An Experimental Study of Group Opinion. Rand Corp, Santa Monica, CA.

Dehaene, S., 1997. The Number Sense: How The Mind Creates Mathematics. Oxford University Press, Oxford, UK.

Delbecq, A.L., A.H. Van de Ven and D.H. Gustafson, 1975. Group Techniques for Program Planning: A Guide to Nominal Group and Delphi Processes. Scott, Foresman & Company, Glenview, IL.

Dyer, J. S., P. C. Fishburn, R. E. Steuer, J.Wallenius, and S. Zionts. 1992. Multiple Criteria Decision Making, Multiattribute Utility Theory: The Next Ten Years. *Management Science* 38 (5):645-654.

Ehrlinger, J., T. Gilovich, and L. Ross, 2005. Peering Into the Bias Blind Spot: People's Assessments of Bias in Themselves and Others. *Personality and Social Psychology Bulletin*, Vol. 31, No. 5, May 2005, p. 1-13.

Ehrmann, J.R. and B. L. Stinson, 1999. Joint Fact-Finding and the Use of Technical Experts. In: *The Consensus Building Handbook: A Comprehensive Guide to Reaching Agreement*, L. E. Susskind, S.McKernan and J.Thomas-Larmer, eds. Sage Publications, Thousand Oaks, CA, pp. 375-399.

Einhorn, H.J., and R M. Hogarth, 1978. Confidence in Judgment: Persistence of the Illusion of Validity, *Psychological Review* 85(3):395-416.

European Commission, 2000. Procedures Guide for Structured Expert Judgment, Nuclear Science and Technology, Directorate-General for Research, EUR 18820EN.

Evans J.S., G.M. Gray, R.L. Sielken Jr, A.E. Smith, C. Valdez-Flores, and J.D. Graham, 1994. Use Of Probabilistic Expert Judgment In Uncertainty Analysis Of Carcinogenic Potency. *Regulatory Toxicology and Pharmacology* 2:15-36.

Feagans, T. and Biller, 1981. Risk Assessment: Describing the Protection Provided by Ambient Air Quality Standards. *The Environmental Professional* 3(3/4):235-247.

Finkel, A., 1990. *Confronting Uncertainty in Risk Management: A Guide for Decision Makers*, Center for Risk Management, Resources for the Future, Washington, DC.

Fischhoff, B., 2003. Judgment And Decision Making. In: *The Psychology of Human Thought*, R. J. Sternberg and E.E.Smith, eds. Cambridge University Press, Cambridge, UK. p.153-187.

Fischhoff, B. and J. S. Downs, 1998. Communicating Foodborne Disease Risk, *Emerging Infectious Diseases* 3(4): 489-495.

Florig H.K., M.G. Morgan, K.M. Morgan, K.E Jenni, B. Fischhoff, P.S. Fischbeck, and M.L DeKay, 2001 A Deliberative Method For Ranking Risks (I): Overview And Test Bed Development. *Risk Analysis* 21:913-21.

Fos P.J., and C.L. McLin, 1990. The Risk Of Falling In The Elderly: A Subjective Approach. *Medical Decision Making* 10:195-200.

Funtowicz, S. O. and J. R. Ravetz, 1990. *Uncertainty and Quality in Science for Policy*. Kluwer Academic Publishers, Dordrecht, The Netherlands.

Garthwaite, P. H., J. B. Kadane and A. O'Hagan, 2005. Statistical Methods for Eliciting Probability Distributions. *Journal of the American Statistical Association* 100: 680-701.

Genest, C. and J.V. Zidek, 1986. Combining Probability Distributions: A Critique and an Annotated Bibliography. *Statistical Science* 36:114-148.

Gigerenzer, G., 1994. Why the distinction between single-event probabilities and frequencies is relevant for psychology (and vice versa). In: *Subjective probability*, G. Wright and P. Ayton, eds. Wiley, New York. pp. 129-162.

Gigerenzer, G., and U. Hoffrage, 1995. How to improve Bayesian reasoning without instruction: frequency formats. *Psychological Review* 102(4):684-704.

Gokhale, A. A., 2001. Environmental Initiative Prioritization with a Delphi Approach: A Case Study. *Environmental Management* 28: 187-193.

Hamm, R.M., 1991. Modeling Expert Forecasting Knowledge for Incorporation into Expert Systems, Institute for Cognitive Studies, University of Colorado, Institute of Cognitive Science (ICS) Tech Report 91-12, Boulder, CO.

Hammond, K. R., R. M. Hamm, J. Grassia and T. Pearson, 1987. Direct Comparison of the Efficacy of Intuitive and Analytical Cognition in Expert Judgment. *IEEE Transactions on Systems, Man and Cybernetics* SMC-17: 753-770.

Harremoes, P., D. Gee, M. MacGarvin, A. Stirling, J. Keys, B. Wynne and S. G. Vaz, 2001. Twelve late lessons. In: *Late lessons from early warning: the precautionary principle 1896-2000*,

P. Harremoes, D.Gee, M.MacGarvin, A.Stirling, B.Wynne and S.G.Vaz, eds. European Environment Agency, Copenhagen, Denmark. pp. 168-194.

Hattis, D., and D. Burmaster, 1994. Assessment of Variability and Uncertainty Distributions for Practical Risk Analyses, *Risk Analysis*, 14(5): 713-730.

