



ENVIRONMENTAL HEALTH PERSPECTIVES

<http://www.ehponline.org>

Standardizing Benchmark Dose Calculations to Improve Science-Based Decisions in Human Health Assessments

Jessica A. Wignall, Andrew J. Shapiro, Fred A. Wright,
Tracey J. Woodruff, Weihsueh A. Chiu, Kathryn Z. Guyton,
and Ivan Rusyn

<http://dx.doi.org/10.1289/ehp.1307539>

Received: 23 August 2013

Accepted: 24 February 2014

Advance Publication: 25 February 2014

Standardizing Benchmark Dose Calculations to Improve Science-Based Decisions in Human Health Assessments

Jessica A. Wignall,¹ Andrew J. Shapiro,¹ Fred A. Wright,² Tracey J. Woodruff,³ Weihsueh A. Chiu,⁴ Kathryn Z. Guyton,⁴ and Ivan Rusyn¹

¹Department of Environmental Sciences and Engineering, and ²Department of Biostatistics, University of North Carolina, Chapel Hill, North Carolina, USA; ³Department of Obstetrics, Gynecology, and Reproductive Sciences, School of Medicine, University of California, Oakland, California, USA; ⁴National Center for Environmental Assessment, United States Environmental Protection Agency, Washington, DC, USA

Address correspondence to Ivan Rusyn, 0031 MHRC, University of North Carolina, Chapel Hill, NC 27599-7431 USA. Telephone: (919) 843-2596. E-mail: iir@unc.edu

Running title: Standardized benchmark dose calculation.

Acknowledgments: The authors wish to thank SRC, Inc. for collating the dose-response information. This work was funded, in part, by grants from NIH (P42-ES005948) and US EPA (STAR-RD83516601). Ivan Rusyn and Tracey Woodruff were supported, in part, by the ORISE fellowship funded by the US EPA.

Disclaimer: The views in this article are those of the authors, and do not necessarily reflect the views or policies of the U.S. EPA.

Competing Financial Interests: The authors declare they have no actual or potential competing financial interests.

Abstract

Background: Benchmark dose (BMD) modeling computes the dose associated with a pre-specified response level. While offering advantages over traditional points of departure (POD), such as no-observed-adverse-effect-levels (NOAELs), BMD methods have lacked consistency and transparency in application, interpretation and reporting in human health assessments of chemicals.

Objectives: We aimed to apply a standardized process for conducting BMD modeling to reduce inconsistencies in model fitting and selection.

Methods: We evaluated 880 dose-response data sets for 352 environmental chemicals with existing human health assessments. We calculated benchmark doses and their lower limit [10% extra risk, or change in the mean equal to 1 standard deviation (SD), BMD/L_{10/1SD}] for each chemical in a standardized way with pre-specified criteria for model fit acceptance. We identified study design features associated with acceptable model fits.

Results: We derived values for 255 (72%) of chemicals. Batch-calculated BMD/L_{10/1SD} values were significantly and highly correlated (R^2 of 0.95 and 0.83, respectively, $n=42$) with points of departure previously used in human health assessments, with values similar to reported NOAELs. Specifically, the median ratio of BMD_{S10/1SD}:NOAELs was 1.96, and the median ratio of BMD_{L10/1SD}:NOAELs was 0.89. We also observed a significant trend of increasing model viability with increasing number of dose groups.

Conclusions: BMD/L_{10/1SD} values can be calculated in a standardized way for use in health assessments on a large number of chemicals/critical effects. This facilitates exploration of health effects across multiple studies of a given chemical, or when chemicals need to be compared, providing greater transparency and efficiency than current approaches.

