

1 Title: Endocrine Disruption: Where have we been, interpretation of data, and lessons learned  
2 from Tier 1

3

4 Running head: Lessons learned from Tier 1

5 Jane P. Staveley <sup>\*†</sup>, Leslie W. Touart <sup>‡</sup>, Keith Solomon <sup>§</sup>, Ellen Mihaich <sup>¶</sup>, Amy Blankinship <sup>#</sup>, and  
6 Gerald Ankley <sup>††</sup>

7 <sup>†</sup>Exponent, 1000 Centre Green Way, Suite 200, Cary, NC, USA, 27513; Telephone: 919-228-

8 6480; Fax: 919-228-6501; e-mail: [jstaveley@exponent.com](mailto:jstaveley@exponent.com); <sup>‡</sup> USEPA, U.S. EPA, Office of

9 Chemical Safety and Pollution Prevention, Washington, DC, USA, [currently Equiparent

10 Consulting, Woodbridge, VA, USA, [les.touart@equiparentconsulting.com](mailto:les.touart@equiparentconsulting.com)]; <sup>§</sup>Centre for

11 Toxicology, University of Guelph, Guelph, Ontario, Canada, [ksolomon@uoguelph.ca](mailto:ksolomon@uoguelph.ca);

12 <sup>¶</sup>Environmental and Regulatory Resources, Durham, NC, USA, [emihaich@nc.rr.com](mailto:emihaich@nc.rr.com); <sup>#</sup>U.S. EPA,

13 Office of Pesticide Programs, Washington, DC, USA, [blankinship.amy@epa.gov](mailto:blankinship.amy@epa.gov); <sup>††</sup>U.S. EPA,

14 Office of Research and Development, Duluth, MN, USA, [ankley.gerald@epa.gov](mailto:ankley.gerald@epa.gov)

15

16

17 \* To whom correspondence may be addressed

18 Corresponding Author:

19 Jane Staveley

20 1000 Centre Green Way, Suite 200, Cary, NC, USA 27513.

21 Email address: [jstaveley@exponent.com](mailto:jstaveley@exponent.com)

22 **ABSTRACT**

23 In response to the requirements of the US EPA's Endocrine Disruptor Screening Program, Tier 1  
24 assays have been performed with a number of pesticides over the past several years. These  
25 assays are designed to be used in concert as a screen for potential interactions with vertebrate  
26 estrogen, androgen, and thyroid systems. The results of the 11 assays in the Tier 1 battery are  
27 then used, along with other lines of evidence, to determine whether a chemical is endocrine-  
28 active and, as a consequence, might be a candidate for Tier 2 testing. An overview of the Tier-1  
29 testing program was presented in Session Two of the Society of Environmental Toxicology and  
30 Chemistry (SETAC) North America Focused Topic Meeting: *Endocrine Disruption Chemical*  
31 *Testing: Risk Assessment Approaches and Implications* (February 4 – 6, 2014). Subsequent  
32 presentations discussed the concept of weight-of-evidence (WoE) and assessment of Tier 1  
33 results in a WoE framework. The importance of scientifically credible, transparent approaches  
34 for conducting WoE analyses was recognized, and approaches for framing the hypotheses,  
35 evaluating the data, assigning weight to different endpoints relative to their diagnostic  
36 effectiveness, and assessing confounding factors were presented. In recognition of the cross-  
37 species conservation of the hypothalamic-pituitary-gonadal axis among vertebrates, a subset of  
38 the Tier-1 *in vivo* assays may be useful for more rapidly screening chemicals for potential  
39 endocrine activity.

40

41 Keywords: Endocrine disruption, Endocrine Disruption Screening Program, testing, weight-of-  
42 evidence

## 43 INTRODUCTION

44 Session Two of the Society of Environmental Toxicology and Chemistry (SETAC) North  
45 America Focused Topic Meeting: *Endocrine Disruption Chemical Testing: Risk Assessment*  
46 *Approaches and Implications* (February 4 – 6, 2014) focused on the experience gained to date  
47 with implementation of the Tier 1 testing of U.S. EPA’s Endocrine Disruptor Screening Program  
48 (EDSP), and how these data can be used to make decisions about the need for further testing.  
49 Leslie Touart presented an overview of the 11 assays in the Tier 1 screening battery. Keith  
50 Solomon discussed the concept of using weight-of-evidence (WoE) in risk assessment, illustrated  
51 by an example on the potential effects of atrazine on fish, amphibians, and reptiles. Ellen  
52 Mihaich described a hypothesis-based weight of evidence framework that was developed to  
53 evaluate experimental data, with a proposed specific use in evaluating results of the Tier 1  
54 screening battery. Amy Blankinship provided an overview of the conceptual basis of the WoE  
55 guidance used by the USEPA to evaluate Tier 1 data for identifying the need for additional (Tier-  
56 2) testing. The session concluded with a presentation by Gary Ankley on an analysis indicating  
57 that it appears possible to use just two of the current Tier-1 tests as initial “gate keeper” assays,  
58 following which chemicals may be exempted from further testing or subjected to additional,  
59 confirmatory analyses with other existing Tier-1 assays.