Hawkins N.C. and J.D. Graham, 1988. Expert Scientific Judgment And Cancer Risk Assessment: A Pilot Study Of Pharmacokinetic Data. *Risk Analysis* 8:615-25.

Hawkins, N.C. and J.S. Evans, 1989. Subjective Estimation of Toluene Exposures: A Calibration Study of Industrial Hygienists. *Applied Industrial Hygiene* 4:61-68.

Hogarth, R.M., 1978. A Note on Aggregating Opinions. *Organizational Behavior and Human Performance* 21:40-46.

Hokkanen, J., R. Lahdelma and P. Salminen, 2000. Multicriteria decision support in a technology competition for cleaning polluted soil in Helsinki. *Journal of Environmental Management* 60:339-348.

Hora, S. and M. Jensen, 2002 Expert Judgement Elicitation. Department of Waste Management and Environmental Protection SSI report ISSN 0282-4434. Available at: http://www.ssi.se/ssi_rapporter/pdf/ssi_rapp_2002_19.pdf.

Hora, S. C., 2004. Probability Judgments for Continuous Quantities: Linear Combinations and Calibration. *Management Science* 50:597-604.

Hora, S.C., 1992. Acquisition of Expert Judgment: Examples form Risk Assessment. *Journal of Energy Engineering* 118(2):136-148.

Hsu, M., M. Bhatt, R. Adolphs, D. Tranel, and C. Camerer, 2005. Neural Systems Responding to Degrees of Uncertainty in Human Decision Making. *Science* 310:1680-1683.

Humphreys, P., O. Svenson, A. Vari, T. Englander, J.Vecsenyi, W. Wagenaar and D. Von Winterfeldt, 1983. Analysing and Aiding Decision Processes. North-Holland Publishing Company, Amsterdam.

Ibrekk, H. and M.G. Morgan, 1987. Graphical Communication of Uncertain Quantities to Non-Technical People. *Risk Analysis* 7:519-529.

IEC, 2004. An Expert Judgment Assessment of the Concentration Response Relationship Between PM_{2.5} Exposure and Mortality, prepared for U.S. EPA Office of Air Quality Planning and Standards. (Available at: <http://www.epa.gov/ttn/ecas/regdata/Benefits/pmexpert.pdf>).

IPCC, 2005. Guidance Notes for Lead Authors of the IPCC Fourth Assessment Report on Addressing Uncertainties, Intergovernmental Panel on Climate Change, July 2005. (Available at <http://www.ipcc.ch/activity/uncertaintyguidancenote.pdf>).

IPCC, 2001. Quantifying Uncertainties in Practice (Chapter 6) *In: IPCC Good Practice Guidance and Uncertainty Management in National Greenhouse Gas Inventories*, Intergovernmental Panel on Climate Change, GPGAUM-Corr.2001.01. (Available at http://www.ipcc-nggip.iges.or.jp/public/gp/english/6_Uncertainty.pdf).

Jamieson, D., 1996. Scientific Uncertainty and the Political Process. *The Annals of the American Academy of Political and Social Science* 545:35-43.

Jelovsek, F. R., D. R. Mattison and J. F. Young, 1990. Eliciting Principles of Hazard Identification from Experts. *Teratology* 42:521-533.

Johnson B.B., 2005. Testing And Explaining A Model Of Cognitive Processing Of Risk Information. *Risk Analysis* 25:631.

Johnson, B.B. and P. Slovic, 1995. Presenting Uncertainty in Health Risk Assessment: Initial Studies of its Effects on Risk Perception and Trust. *Risk Analysis* 15:485-494.

Jouini, M.N. and R.T. Clemen, 1996. Copula Models for Aggregating Expert Opinion. *Operations Research* 44(3):444-457.

Kaplan, S., 1992. Expert Information” versus “Expert Opinions.” Another Approach to the Problem of Eliciting /Combining /Using Expert Knowledge in Probabilistic Risk Analysis. *Reliability Engineering and System Safety* 35:61-72.

Karelitz, T. M., M. K. Dhimi, D. V. Budescu, and T. S. Wallsten, 2002. Toward a Universal Translator of Verbal Probabilities. In Proceedings of the 15th International Florida Artificial Intelligence Research Society (FLAIRS) Conference. AAAI Press, 298-502.

Keeney, R. L. and D. von Winterfeldt, 1991. Eliciting Probabilities from Experts in Complex Technical Problems. *IEEE Transactions on Engineering Management* 38(3):191-201.

Keith, D. W., 1996. When is it appropriate to combine expert judgments? *Climatic Change* 33:139-143.

Krayer von Krauss, M.P., E.A. Cashman, and M.J. Small, 2004. Elicitation of Expert Judgments of Uncertainty in the Risk Assessment of Herbicide-Tolerant Oilseed Crops, *Risk Analysis* 24(6):1515-1527.

Krupnick, A., R. Morgenstern, M. Batz, P. Nelson, D. Burtraw, J-S Shih, and M. McWilliams, 2006. Not a Sure Thing: Making Regulatory Choices under Uncertainty, Resources for the Future, Washington, D.C.