Introduction

Public health agencies (e.g., the U.S. Environmental Protection Agency [EPA] and California EPA) conduct health assessments of environmental chemicals to determine the likelihood of human health hazard and to establish levels of exposure considered as health protective. To derive quantitative toxicity values (i.e., cancer slope factors or reference doses/concentrations) for comparison to environmental exposure levels, the relationship between a dose/concentration of a chemical and a health outcome is characterized (U.S. EPA 2012b). Data from occupational cohorts or from studies in experimental animals are typically used for this purpose (National Research Council 1983). The first step in developing toxicity values is identifying, for each data set, a POD dose from which extrapolations to environmentally relevant doses are made.

PODs traditionally used in non-cancer health effect assessments are NOAELs or lowest-observed-adverse-effect-levels (LOAELs) (U.S. EPA 2012b). NOAELs and LOAELs are limited to the dose groups tested in a particular study and are not informed by the shape of the dose-response relationship (Barnes and Dourson 1988; Travis et al. 2005). BMD modeling, a process of fitting a model to dose-response data, estimates a POD that is associated with a predefined level of biological response (i.e., the benchmark response [BMR]) (Crump 1984). BMD modeling addresses some limitations of NOAELs and LOAELs in that BMDs account for the shape of the dose-response curve, are more independent of study design elements such as dose choice or spacing, and can be more easily compared across multiple chemicals. In addition, estimating the BMD lower limit (BMDL) informs uncertainty in risk estimates. However, not all dose-response data sets are amenable to BMD modeling, for example, when group sizes are very small, but otherwise reflect the species of choice (e.g., as is often the case with dog studies).

BMD modeling is traditionally conducted on a chemical-by-chemical basis, with variability introduced during selection of critical endpoints, BMR values and models used to compute BMDs, as well as in evaluating model fit (Travis et al. 2005; U.S. EPA 2012b). For example, the biological significance of a given magnitude of change can differ among endpoints, especially when they range in severity. Thus, while the choice of BMR may vary from chemical to chemical and study to study, we investigated ways to standardize BMD methodology to increase consistency in POD derivation, reduce complexity, and improve efficiency.

A large database of developmental toxicity studies was used previously to derive BMD estimates (Allen et al. 1994a, b) to demonstrate that a standardized approach to dose-response modeling is advantageous. Using a limited set of data and models it was shown that BMDs based on a 5% extra risk response were within an order of magnitude of statistically derived NOAELs. In this study, we expand upon this previous work by applying a standardized process for conducting BMD modeling to 880 dose-response data sets for 352 environmental chemicals extracted from publicly available human health assessments. Using standard approaches, as recommended by U.S. EPA (2012a), we evaluate multiple endpoints and identify features of animal study methods that may influence their utility for BMD modeling.

Methods

Data sets

EPA Integrated Risk Information System (IRIS) (U.S. EPA 2013a), EPA Office of Pesticide Programs (U.S. EPA 2013b), EPA Superfund Regional Screening Levels (RSL) (U.S. EPA 2013d), and California EPA (OEHA 2013) were surveyed for publicly available information on chemicals with human health assessments. Superfund RSL also included toxicity values from

EPA Provisional Peer Reviewed Toxicity Value (U.S. EPA 2013c), Centers for Disease Control and Prevention's Agency for Toxic Substances and Disease Registry (ATSDR 2013), and EPA Health Effects Assessment Summary Tables (U.S. EPA 2011a). We collected both non-cancer and cancer toxicity values [reference doses (RfDs), reference concentrations (RfCs), oral slope factors, inhalation unit risks, and cancer potency values], and PODs that were used to derive the toxicity values, where applicable [NO(A)ELs, LO(A)ELs, and BMD/Cs used to derive RfDs and RfCs].

For each toxicity value, we extracted the dose-response data from the critical study used in the human health assessment. For each chemical, we obtained the name and a unique chemical identifier in the form of the Chemical Abstracts Service Registry Number (CASRN). The chemicals and their associated toxicity values, PODs, dose-response data and calculated BMD/Ls are available for download from <http://comptox.unc.edu/bmddata.php> (UNC 2013).