## 60 SESSION PRESENTATION SUMMARIES

61 *USEPA Endocrine Disruptor Screening Program (EDSP) Tier-1 Battery Overview by: Leslie*  
62 *Touart*

63 The suite of 11 Tier-1 EDSP assays is specifically designed to detect chemicals with the  
64 potential to interact with the estrogen, androgen, and thyroid (EAT) systems in vertebrates,

65 through mechanisms such as activation and antagonism of target nuclear hormone receptors, and  
66 inhibition of hormone synthesis (<http://www.epa.gov/endo/>). Given the complex interactive  
67 nature of the endocrine system, if the objective is to comprehensively detect their potential to  
68 disrupt endocrine regulated processes, it is clear that chemicals should be tested for apical effects  
69 (e.g., the ability to alter growth, development, or reproductive processes) and their potency in *in*  
70 *vitro* assays of receptors and synthesis of sex steroids. A battery of screening tests has been  
71 developed which includes a range of taxonomic groups and sufficient diversity of endpoints to  
72 maximize sensitivity and minimize false negatives. There are five *in vitro* assays focused on  
73 binding to and transactivation of the estrogen receptor, binding to the androgen receptor, and  
74 inhibition of synthesis of sex steroids. There are six *in vivo* Tier-1 screens, four utilizing rats  
75 (uterotrophic and Hershberger assays; male and female pubertal assays), one with the fathead  
76 minnow (fish short-term reproduction assay; FSTRA), and one with the amphibian *Xenopus*  
77 *laevis* (amphibian metamorphosis assay; AMA). Although each of the Tier-1 assays provides  
78 unique data, the suite was purposefully designed to result in some redundancy with respect to  
79 detecting endocrine pathways of concern (Table 1). The *in vitro* assays provide sensitivity and  
80 mechanistic clues, while the *in vivo* assays provide for integrative responses and metabolism and  
81 distribution considerations. The results of the Tier-1 battery are to be interpreted in a WoE  
82 context, rather than the sum of positive and negative assays. Some endpoints are more  
83 diagnostic/specific than others, and effects seen in multiple endpoints and multiple assays carry  
84 the most weight. There are two possible interpretations of the outcome of the Tier-1 battery:  
85 either the potential for EAT activity exists, which warrants analysis in Tier-2 testing, or there is  
86 low or no potential for EAT activity. A FIFRA Science Advisory Panel meeting held in 2008 to  
87 review the Tier-1 screening battery concluded that, based on the state of the science at the time,

88 the set of assays were an appropriate starting point to detect potential endocrine disruptors and  
89 should continue to be refined and developed. In summary, multiple assays are required to  
90 comprehensively screen endocrine, androgen, and thyroid hormone systems. The *in vitro* assays  
91 are suitable for well-understood mechanisms (e.g., receptor binding), while the *in vivo* assays  
92 with intact hypothalamic-pituitary-gonadal (HPG) and hypothalamic-pituitary-thyroidal (HPT)  
93 axes are useful for efficiently screening complex processes. The totality of the results of the Tier-  
94 1 screening battery are needed to support WoE conclusions about the potential of a chemical to  
95 interact with vertebrate EAT systems.

96

97 *Use of weight of evidence for characterizing adverse outcome pathways in risk assessment by:*

98 ***Keith Solomon***

99 Information and data on chemicals from studies published in the open literature are  
100 increasingly being used for assessment purposes by regulatory agencies in many jurisdictions,  
101 including North America and Europe. Because most of these studies are not conducted to the  
102 Good Laboratory Practices (GLP) standards as required by regulatory agencies, there is a need to  
103 assess their quality and relevance in light of the regulatory endpoints being considered. To aid in  
104 interpretation and to use these data in regulatory decision-making, they need to be integrated into  
105 lines of evidence that inform adverse outcome pathways (AOPs) and lines of evidence related to  
106 apical endpoints such as survival, growth, development, and reproduction.

107 There are important differences between studies published in the open literature and those  
108 conducted under GLP guidelines for regulatory agencies. Published studies often are  
109 incompletely documented, raw data are rarely available, and many studies do not follow  
110 standardized protocols. In addition, studies used in reviews and meta-analyses may be subjected

111 to selection bias or, in a worse case, there may be selection bias where negative (no effect)  
112 results are not published (Walker et al. 2008). In contrast, studies conducted under GLP with  
113 Quality Assurance and Quality Control (QA/QC) are required by regulation in many  
114 jurisdictions, are completely documented, the raw data are available, most studies are conducted  
115 using standardized protocols, and there is no publication bias; all observations are documented.  
116 For this reason, GLP studies with QA/QC cannot always be directly compared to or combined  
117 with published studies for a WoE analysis for decision making. A WoE analysis of the potential  
118 effects of atrazine on fish, amphibians, and reptiles (Van Der Kraak et al. 2014) was conducted  
119 using quantitative methods to characterize the strength and relevance of published and GLP  
120 studies. This brief overview describes a subset of data taken from Van Der Kraak et al. (2014)  
121 with a specific focus on reproductive outcomes in fish, amphibians, and reptiles.

122 In this example, the strength of the experimental methods and the ecological relevance of  
123 the observed responses from over 2000 studies and experiments were scored. The detailed  
124 methods of scoring are reported (Van Der Kraak et al. 2014) and are not repeated here. Briefly,  
125 the strength of the methods was scored based on various aspects of the studies, such as the  
126 experimental design and conduct, the use of appropriate controls, measures of exposures, the  
127 inclusion of environmentally realistic concentrations, number of concentrations, quality control,  
128 and transparency of data. These criteria are similar to those suggested by Klimisch et al. (1997).  
129 The relevance of the each response was assessed by scoring statistical significance, concentration  
130 or dose-response, its relevance to an appropriate apical endpoint, and a biologically plausible  
131 mechanism. The WoE process was inclusive and no studies were excluded, except those with  
132 mixtures where the individual components were not tested individually. Results were presented  
133 graphically where strength and relevance were shown separately for easy interpretation and are

134 supported with details of the experimental procedures (see SI provided with Van Der Kraak et al.  
135 2014).