Kunzli N., S. Medina, R. Kaiser, P. Quenel, F. Horak, M. Studnicka, 2001. Assessment of Deaths Attributable to Air Pollution: Should We Use Risk Estimates based on Time Series or on Cohort Studies? *American Journal of Epidemiology* 153:1050-5.

- Levin, R., S. O. Hansson, and C. Rudén, 2004. Indicators of uncertainty in chemical risk assessments. *Regulatory Toxicology and Pharmacology* 39(1):33-43. (Available at: <http://www.sciencedirect.com/science/article/B6WPT-4B6682R-1/2/a7e2179d50dea279f99f696e4e4de4a7>)
- Libby, R. and R. K. Blashfield, 1978. Performance of a Composite as a Function of the Number of Judges. *Organizational Behavior and Human Performance*, 21:121-129.
- Linkov, I., A. Varghese, S. Jamil, T. P. Seager, G. Kiker and T. Bridges, 2004. Multi-criteria Decision Analysis: A Framework for Structuring Remedial Decisions at Contaminated Sites. In: *Comparative Risk Assessment and Environmental Decision Making*, I. Linkov and A. Ramadan, eds. Kluwer, New York., pp. 15-54.
- Linstone, H. A. and M. Turoff, 1975. *The Delphi Method: Techniques and Applications*. Addison-Wesley, Reading, MA. (Available at: <http://www.is.njit.edu/pubs/delphibook/>)
- Martin, S.A., T.S. Wallsten, and N.D. Beaulieu, 1995. Assessing the Risk of Microbial Pathogens: Application of a Judgment-Encoding Methodology. *Journal of Food Protection* 58(3):289-295.
- Mazur, A., 1973. Disputes Between Experts. *Minerva* 11: 243-262.
- McKeown, B. and D. Thomas, 1988. *Q Methodology*. Sage Publications, Newbury Park, CA.
- Merkhofer, M. W. and R. L. Keeney, 1987. A Multiattribute Utility Analysis of Alternative Sites for the Disposal of Nuclear Waste. *Risk Analysis* 7:173-194.
- Meyer, M. A. and J. M. Booker, eds. 2001. *Eliciting and Analyzing Expert Judgment: A Practical Guide*. Society for Industrial and Applied Mathematics; American Statistical Association, Philadelphia, PA.
- Morgan, K., 2005. Development of a Preliminary Framework for Informing the Risk Analysis and Risk Management of Nanoparticles. *Risk Analysis* 25(6):1621-1635.
- Morgan K.M., D.L. DeKay, P.S. Fischbeck, M.G. Morgan, B. Fischhof, and H.K. Florig, 2001. A deliberative method for ranking risks (II): Evaluation of validity and agreement among risk managers. *Risk Analysis* 21:923-37.
- Morgan, M.G., S.C. Morris, M. Henrion, D.A.L. Amaral and W.R. Rish, 1984. Technical Uncertainty in Quantitative Policy Analysis: A Sulfur Air Pollution Example. *Risk Analysis* 4:201-216.
- Morgan M.G. and D.W. Kieth, 1995. Subjective Judgments by Climate Experts. *Environmental Science and Technology* 29(10):468A-476A.
-

Morgan, M.G., 1998. Uncertainty in Risk Assessment. *Human and Ecological Risk Assessment* 4(1): 25-39.

Morgan, M.G., B. Fischhoff, A. Bostrom, and C.J. Atman, 2002. Risk Communication. A Mental Models Approach. Cambridge University Press, Cambridge, UK.

Morgan, M.G. and M. Henrion, 1990. Uncertainty: A Guide to Dealing with Uncertainty in Quantitative Risk and Policy Analysis. Cambridge University Press, Cambridge, MA.

Moss, R. and S.H. Schneider, 2000. Uncertainties – Guidance Papers on the Cross Cutting Issues of the Third Assessment Report of the IPCC. In: Title??, Pachauri, R., Taniguchi, R., and Tanaka, K. eds., World Meteorological Organisation, Geneva, Switzerland. p.??

North, W., B.R. Judd, and J.P. Pezier, 1974. “New Methodology for Assessing the Probability of Contaminating Mars,” *Life Sciences and Space Research* 13, pp. 103-109.

National Research Council (NRC), 1994. Science and Judgment in Risk Assessment. National Academy Press, Washington, DC.

Nauta, M., I. van der Fels-Klerx, and A. Havelaar, 2005. A Poultry-Processing Model for Quantitative Microbiological Risk Assessment. *Risk Analysis* 25(1):85-98.

NRC. 1999. Upgrading the Space Shuttle, National Academy Press, Washington, DC.

NRC. 2002. Estimating the Public Health Benefits of Proposed Air Pollution Regulations. National Academy Press, Washington, DC.

O'Hagan, A., 2005. Research in Elicitation. Research Report No. 557/05, Department of Probability and Statistics, University of Sheffield. Invited article for a volume entitled Bayesian Statistics and its Applications.

Otway, H. and D. von Winterfeldt, 1992. Expert Judgment in Risk Analysis and Management: Process, Context and Pitfalls. *Risk Analysis* 12(1):83-93.

Parenté, F. J. and J. K. Anderson-Parenté, 1987. Delphi Inquiry Systems. In: Judgmental Forecasting, G. Wright and P. Ayton (eds.), Wiley, Chichester, England. pp.129-156.

Pike, W.A., 2004. Modeling Drinking Water Quality Violations with Bayesian Networks, *Journal of the American Water Resources Association*, December::1563-1578.

Pope, C.A.,III, R.T. Burnett, M.J. Thun, E.E. Calle, D. Krewski, K. Ito, and G.D. Thurston, 2002. Lung Cancer Cardiopulmonary Mortality and Long-term Exposure to Fine Particulate Air Pollution. *Journal of the American Medical Association* 287(9): 1132-1141.

Pronin, E., T. Gilovich, and L. Ross, 2004. Objectivity in the Eye of the Beholder: Divergent Perceptions of Bias in Self Versus Others. *Psychological Review*, Vol. 111, No. 3, 781-799.

Renn, O., 1986. Decision Analytic Tools for Resolving Uncertainty in the Energy Debate. *Nuclear Engineering and Design* 93:167-179.

Renn, O., 1999. A Model for an Analytic-Deliberative Process in Risk Management. *Environmental Science and Technology* 33:3049-3055.

Renn, O., 2001. The Need For Integration: Risk Policies Require The Input From Experts, Stakeholders And The Public At Large. *Reliability Engineering and System Safety* 72:131-135.

Richmond, H.M., 1991. Overview of Decision Analytic Approach to Noncancer Health Risk Assessment. Paper No. 91-173.1 presented at Annual Meeting of the Air and Waste Management Association; Vancouver, BC.

RIVM, 2003. Uncertainty Analysis for NOx Emissions from Dutch passenger cars in 1998, RIVM Report 550002004/2003.

Rosenbaum, A.S., R.L. Winkler, T.S. Wallsten, R.G. Whitfield, and H.M. Richmond, 1995. An Assessment Of The Risk Of Chronic Lung Injury Attributable To Long-Term Ozone Exposure. *Operations Research* 43(1):19-28.

RTI, 2004. Peer Review of Expert Elicitation, memo to EPA Office of Air Quality Planning and Standards, Research Triangle Park, NC.

Saaty, T. L., 1990a. Multicriteria Decision Making: The Analytic Hierarchy Process. RWS Publications, Pittsburgh, PA.

Saaty, T. L., 1990b. The Analytic Hierarchy Process: Planning, Priority Setting, Resource Allocation. RWS Publications, Pittsburgh, PA.

Shackle, G. L. S., 1972a. Economic Theory and the Formal Imagination. In: *Epistemics and Economics: A critique of economic doctrines*. Cambridge University Press, Cambridge, UK. pp.3-24.

Shackle, G. L. S., 1972b. Languages for Expectation. In: *Epistemics and Economics: A critique of economic doctrines*. Cambridge University Press, Cambridge, UK. pp.364-408.

Shackle, G. L. S., 1972c. The Science of Imprecision. In: *Epistemics and Economics: A Critique Of Economic Doctrines*. Cambridge University Press, Cambridge, UK. pp.359-363.

Shrader-Frechette, K. S., 1991. Risk and Rationality: Philosophical Foundations for Populist Reforms. University of California Press, Berkeley, CA.

Slovic, P., 1986. Informing And Educating The Public About Risk. *Risk Analysis*. 6(4):403-415.

Slovic, P., B. Fischhoff, and S. Lichtenstein, 1988. Response Mode, Framing, and Information-Processing Effects in Risk Assessment. In: Decision Making: Descriptive, Normative, and Prescriptive Interactions, D. E. Bell, H. Raiffa and A. Tversky, eds. Cambridge University Press, Cambridge, UK. pp.152-166.

Slovic, P., B. Fischhoff, and S. Lichtenstein, 1979. Rating the Risks. *Environment* 21(4):14-20 and 36-39 (Reprinted in P. Slovic (ed.), *The Perception of Risk*, London, Earthscan, 2000).

Spetzler, C.S. and C.A.S. Staehl von Holstein, 1975. Probability Encoding in Decision Analysis. *Management Science* 22:3.

SRI, 1978. *The Use of Judgmental Probability in Decision Making*. SRI Project 6780. Prepared for U.S. EPA, OAQPS, RTP, NC.

Stahl, C. H. and A. J. Cimorelli, 2005. How Much Uncertainty is Too Much and How Do We Know? A Case Example of the Assessment of Ozone Monitoring Network Options. *Risk Analysis* 25:1109-1120.

Stahl, C. H., A. J. Cimorelli and A. H. Chow, 2002. A New Approach to Environmental Decision Analysis: Multi-criteria Integrated Resource Assessment (MIRA). *Bulletin of Science, Technology, and Society* 22:443-459.

Stephenson, W., 1935. Correlating Persons Instead Of Tests. *Character and Personality* 4:17-24.

Stephenson, W., 1935. Technique In Factor Analysis. *Nature* 136(3434): 297.

Stephenson, W., 1953. The Study of Behavior: Q-Technique and its Methodology. University of Chicago Press, Chicago.

Stiber, N.A.; M. Pantazidou, M.J. Small, 1999. Expert System Methodology For Evaluating Reductive Dechlorination At TCE Sites. *Environmental Science & Technology* 33:3012-3020.

Stiber, N.A.; M.J. Small; and M. Pantazidou, 2004. Site-specific Updating and Aggregation of Bayesian Belief Network Models for Multiple Experts. *Risk Analysis* 24(6):1529-1538.