136 AOPs (Ankley et al. 2010) are used to characterize links between responses at lower  
137 levels of biological organization and apical endpoints such as survival, growth, development, and  
138 reproduction (Figure 1). AOPs provide the framework for extrapolation of effects to other  
139 organisms/taxa or to identify reliable and robust biomarkers that can be used in place of the  
140 apical endpoint. Responses from multiple studies at each level of an AOP can be subjected to  
141 WoE analysis. If one or more apical endpoints (4 and 5 in Figure 1) have been characterized  
142 under WoE, and the combination of these indicates no or *de minimis* effects at environmentally  
143 relevant exposures, an analysis of AOPs is not needed. In this case any effects observed at lower  
144 levels of organization are “trumped” or negated by lack of effect on apical endpoints and those at  
145 lower levels are most likely only bioindicators of exposure or adaptive response. However, if  
146 one or more of the apical endpoints indicates relevant effects at environmentally-relevant  
147 exposures, then a characterization of AOP might be useful to better understand the response.  
148 Because responses in an AOP are concatenated, a break in the chain at any point in the pathway  
149 (illustrated by the red X in Figure 1) provides evidence that the responses are not important for  
150 apical effects and that regulatory action would not be needed.

151 To illustrate the combination of AOPs with WoE analysis, reproductive responses to  
152 atrazine in fish, amphibians, and reptiles were combined in graphs showing the mean scores for  
153 multiple responses and their uncertainty (see Van Der Kraak et al. (2014) for details) in four  
154 links of an AOP. These links in the AOP chain were at the biochemical (A), cellular (B), organ  
155 (C), and organism (D) levels (Figure 2). The organismal level is apical. As can be seen from the  
156 graphics in Figure 2 (A to D), the mean values for relevance of all the responses in the AOP

157 chain cluster at the low end of the relevance scale. The means and uncertainty of the scores  
158 provide the basis for testing risk hypotheses in the WoE framework. These analyses suggested  
159 that there was a *de minimis* risk of adverse effects at all levels of the AOP. Strictly speaking, the  
160 lack of effects at the organismal level would negate the need for AOP analysis but the example is  
161 illustrative of the robustness of the response as effects at all levels of the AOP are of low  
162 relevance. This provides greater assurance that the lack of response is real and not just due to a  
163 lack of data or measures at different levels of organization. As is indicated by the error-bars  
164 (Figure 2), there was less uncertainty in the scores for relevance than the scores for strength. The  
165 scores from strength for these responses (see details in Van Der Kraak et al. 2014) ranged from  
166 low to high but the high-strength scores were consistent in indicating very low or *de minimis*  
167 relevance.

168 In conclusion, the use of a formal, well described, transparent, and quantitative process  
169 for WoE provides a helpful tool for conducting risk assessment. It is more objective and, when  
170 combined with analysis of AOPs, provides more clarity and understanding of the significance of  
171 effects. The example provided is directed specifically to reproduction but the process is  
172 applicable to areas other than risk assessment; however, different and response-specific methods  
173 of scoring may be needed.

174

175 ***"Weighing" the Evidence: Relevance and Transparency in the Evaluation of Endocrine***

176 ***Activity by: Ellen Mihaich***

177 A comprehensive, hypothesis-based weight of evidence (HB-WoE) framework was  
178 developed to be applicable to any determination relying on experimental data, with a proposed  
179 specific formulation for evaluating results of the U.S. EPA's Tier-1 Endocrine Screening Battery



180 (ESB) (Borgert et al. 2011a). The framework requires that before any WoE determinations are  
181 considered, each experimental endpoint be weighted according to its relevance for deciding each  
182 of 8 hypothesis addressed by the ESB. These hypotheses test whether or not the chemical under  
183 evaluation has the potential to act as an (anti)-estrogen, (anti)-androgen, (anti)-thyroid, or induce  
184 or inhibit steroidogenesis. The purpose of an *a priori* relevance weighting is to ensure a level of  
185 transparency and objectivity exceeding that possible from WoE processes claiming a basis in  
186 professional judgment alone. Ideally, quantitative relevance weighting (Wrel) values would be  
187 derived from data revealing the positive and negative predictive value of the various endpoints  
188 for the hypotheses addressed by the ESB assays. Because the ESB assays have not been  
189 validated to that level (Borgert et al. 2011b), obviating the derivation of quantitative Wrels, this  
190 method provides for endpoints to be ranked according to 4 categories based on interpretations of  
191 relevant literature (Borgert et al. 2014). Although these Wrel rankings necessarily involve  
192 professional judgment, their *a priori* derivation based on a defined rationale (Borgert et al. 2014)  
193 enhances transparency nonetheless and renders any WoE determinations based on them  
194 amenable to methodological scrutiny according to basic scientific premises. To make WoE  
195 determinations for a particular substance, the framework requires combining Wrel  
196 values/rankings for each hypothesis with response weightings (Wres) derived from the ESB data.

197 The method has been more fully described by Borgert et al. (2014). Wrels were  
198 determined by ranking the endpoints by hypothesis according to the following definitions below.  
199 Although no hypothesis can be decided on the results of a single assay, “interpretable” means  
200 that the results for an endpoint provide information relevant to the hypothesis, without  
201 clarification from other endpoints. Whether a hypothesis is supported requires consideration of  
202 results from all relevant (#1, #2, #3) assays and endpoints. Rank 1 endpoints are typically *in vivo*

203 endpoints, specific & sensitive for the hypothesis and interpretable without other endpoints.  
204 Rank 2 includes many *in vitro* endpoints that are sensitive and specific, but less informative than  
205 Rank 1. Rank 3 includes many apical *in vivo* endpoints that are relevant for the hypothesis, but  
206 are only corroborative of Rank #1 and #2 endpoints. Rank 4 endpoints were considered not  
207 relevant for the hypothesis.

208 Data for the test chemicals are evaluated for each hypothesis individually, beginning with  
209 Rank 1 and continuing through Rank 3 endpoints. The response to Rank 1 endpoints guides the  
210 evaluation and interpretation of information from lower-ranked endpoints. Responses in Rank 1  
211 are a preliminary indication that the hypothesis is or is not supported. Rank 2 endpoints are then  
212 evaluated, with consistent positive responses among Rank 1 and 2 endpoints considered  
213 sufficient support, and consistent negative responses considered refutation of the hypothesis.  
214 Rank 3 endpoints are then consulted for consistency and, together with the strength of response  
215 (Wres) in Rank 1 and 2 endpoints, temper or strengthen the conclusion. The interpretation  
216 becomes more complex if Rank 2 endpoints are inconsistent with negative results in Rank 1  
217 endpoints. In this case, the strength of the response in Rank 2 endpoints becomes even more  
218 critical, as does an evaluation of Rank 3 endpoints, along with a consideration of the potential  
219 reasons that Rank 1 endpoints might not respond. Some overarching guidelines for interpretation  
220 can be established. Rank 1 endpoints cannot be dismissed for inconsistency with Rank 2. Rank  
221 3 endpoints, in contrast, provide little useful information other than as corroboration for findings  
222 in Ranks 1 and 2. Situations in which Rank 2 and 3 are consistent, but inconsistent with Rank 1  
223 endpoints present the greatest challenge, and no general statements can be made.

224 Published data from genistein was used to illustrate the application of this WoE  
225 framework and process for determining the potential for genistein to act as an estrogen agonist.  
226 Genistein is an isoflavone present in plant foods like soy, fava beans, and clover. Phytoestrogens  
227 like genistein are known to cause effects on reproduction in female ruminants, such as sheep and  
228 cattle (Adams 1995), and have been well studied to understand potential impacts on humans  
229 given the number of populations using a diet high in soy. For brevity, summary results are  
230 presented only for the estrogen agonist hypothesis in Table 2. In this example, although there are  
231 studies that provide some conflicting results (data not shown), the overall weight of the evidence  
232 of the data for genistein would support the estrogen agonist hypothesis. While few studies use  
233 positive controls because of animal use concerns, and specific positive controls would be needed  
234 to address each hypothesis being tested, some studies with genistein have employed compounds  
235 such as ethinyl estradiol (Kim et al. 2005) which allows for an estimation of estrogenic potency.  
236 Each additional hypothesis and the appropriately ranked endpoints would be considered  
237 separately; more detail on endpoint ranking can be found in Borgert et al. (2014).

238 This HB-WoE framework has been criticized for excessive detail, burdensome number  
239 and impossible requirements for quantitative rankings, and excessive time required to complete  
240 the process. As shown here, these criticisms are unfounded. The HB-WoE framework (Borgert  
241 et al. 2011a) provides a means for transparent, objective conclusions about ESB results, and  
242 moreover, streamlines the evaluation by allowing the analyst to appropriately allocate time and  
243 attention to the most definitive information. Although it is not yet possible to attain the goal of  
244 data-derived quantitative  $W_{rel}$  and  $W_{res}$  values, use of explicit  $W_{rel}$  rankings, derived *a priori*  
245 and applied similarly for each hypothesis, helps to ensure transparency and consistency, a feature  
246 absent from WoE approaches based solely on professional judgment. Despite an absence of

247 positive and negative control data in some ESB assays, Wres information can often be gleaned  
248 from Rank 1 and some Rank 2 endpoints, including an estimate of potency differences. The HB-  
249 WoE framework provides for efficient processing and interpretation of ESB data by considering  
250 the results of Rank 1 through 3 endpoints in consecutive order for each hypothesis. It provides  
251 for a systematic method for identifying and resolving inconsistencies in results from ESB and  
252 other scientifically relevant information and obviates a need to consider less definitive  
253 information unless it could help to resolve an ambiguous interpretation.

254

255 ***Weight of Evidence: Evaluating Results from Tier-1 Screening for the U.S. EPA Endocrine***  
256 ***Disruptor Screening Program by: Amy Blankinship***

257 In 2011, the United States Environmental Protection Agency's Office of Chemical Safety  
258 and Pollution Prevention (EPA/OCSPP) published a guidance document for the Endocrine  
259 Disruptor Screening Program (EDSP) which presented a weight of evidence (WoE) approach for  
260 evaluating Tier-1 screening data for identifying the need for additional (Tier-2) testing (USEPA  
261 2011). The function of the EDSP Tier-1 screening process is to identify chemicals that have the  
262 potential to interact with the estrogen, androgen, or thyroid (E, A, or T) pathways and evaluate  
263 the need for additional testing. The WoE guidance document provides general guidance in  
264 support of EPA efforts to integrate and interpret data submitted in response to orders for Tier-1  
265 screening; however, the guidance is not considered binding and reviewers may deviate from the  
266 guidance where circumstances warrant. As described in the guidance document, the WoE  
267 process identifies how the individual lines of evidence are assembled and integrated along two  
268 concepts (*i.e.*, complementarity and redundancy) within the conceptual framework of an adverse

269 outcome pathway (AOP). Broadly, there are four main steps outlined in the guidance which  
270 provide the foundation for WoE evaluations. The first step is to assemble and evaluate the  
271 individual studies for their scientific quality and relevance in evaluating potential endocrine  
272 interaction(s). The second step is to integrate the data along different levels of biological  
273 organization while examining the extent of concordance (robustness) of complementarity (*i.e.*,  
274 the concordance of endpoints within an assay that measures multiple endpoints) and redundancy  
275 (*i.e.*, the concordance of endpoints/responses across assays) in the observed responses across  
276 these different levels of biological organization. The third step is to then characterize the main  
277 lines of evidence as well as any conclusions. Finally, the last step is to evaluate whether  
278 additional testing is needed based on the evidence and conclusions described above.