Teigen, K. H., 1994. Variants Of Subjective Probabilities: Concepts, Norms, And Biases. In: G. Wright and P. Ayton, eds., *Subjective probability*, Wiley, New York.. pp.211-238.

Thompson, K.M. and D.L. Bloom, 2000. Communication of Risk Assessment Information to Risk Managers. *J. Risk Research* 3(4):333-352.

Titus, J.G., and V. Narayanan, 1996. The Risk of Sea Level Rise: A Delphic Monte Carlo analysis in which twenty researchers specify subjective probability distributions for model coefficients within their areas of experience, *Climatic Change*, 33:151-212, 1996.

Trauth, K. M., S.C. Hora, and R.V. Guzowski, 1993. Expert Judgment on Markers to Deter Inadvertent Human Intrusion into the Waste Isolation Pilot Plant, Report SAND92-1382, Sandia National Laboratories, Albuquerque, NM.

Tufte, E.R., 1983. *The Visual Display of Quantitative Information*. Graphics Press, Cheshire, CT.

Tversky, A. and D. Kahneman, 1974. Judgments Under Uncertainty: Heuristics and Biases. *Science* 185:1124-1131.

U.S. Department of Transportation / Federal Railroad Administration, 2003. Human Reliability Analysis in Support of Risk Assessment for Positive Train Control, Chapter 2, Approach to Estimation of Human Reliability in Train Control System Studies. DOT/FRA/ORD-03/15. Office of Research and Development, Washington, DC.

U.S. Environmental Protection Agency (USEPA), 1991. Risk Assessment Guidance for Superfund: Volume I – Human Health Evaluation Manual, Supplement to Part A: Community Involvement in Superfund Risk Assessments.

USEPA, 1992. Guidance on Risk Characterization for Risk Managers and Risk Assessors, Signed by Deputy Administrator F. Henry Habicht II. February 26, 1992.

USEPA, 1994. Seven Cardinal Rules of Risk Communication. EPA/OPA/87/020. Office of Policy Analysis, Washington, D.C.

USEPA, 1995a. Policy for Risk Characterization, Signed by Administrator Carol M. Browner. March 21, 1995.

USEPA, 1995b. The Probability of Sea Level Rise, Office of Policy, Planning and Evaluation. Washington, DC, EPA 230-R-95-008.

USEPA, 1997. Guiding Principles for Monte Carlo Analysis. EPA/630/R-97/001. Risk Assessment Forum, Office of Research and Development, Washington, DC. (Available at: <http://cfpub.epa.gov/ncea/raf/recordisplay.cfm?deid=29596>).

USEPA, 1998. Superfund Community Involvement Handbook and Toolkit. EPA 540-R-98-007. Office of Emergency and Remedial Response, Washington, D.C.

USEPA, 1999a. Report of the Workshop on Selecting Input Distributions for Probabilistic Assessments, Risk Assessment Forum, Washington, D.C. EPA/630/R-98/004.

USEPA, 1999b. Risk Assessment Guidance for Superfund: Volume I – Human Health Evaluation Manual, Supplement to Part A: Community Involvement in Superfund Risk Assessments. EPA/540/R-98/042, Washington, D.C. (Available at: http://www.cepis.ops-oms.org/tutorial6/fulltext/ci_ra.pdf).

USEPA, 1999c. Superfund Risk Assessment and How You Can Help: An Overview Videotape. September 1999 (English version) and August 2000 (Spanish version). English Version: EPA-540-V-99-003, OSWER Directive No. 93285.7.29B. Spanish Version (northern Mexican); EPA-540-V000-001, OSWER Directive No. 9285.7-40. Available Through NSCEP: (800) .490-9198 or (613) 489-8190.

USEPA, 2000a. Risk Characterization Handbook. EPA 100-B-00-002. Science Policy Council, Office of Science Policy, Office of Research and Development, Washington, DC. (Available at: <http://www.epa.gov/osa/spc/pdfs/rhandbk.pdf>).

USEPA, 2000b. Guidelines for Preparing Economic Analyses. EPA 240-R-00-003. Office of the Administrator, Washington D.C. (Available at: <http://yosemite.epa.gov/ee/epa/eed.nsf/webpages/Guidelines.html>).

USEPA, 2000c. Options for Development of Parametric Probability Distributions for Exposure Factors. EPA/600/R-00/058. National Center for Environmental Assessment, Office of Research and Development Washington, D.C. (Available at: <http://www.epa.gov/NCEA/pdfs/paramprob4ef/chap1.pdf>).

USEPA, 2001a. Risk Assessment Guidance for Superfund: Volume III – Part A. Process for Conducting Probabilistic risk Assessment. EPA 540-R-02-002, OSWER 9285.7-45. Office of Solid Waste and Emergency Response, Washington, D.C. (Available at: <http://www.epa.gov/oswer/riskassessment/rags3adt/index.htm>).

USEPA, 2001b. Early and Meaningful Community Involvement. Directive No. 90230.0-99, October 12. Office of Solid Waste and Emergency Response, Washington, D.C. (Available at: <http://www.epa.gov/superfund/resources/early.pdf>).