279 As mentioned, the first step is to assemble and evaluate the available scientific data. Data  
280 for the EDSP Tier-1 WoE evaluation falls into one of two categories: 1) EDSP Tier-1 data, and  
281 2) other scientifically relevant information (OSRI). The EDSP Tier-1 data represent a battery of  
282 11 assays consisting of *in vitro* and mammalian and wildlife *in vivo* assays. OSRI may include  
283 published literature studies as well as studies conducted under USEPA (often referred to as Part-  
284 158 data) or OECD guidelines submitted in support of registration of pesticides or other  
285 chemicals. Each study is evaluated for scientific quality and relevance for informing interactions  
286 with the E, A, or T pathway. Additionally, the concordance or consistency (complementarity) of  
287 the responses in the individual study is evaluated. For the Tier-1 *in vivo* assays, often multiple  
288 endpoints are measured in each assay. Decision logic trees were developed for some Tier-1 *in*  
289 *vivo* assays in an effort to help guide the investigator/reviewer in interpreting results across  
290 multiple endpoints within an assay (Ankley and Jensen 2014; USEPA 2009). Evaluation of the  
291 potential confounding effects of overt toxicity in the study as well as the relative degree of

292 diagnostic utility of a specific endpoint for discerning whether or not the chemical has interacted  
293 with the endocrine system are considered. The collective response of the individual endpoints, as  
294 well as the conditions under which they were expressed, are considered when evaluating an  
295 overall indication of potential interaction as measured by the study.

296         The second step in this WoE process is to formulate hypotheses and integrate the  
297 available data along different levels of biological organization. Two key elements in the  
298 integration of data as well as characterizing the extent to which the available data support a  
299 hypothesis that a chemical has the potential to interact with E, A, or T pathways are the concepts  
300 of complementarity and redundancy. These two concepts provide a basis for considering the  
301 plausibility, coherence, strength, and consistency of the body of evidence. The current EDSP  
302 Tier-1 screening assays are meant to evaluate whether or not a chemical can interact with E, A  
303 and T consisting of different levels of biological organization from a molecular initiating event  
304 such as receptor binding through potential adverse effects in apical endpoints such as sexual  
305 development and fecundity at the whole organism level. Transitions to higher levels of biological  
306 organization can indirectly provide information on potential compensatory capabilities of an  
307 individual organism.

308         After the data have been assembled and integrated, the third step is to characterize the  
309 main lines of evidence along with the conclusions; this characterization involves three  
310 components. The first component is whether the data provide relevant, robust and consistent  
311 evidence in terms of complementarity and redundancy as well as biological plausibility. Second,  
312 is at what level of biological organization were the responses observed and whether organisms  
313 exhibit compensatory responses at higher-levels of biological organization. Finally, an

314 evaluation of under what conditions did the responses occur including discussions regarding  
315 whether the responses were observed in the presence of overt or systemic toxicity. The presence  
316 of overt and/or systemic toxicity introduces uncertainty in the ability to distinguish effects  
317 specifically related to an endocrine-mediated effect from a non-endocrine toxic response. This  
318 uncertainty in distinguishing whether the responses were endocrine-mediated was discussed at  
319 the FIFRA Scientific Advisory Panel (SAP) meeting in July 2013 that evaluated scientific issues  
320 associated with the WoE evaluation of the EDSP Tier-1 screening process. The SAP stated that ,  
321 *“In summary, the Panel agreed that little, if any, weight should be placed on signs of endocrine*  
322 *disruption in the presence of overt toxicity. All effects in endocrine sensitive tissues should be*  
323 *evaluated in terms of primary interactions with the endocrine system vs. secondary effects*  
324 *related to toxicity in non-endocrine organs or overall disruptions in homeostasis”* (Schlenk and  
325 Jenkins, 2013; Page 12; SAP 10/30/2013). Therefore, EPA considers multiple lines of evidence  
326 in including the observed responses in the Tier-1 assays and OSRI in the context of a chemical’s  
327 physical/chemical properties and its known modes of action in its overall characterization of a  
328 chemical’s potential to interact with the E, A or T pathway. Adequately addressing these three  
329 main questions is fundamental to the WoE process and in determining whether additional data  
330 are needed.

331 In addition to characterizing the WoE, reviewers also consider: 1) uncertainties and their  
332 potential impact to conclusions; 2) discussion of key studies; 3) description of inconsistent or  
333 conflicting data; 4) overall strength of evidence supporting a conclusion; and, 5) what, if any,  
334 additional data are needed and why. Assessing the need for additional data is based on a case-  
335 by-case analysis which will include integration of existing knowledge on the chemical including  
336 relevant hazard and exposure information. In summary, the evaluation of the EDSP Tier-1

337 screening process and ultimate decision for any additional testing is based on a totality of the  
338 scientific evidence.

339 ***Cross-Species Conservation of Endocrine Pathways Provides a Basis for Reevaluation of***  
340 ***EDSP Tiered Testing Paradigm: by Gerald Ankley***

341 Many structural and functional aspects of the HPG axis are known to be highly  
342 conserved, but the relative significance of this from a regulatory toxicology perspective has  
343 received comparatively little attention. High-quality data generated through development and  
344 validation of Tier-1 tests for the USEPA Endocrine Disruptor Screening Program (EDSP) offer a  
345 unique opportunity to compare responses of mammals versus fish to chemicals that may affect  
346 shared pathways within the HPG axis. The analysis described by Ankley and Gray (2013)  
347 focused on data generated with model chemicals that act (primarily) as estrogen receptor  
348 agonists (17 $\alpha$ -ethynylestradiol, methoxychlor, bisphenol A), androgen receptor agonists  
349 (methyltestosterone, 17 $\beta$ -trenbolone), androgen receptor antagonists (flutamide, vinclozolin,  
350 p,p'-DDE) or inhibitors of different steroidogenic enzymes (ketoconazole, fadrozole, fenarimol,  
351 prochloraz). All 12 chemicals had been tested in the EDSP fish short-term reproduction assay  
352 (FSTRA) and in one or more of the four *in vivo* Tier-1 screens with rats (Uterotrophic,  
353 Hershberger, male and female pubertal assays). In most cases there was high concordance  
354 between the fish and rat assays with respect to identifying chemicals that impacted specific HPG  
355 pathways of concern, with the test chemicals producing positive results in the fish and one or  
356 more of the rat assays. However, some assays were clearly superior to others in terms of  
357 detecting specific pathways; for example, the effects of inhibitors of steroid hormone synthesis  
358 were most obvious in the FSTRA, whereas the activity of androgen receptor antagonists were



359 clearest in the Hershberger and male pubertal assays. Based on this analysis it appears possible to  
360 use just two of the current Tier-1 tests, the FSTRA and the male pubertal assay, to ensure full  
361 coverage of HPG axis pathways of concern. Specifically, these two tests could serve as initial  
362 “gate keeper” assays, following which chemicals may be exempted from further testing  
363 (negatives) or (when positive) subjected to additional, confirmatory analyses with other existing  
364 Tier-1 assays. This would greatly enhance throughput of chemicals through initial testing, both  
365 in terms of resource utilization and timing.

## 366 **ACKNOWLEDGMENTS**

367 The authors would like to thank the following collaborators: Christopher J. Borgert, Mark  
368 Hanson, Alan Hosmer, Werner Kloas, and Glen Van Der Kraak.

## 369 **DISCLAIMER**

370 The views and statements expressed in this paper are those of the authors alone. The views  
371 or statements expressed in this publication do not necessarily represent the views of the  
372 organisations to which the authors are affiliated, and those organisations cannot accept any  
373 responsibility for such views or statements.

374 The manuscript has been subjected to review by the National Health and Environmental  
375 Effects Research Laboratory and the Office of Chemical Safety and Pollution Prevention and  
376 approved for publication. Approval does not signify that the contents reflect the views of the  
377 USEPA and mention of trade names or commercial products does not constitute endorsement or  
378 recommendation for use by USEPA.

## 379 **REFERENCES**

- 380 Adams NR. 1995. Detection of the effects of phytoestrogens on sheep and cattle. *J Anim Sci*  
381 73:1509-1515.
- 382 Ankley GT, Bennett RS, Erickson RJ, Hoff DJ, Hornung MW, Johnson RD, Mount DR, Nichols  
383 JW, Russom CL, Schmieder PK, Serrano JA, Tietge JE, Villeneuve DL. 2010. Adverse  
384 outcome pathways: A conceptual framework to support ecotoxicology research and risk  
385 assessment. *Environ Toxicol Chem* 29:730-741.
- 386 Ankley, GT, Gray LE. 2013. Cross-species conservation of endocrine pathways: A critical  
387 analysis of Tier 1 fish and rat screening assays with 12 model chemicals. *Environ.*  
388 *Toxicol. Chem.* 32, 1084-1087.
- 389 Ankley GT, Jensen KM. 2014. A novel framework for interpretation of data from the fish short-  
390 term reproduction assay (FSTRA) for the detection of endocrine-disrupting chemicals.  
391 *Environ. Toxicol. Chem.* DOI: 10.1002/etc.2708.
- 392 Borgert CJ, Mihaich EM, Ortego LS, Bentley KS, Holmes CM, Levine SL, Becker RA. 2011a.  
393 Hypothesis-driven weight of evidence framework for evaluating data within the U. S.  
394 EPA's Endocrine Disruptor Screening Program. *Reg Toxicol Pharmacol* 61:185-191.
- 395 Borgert CJ, Mihaich EM, Quill TF, Marty MS, Levine SL, Becker RA. 2011b. Evaluation of  
396 EPA's Tier 1 Endocrine Screening Battery and recommendations for improving the  
397 interpretation of screening results. *Reg Toxicol Pharmacol* 59:387-411.
- 398 Borgert, CJ, Stuchal LD, Mihaich EM, Becker RA et al., 2014. Relevance weighting of Tier 1  
399 endocrine screening endpoints by rank order. *Birth Defects Res (Part B)*, 101, 90–113.
- 400 Kim, HS, Kang TS, Kang IH, Kim TS, Moon HJ, Kim IY, et al. 2005. Validation study of  
401 OECD rodent uterotrophic assay for the assessment of estrogenic activity in sprague-

402 dawley immature female rats. *Journal of Toxicology and Environmental Health. Part A*  
403 68(23-24): 2249-62.

404 Klimisch H-J, Andreae M, Tillmann U. 1997. A systematic approach for evaluating the quality  
405 of experimental toxicological and ecotoxicological data. *Reg Toxicol Pharmacol* 25:1-5.

406 Schlenk D., Jenkins F. 2013. Endocrine Disruptor Screening Program (EDSP) Tier 1 Screening  
407 Assays and Battery Performance. US EPA FIFRA SAP Minutes No. 2013-03. May 21-  
408 23, 2013. Washington, DC.

409 [USEPA] U.S. Environmental Protection Agency. 2009. Endocrine disruptor screening program  
410 test guidelines - OCSPP 890.1250: Estrogen Receptor Binding Assay Using Rat Uterine  
411 Cytosol (ER-RUC). EPA 740-C-09-005.

412 [USEPA] U.S. Environmental Protection Agency. 2009. Endocrine Disruptor Screening Program  
413 Test Guidelines - OCSPP 890.1300: Estrogen Receptor Transcriptional Activation  
414 (Human Cell Line (HeLa-9903). EPA 740-C-09-006.

415 [USEPA] U.S. Environmental Protection Agency. 2009. Endocrine Disruptor Screening Program  
416 Test Guidelines - OCSPP 890.1150: Androgen Receptor Binding (Rat Prostate Cytosol).  
417 EPA 640-C-09-003.

418 [USEPA] U.S. Environmental Protection Agency. 2009. Endocrine Disruptor Screening Program  
419 Test Guidelines - OPPTS 890.1550: Steroidogenesis (Human Cell line - H295R). EPA  
420 640-C-09-003.