USEPA, 2002. Guidelines for Ensuring and Maximizing the Quality, Objectivity, Utility, and Integrity of Information Disseminated by the Environmental Protection Agency. (available at: <http://www.epa.gov/oei/qualityguidelines>)

USEPA, 2002b. Risk Communication in Action: Environmental Case Studies. Washington, D.C.. (available at: <http://www.epa.gov/ordntrnt/ORD/NRMRL/pubs/625r02011/625r02011.htm>).

USEPA, 2004. Final Regulatory Analysis: Control of Emissions from Nonroad Diesel Engines. EPA 420-R-04-007. Office of Transportation and Air Quality, Washington, D.C. (Available at: <http://www.epa.gov/nonroad-diesel/2004fr/420r04007a.pdf>).

USEPA, 2005. Guidelines for Carcinogen Risk Assessment. EPA/630/P-03-001B. Risk Assessment Forum, Washington, D.C. (Available at: <http://cfpub.epa.gov/ncea/cfm/recordisplay.cfm?deid=116283>).

USEPA, 2006. Peer Review Handbook, 3rd Edition, Review Draft. Washington, D.C. (available at: <http://intranet.epa.gov/ospintra/scipol/prhndbk05.doc>).

U.S. Nuclear Regulatory Commission (USNRC), 1996. Branch Technical Position on the Use of Expert Elicitation in the High-Level Radioactive Waste Program. NUREG-1563. Division of Waste Management, Office of Nuclear Material Safety and Standards, Washington, D.C.
USNRC, 1997. Recommendations for Probabilistic Seismic Hazard Analysis: Guidance on Uncertainty and Expert Use, prepared by the Senior Seismic Hazard Analysis Committee, NUREG/CR-6372, UCRL-ID-122160, Vol. 1 and 2, Washington, DC.

U.S. Office of Management and Budget (USOMB), 2002. Guidelines for Ensuring and Maximizing the Quality, Objectivity, Utility, and Integrity of Information Disseminated by Federal Agencies, Federal Register, February 22, 2002, Vol. 67, No. 36: 8452-8460. (Available at: <http://www.whitehouse.gov/OMB/fedreg/reproducible2.pdf>).

USOMB, 2003c. *OMB Circular A-4*, Regulatory Analysis, To the Heads of Executive Agencies and Establishments, Office of Information and Regulatory Affairs, Washington, DC. September 17, 2003. (available at: <http://www.whitehouse.gov/omb/circulars/a004/a-4.pdf>).

USOMB, 2006. Proposed Risk Assessment Bulletin, Office of Management and Budget, Office of Information and Regulatory Affairs, January 9. (Available at: http://www.whitehouse.gov/omb/inforeg/proposed_risk_assessment_bulletin_010906.pdf).

V. D. Fels-Klerx, H.J. Ine. L. H. J. Goossens, H. W. Saatkamp and S. H. S. Horst, 2002. Elicitation of Quantitative Data from a Heterogeneous Expert Panel: Formal Process and Application in Animal Health. *Risk Analysis* 22:67-81.

Walker K.D., P. Catalano, J.K. Hammitt, J.S. Evans, 2003. Use Of Expert Judgment In Exposure Assessment: Part 2. Calibration Of Expert Judgments About Personal Exposures To Benzene. *Journal of Exposure Analysis and Environmental Epidemiology*. 13:1-16.

Walker K.D., J.S. Evans, D. MacIntosh, 2001. Use Of Expert Judgment In Exposure Assessment. Part I. Characterization Of Personal Exposure To Benzene. *Journal of Exposure Analysis and Environmental Epidemiology* 11:308-322.

Walker, K.D., 2004. Memo: Appropriate Number of Experts for the PM EJ Project. Memo to Jim Neumann, Henry Roman, and Tyra Gettleman, IEC, November 11.

Wallsten, T.S. and D.V. Budescu, 1983. Encoding Subjective Probabilities: A Psychological and Psychometric Review. *Management Science* 29(2):151-173.

Wallsten, T.S., B.H. Forsyth, and D.V. Budescu, 1983. Stability and Coherence of Health Experts' Upper and Lower Subjective Probabilities about Dose-Response Functions. *Organizational Behavior and Human Performance* 31:277-302.

Wallsten, T.S. and R.G. Whitfield, 1986. Assessing the Risks to Young Children of Three Effects Associated with Elevated Blood Levels, ANL/AA-32, Argonne National Laboratory, Argonne, IL.

Wallsten, T.S., D.V. Budescu, I. Erev, and A. Diedrich, 1997. Evaluating and Combining Subjective Probability Estimates. *Journal of Behavioral Decision Modeling* 10:243-268.

Warren-Hicks, W.J. and D. R. J. Moore, (eds.), 1998. Uncertainty Analysis in Ecological Risk Assessment. SETAC Press, Pensacola, FL.

Whitfield, R.G. and T.S. Wallsten, 1989. A Risk Assessment for Selected Lead-Induced Health Effects: An Example of a General Methodology. *Risk Analysis* 9(2):197-208.

Whitfield, R.G., T.S. Wallsten, R.L. Winkler, H.M. Richmond, S.R. Hayes, and A.S. Rosenbaum, 1991. Assessing the Risk of Chronic Lung Injury Attributable to Ozone Exposure, Argonne National Laboratory, Report No. ANL/EAIS-2 (July).

Wilson, J. D., 1998. Default and Inference Options: Use in Recurrent and Ordinary Risk Decisions Discussion Paper 98-17, February, Resources for the Future, Washington, D.C.

Winkler, R.; T.S. Wallsten, R.G. Whitfield, H.M. Richmond and A. Rosenbaum, 1995. An Assessment of the Risk of Chronic Lung Injury Attributable to Long-Term Ozone Exposure. *Operations Research* 43(1):19-28.

Supplemental References Regarding Risk Communication and Public Perception.

Covello, V.T. 1987. Decision Analysis and Risk Management Decision Making: Issues and Methods. *Risk Anal* 7(2):131-139.

Fischhoff, B. 1995: Risk Perception and Communication Unplugged: Twenty Years of Process. *Risk Anal.* 15(2):137-145.

Fischhoff, B. 1998. Communicate unto others. *Reliab. Eng. Syst. Saf.* 59:63-72.

Hora, S.C. 1992. Acquisition of Expert Judgment: Examples form Risk Assessment. *J. Energy Eng.* 118(2):136-148.

Ibrekk, H. and M. G. Morgan. 1987. Graphical Communication of Uncertain Quantities to Non-Technical People. *Risk Anal.* 7:519-529.

Johnson, B.B. and P. Slovic, 1995. Presenting Uncertainty in Health Risk Assessment: Initial Studies of its Effects on Risk Perception and Trust. *Risk. Anal.* 15:485-494.

Kaplan. S. 1992. "Expert Information" versus "Expert Opinions." Another Approach to the Problem of Eliciting /Combining /Using Expert Knowledge in Probabilistic Risk Analysis. *Reliab. Eng. Syst. Saf* 35:61-72.

Morgan, M.G., A. Bostum, L. Lave, and C. J. Atman. 1992. Communicating Risk to the Public. *Environ. Sci. Technol.* 26(11):2048-2056.

Ohanian. E.V., J. A. Moore, J.R. Fowle, etc. Workshop Overview. 1997. Risk Characterization: A Bridge to Informed Decision Making. *Fundam. Appl. Toxicol.* 39:81-88.

Thompson. K.M. and D.L. Bloom. 2000. Communication of Risk Assessment Information to Risk Managers. *J. Risk. Res.* 3(4):333-352.

APPENDIX: FACTORS TO CONSIDER WHEN MAKING PROBABILITY JUDGMENTS

Introduction

Uncertainty is often associated with conclusions that we draw from research and more generally in our everyday thinking. When the data desired to support a particular decision do not yet exist, are sparse, of poor quality, or of questionable relevance to the problem at hand, subjective judgment comes into play. Formal elicitation of subjective judgments, often conducted with experts in the particular field, attempts to integrate what is known with what is not known about a particular quantity into a comprehensive, probabilistic characterization of uncertainty.

Many sources often contribute to uncertainty about any given issue, and it is generally difficult for most people to consider and integrate them all. When an expert is asked to make probability judgments on socially important matters, it is particularly important that he or she consider the relevant evidence in a systematic and effective manner and provide judgments that represent his or her opinions well.

Several academic traditions – decision analysis, human factor cognitive sciences, experimental psychology, and expert systems analysis – have sought to understand how to elicit probabilistic judgments from both lay people and experts in a reliable way. Researchers have amassed a considerable amount of data concerning the way people form and express probabilistic judgments. The evidence suggests that, when considering large amounts of complex information, most people use heuristics (i.e., simplifying rules) and demonstrate certain cognitive biases (i.e., systematic distortions of thought). These heuristics and biases can lead to systematic biases in the judgments and errors of over-and under-confidence. In particular, many studies indicate that experts and lay people alike tend to be overconfident. In probabilistic assessments of uncertainty, overconfidence manifests itself as placing higher probabilities on being correct (or narrower confidence intervals around a prediction) than measures of performance ultimately warrants. Such errors in judgments may have important implications for decisions that depend on them.

The purpose of this paper is to make you aware of these heuristics and biases. We

will first review the most widespread heuristics and biases, and then offer some suggestions to help you mitigate their effects.

Heuristics and Biases Involved in Expert Judgment

Sequential Consideration of Information

Generally, the order in which evidence is considered influences the final judgment, although logically that should not be the case. Of necessity, pieces of information are considered one by one in a sequential fashion. However, those considered first and last tend to dominate judgment. In part, initial information has undue influence because it provides the framework that subsequent information is then tailored to fit. For example, people usually search for evidence to confirm their initial hypotheses; they rarely look for evidence that weighs against them. The latter evidence has an undue effect simply because it is fresher in memory.

Related to these sequential effects is the phenomenon of “Anchoring and Adjustment.” Based on early partial information, individuals typically form an initial probability estimate, the “anchor”, regarding the event in question. They then make adjustments to this judgment as they consider subsequent information. Such adjustments tend to be too small. In other words, too little weight is attached to information considered subsequent to the formation of the initial judgment.

Effects of Memory on Judgment

It is difficult for most people to conceptualize and make judgments about large, abstract universes or populations. A natural tendency is to recall specific members and then to consider them as representative of the population as a whole. However, the specific instances often are recalled precisely because they stand out in some way, such as being familiar, unusual, especially concrete, or of personal significance. Unfortunately, the specific characteristics of these singular examples are then attributed, often incorrectly, to all the members of the population of interest. Moreover, these memory effects are often combined with the sequential phenomena discussed earlier. For example, in considering evidence regarding the relationship between changes in ambient PM_{2.5} and premature mortality, you might naturally think first of a study you recently read or one that was unusual and therefore stands out. The tendency might then be to

treat the recalled studies as typical of the population of relevant research and ignore important differences among studies. Subsequent attempts to recall information could result in thinking primarily of evidence consistent with the initial items considered.