421 [USEPA] U.S. Environmental Protection Agency. 2009. Endocrine Disruptor Screening Program  
422 Test Guidelines - OPPTS 890.1200: Aromatase (Human Recombinant). EPA 740-C-09-  
423 004.

424 [USEPA] U.S. Environmental Protection Agency. 2009. Endocrine disruptor screening program  
425 test guidelines—OCSPP 890.1600: Uterotrophic assay. EPA 740/C-09-0010.  
426 Washington, DC.

427 [USEPA] U.S. Environmental Protection Agency. 2009. Endocrine disruptor screening program  
428 test guidelines—OCSPP 890.1400: Hershberger bioassay. EPA 740/C-09-008.  
429 Washington, DC.

430 [USEPA] U.S. Environmental Protection Agency. 2009. Endocrine disruptor screening program  
431 test guidelines—OCSPP 890.1500: Pubertal development and thyroid function in intact  
432 juvenile/peripubertal male rats. EPA 740/C-09/012. Washington, DC.

433 [USEPA] U.S. Environmental Protection Agency. 2009. Endocrine disruptor screening program  
434 test guidelines—OCSPP 890.1450: Pubertal development and thyroid function in intact  
435 juvenile/peripubertal female rats. EPA 740/C-09/009. Washington, DC.

436 [USEPA] U.S. Environmental Protection Agency. 2009. Endocrine disruptor screening program  
437 test guidelines—OCSPP 890.1350: Fish short-term reproduction assay. EPA 740/C-  
438 09/007. Washington, DC.

439 [USEPA] U.S. Environmental Protection Agency. 2009. Endocrine disruptor screening program  
440 test guidelines—OCSPP 890.1100: Amphibian Metamorphosis assay. EPA 740/C-  
441 09/002. Washington, DC.

442 [USEPA] United States Environmental Protection Agency. 2011. Endocrine Disruptor Screening  
443 Program, Weight-of-Evidence: Evaluating Results of EDSP Tier 1 Screening to Identify  
444 the Need for Tier 2 Testing. Office of Chemical Safety and Pollution Prevention.  
445 September 14, 2011.

446 Van Der Kraak GJ, Hosmer AJ, Hanson ML, Kloas W, Solomon KR. 2014. Effects of atrazine  
447 in fish, amphibians, and reptiles: An analysis based on quantitative weight of evidence.  
448 *Crit Rev Toxicol* 44(S5):1-66.

449 Walker E, Hernandez AV, Kattan MW. 2008. Meta-analysis: Its strengths and limitations.  
450 *Cleveland Clin J Med* 75:431-439.

451

452

453 Figure 1: Graphical illustration of an adverse outcome pathway. Outcomes at levels 4 and 5 are  
454 apical.

455 Figure 2: Illustration of the combination links in the AOP for reproduction for atrazine in fish,  
456 amphibians, and reptiles. The symbols indicate the mean score for relevance and strength and the  
457 vertical and horizontal bars  $2 \times \text{SE}$  of the mean score (from data in Van der Kraak et al. 2014)  
458

459 Table 1. Ability of the Tests in the Tier 1 Battery to Detect Endocrine Activity

<b>Estrogen, Androgen, Thyroid, and Steroidogenesis Pathways</b>	<b>Derivation of Detection Ability</b>
<b>Estrogenic Activity</b>	ER Binding and ERTA Uterotrophic Female Pubertal Fish Short-Term Reproduction Assay
<b>Anti-estrogenic Activity</b>	ER Binding Female Pubertal Fish Short-Term Reproduction Assay
<b>Androgenic Activity</b>	AR Binding Hershberger Male Pubertal Fish Short-Term Reproduction Assay
<b>Anti-androgenic Activity</b>	AR Binding Hershberger Male Pubertal Fish Short-Term Reproduction Assay
<b>Modulation of Steroidogenesis</b>	Steroidogenesis and Aromatase Assays Male and Female Pubertals Fish Short-Term Reproduction Assay
<b>Modulation of Aromatase</b>	Steroidogenesis and Aromatase Assays Female Pubertals Fish Short-Term Reproduction Assay
<b>Altered Hypothalamic-Pituitary Function</b>	Male and Female Pubertals Fish Short-Term Reproduction Assay Amphibian Metamorphosis Assay
<b>Anti-thyroid Activity</b>	Male and Female Pubertals Amphibian Metamorphosis Assay
<b>Thyromimetic Activity</b>	Amphibian Metamorphosis Assay

460

461

462 Table 2: Summary of Hypothesis-Based WoE Evaluations for Genistein for the Estrogen Agonist  
463 Hypothesis

	Rank 1	Rank 2	Rank 3
<b>Genistein</b>	Vitellogenin in male fish inconsistent (possibly due to route of exposure) [a,b]  Uterotrophic assays positive [c]	ERTA activation [d]; observed fish histopath [b], some changes in rat testes [e], some female pubertal changes [e].	ER binding positive; corroborative observations in pubertal endpoints [e]; steroid hormone changes in fish [b].

464 [a] Zhang, L., Khan, I. A., & Foran, C. M. (2002). Characterization of the estrogenic response to genistein in Japanese medaka  
465 (*Oryzias latipes*). *Comparative Biochemistry and Physiology. Toxicology & Pharmacology* : CBP, 132(2), 203-11.

466 [b] Bennetau-Pelissero, C., Breton B, B., Bennetau, B., Corraze, G., Le Menn, F., Davail-Cuisset, B., et al. (2001). Effect of  
467 genistein-enriched diets on the endocrine process of gametogenesis and on reproduction efficiency of the rainbow trout  
468 *Oncorhynchus mykiss*. *General and Comparative Endocrinology*, 121(2), 173-87.