Estimating Reliability of Information

People tend to overestimate the reliability of information, ignoring factors such as sampling error and imprecision of measurement. Rather, they summarize evidence in terms of simple and definite conclusions, which causes them to be overconfident in their judgments. This tendency is stronger when one has a considerable amount of intellectual and/or personal involvement in a particular field. In such cases, information is often interpreted in a way that is consistent with one's beliefs and expectations, results are overgeneralized, and contradictory evidence is ignored or underestimated.

Relation between Event Importance and Probability

Sometimes the importance of events, or their possible costs or benefits, influences judgments about the uncertainty of the events when rationally, importance should not affect probability. In other words, one's attitudes towards risk tend to affect one's ability to make accurate probability judgments. For example, many physicians tend to overestimate the probability of very severe diseases because they feel it is important to detect and treat them; similarly, many smokers underestimate the probability of adverse consequences of smoking because they feel that the odds do not apply to themselves personally.

Assessment of Probabilities

Another limitation is related to one's ability to discriminate between levels of uncertainty and to use the appropriate criteria of discrimination for different ranges of probability. One result is that people tend to assess both extreme and midrange probabilities in the same fashion, usually doing a poor job in the extremes. It is important to realize that the closer to the extremes (either 0 or 1) that one is assessing probabilities, the greater the impact of small changes. It helps here to think in terms of odds as well as probabilities. Thus, for example changing a probability by 0.009 from 0.510 to 0.501 leaves the odds almost unchanged, but the same change from 0.999 to 0.990 changes the odds by a factor of about 10 from 999:1 to 99:1.

Recommendations

Although extensive and careful training would be necessary to eliminate all the problems mentioned above, some relatively simple suggestions can help minimize them. Most important is to be aware of natural cognitive biases and to try consciously to avoid them.

To avoid sequential effects, keep in mind that the order in which you think of information should not influence your final judgment. It may be helpful to actually note on paper the important facts you are considering and then to reconsider them in two or more sequences, checking the consistency of your judgments. Try to keep an open mind until you have gone through all of the evidence, and don't let the early information you consider sway you more than is appropriate.

To avoid adverse memory effects, define various classes of information that you deem relevant, and then search your memory for examples of each. Don't restrict your thinking only to items that stand out for specific reasons. Make a special attempt to consider conflicting evidence and to think of data that may be inconsistent with a particular theory. Also, be careful to concentrate on the give probability judgment, and do not let your own values (how you would make the decision yourself) affect those judgments.

To accurately estimate the reliability of information, pay attention to such matters as sample size and the power of the statistical tests. Keep in mind that data are probabilistic in nature, subject to elements of random error, imprecise measurements, and subjective evaluation and interpretation. In addition, the farther one must extrapolate, or generalize, from a particular study to a situation of interest, the less reliable is the conclusion may be and the less certainty should be attributed to it. Rely more heavily on information that you consider more reliable, but do not treat it as absolute truth.

Keep in mind that the importance of an event or an outcome should not influence its judged probability. It is rational to let the costliness or severity of outcome influence the point at which action is taken with respect to it, but not the judgment that is made about the outcome's likelihood. Finally, in making probability judgments, think primarily in terms of the measure (probability or odds) with which you feel more

comfortable, but sometimes translate to the alternative scale, or even to measures of other events (e.g., the probability of the event not happening). When estimating very small or very large likelihoods, it is usually best to think in terms of odds, which are unbounded, instead of probabilities, which are bounded. For example, one can more easily conceptualize odds of 1:199 than a probability of 0.005.

Appendix B: Glossary

Bayesian Analysis: Statistical analysis that describes the probability of an event as the degree of belief or confidence that a person has, given some state of knowledge, that the event will occur. Bayesian Monte Carlo combines a prior probability distribution and a likelihood function to yield a posterior distribution. Also called subjective view of probability, in contrast to the frequentist view of probability.

Expert Elicitation: A systematic process of formalizing and quantifying, in terms of probabilities, experts' judgments about uncertain quantities, events, or relationships.

Expert Judgment: An inferential opinion of a specialist or group of specialists within an area of their expertise. Expert judgment (alternatively referred to as professional judgment) may be based on an assessment of data, assumptions, criteria, models and parameters in response to questions posed in the relevant area of expertise.

Frequentist: A term referring to classical statistics in which the probability of an event occurring is defined as the frequency of occurrence measured in an observed series of repeated trials.

Likelihood Function: A term from Bayesian statistics referring to a probability distribution that expresses the probability of observing new information given that a particular belief is true.

Stochastic Process: A process involving random variables, and characterized by variability in space or time.

Uncertainty: Lack of knowledge about specific variables, parameters, models or other factors. Examples include limited data regarding concentrations of environmental contaminants. Some forms of uncertainty may be reduced through further study.

Variability: True heterogeneity or diversity in characteristics among members of a population or one individual over time.