469 [c] Kim, H. S., Kang, T. S., Kang, I. H., Kim, T. S., Moon, H. J., Kim, I. Y., et al. (2005). Validation study of OECD rodent  
470 uterotrophic assay for the assessment of estrogenic activity in Sprague-Dawley immature female rats. *Journal of Toxicology and  
471 Environmental Health. Part A*, 68(23-24), 2249-62.

472 [d] Ranhotra, H. S. & Teng, C. T. (2005). Assessing the estrogenicity of environmental chemicals with a stably transfected  
473 lactoferrin gene promoter reporter in HELA cells. *Environmental Toxicology and Pharmacology*, 20(1), 42-7.

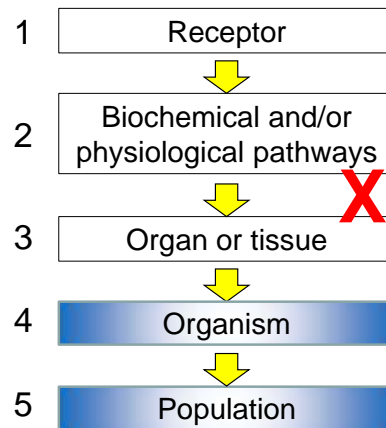
474 [e] Delclos, K. B., Bucci, T. J., Lomax, L. G., Latendresse, J. R., Warbritton, A., Weis, C. C., et al. (2001). Effects of dietary  
475 genistein exposure during development on male and female CD (Sprague-Dawley) rats. *Reproductive Toxicology* (Elmsford,  
476 N.Y.), 15(6), 647-63.

477



478 Figure 1: Graphical illustration of an adverse outcome pathway. Outcomes at levels 4 and 5 are  
479 apical.

480



481

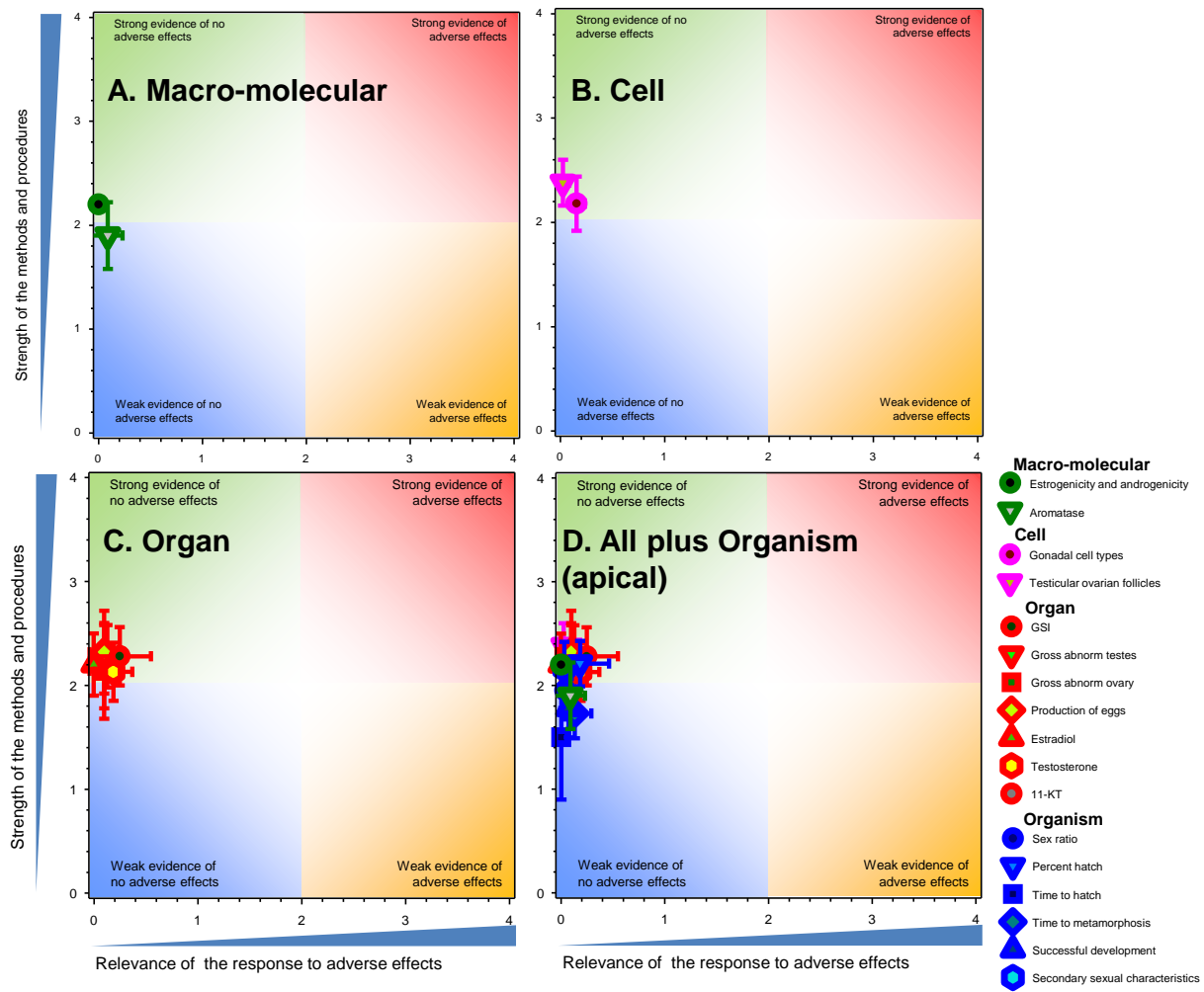
482

483

484

485

486 Figure 2: Illustration of the combination links in the AOP for reproduction for atrazine in fish,  
 487 amphibians, and reptiles. The symbols indicate the mean score for relevance and strength and the  
 488 vertical and horizontal bars 2xSE of the mean score (from data in Van der Kraak et al. 2014)



489