Chemical structure curation

Chemicals lacking CASRN were removed (e.g., mixtures such as “coke emissions”). CASRN were used to retrieve chemical structure information in the form of simplified molecular-input line-entry system codes (Weininger et al. 1989) which were converted to structure-data files using KNIME (Berthold et al. 2007). A rigorous chemical structure curation protocol (Fourches et al. 2010) was applied to ensure that the chemical structures were standardized and that mixtures and chemicals for which descriptors cannot be calculated (i.e., inorganics, organometallics) were removed.

BMD/L calculation

BMDs and BMDLs were calculated in a consistent fashion using BMDS Wizard (Beta Version 1.6.1) (ICF International 2012) and BMD Software (BMDS, Version 2.3.1). Specifically, we applied automated rules with no manual interpretation of results with respect to the following: (i) selection of the benchmark response (BMR) value; (ii) choice of the model(s); (iii) model fitting criteria; (iv) computation of the BMDL; and (v) reporting of BMD and BMDL values. All automated rules were consistent with BMD modeling guidelines (U.S. EPA 2012a). The results are hereto referred to as “batch-calculated” BMDs and BMDLs.

The BMDS Wizard program was used to automatically run BMDS. This program also recommended BMD/Ls for the collected dose-response data, based on the best-fitting model selected according to decision logic determined prior to modeling. The model decision logic and the criteria used to determine each model’s viability, based on adequacy of the fit of the model to the data are specified in Supplemental Material, Table S1. That is, after models are fit to the dose-response data, the tests listed in Supplemental Material, Table S1 were used to assign model fits of the dose-response data to Unusable, Questionable, or Viable categories by BMDS Wizard. As described in Figure 1, only Viable model outputs are used in the remainder of this analysis. We termed such Viable models “successful”.

Data sets were grouped according to dose-response type (continuous, dichotomous, or dichotomous-cancer), which guided the choice of BMRs and the types of models used to calculate BMDs. All models specified in the BMD modeling guidelines (U.S. EPA 2012a) were run for the appropriate data type (Table 1). Several additional model types that take into account more advanced biology, such as nested dichotomous, background-dose, background-response,

repeated response, concentration/time, and multi-tumor models were not within the scope of this project.

The BMR levels associated with the batch-calculated BMD/Ls (termed BMD/BMDL_{10/1SD} throughout) were standardized only according to the mathematical representation of the response data (continuous or dichotomous) and following the recommendations outlined in BMD Guidance (U.S. EPA 2012a). A 10% BMR was used for dichotomous data, and a “change in the mean equal to one control SD” BMR was used for continuous data. These two BMR levels are the standard reporting levels for each dose-response type, and do not necessarily represent equivalent values. However, Crump (1995) found that using a one control standard deviation change for continuous endpoint gives an excess risk of approximately 10% for the proportion of individuals below the 2nd percentile or above the 98th percentile of controls for normally distributed effects. Tailoring of BMR levels to the specific type or severity of the endpoint measured may depend on the decision-making context for which the BMD results will be used, and was therefore beyond the scope of this study.

The final model and associated BMD and BMDL for each dose-response set was selected according to the following criteria. The Viable model with the lowest Akaike’s Information Criterion (AIC) was always selected if the BMDLs were “sufficiently close”, i.e., there was no more than a 3-fold difference between lowest and highest BMDL for Viable models (Davis et al. 2011). Otherwise, the model with the lowest BMDL was selected. If no models were Viable, the highest dose(s) were removed and the models were re-run in cases where at least 3 (including control) doses remained. If two or more models had the same lowest AIC value and the BMD and BMDL values were different, the averages of the BMDs and BMDLs of those models were recorded. This final step is not done automatically by the BMDS Wizard. After completion of a

modeling run of a dose-response data set and BMR, for all successful models we recorded the BMD, BMDL, and any applicable model warnings or notes (based on passing or failing the tests listed in the decision logic reported in Supplemental Material, Table S1). If no model was successful, the dose-response data set was noted as having failed BMD modeling.

BMD/L selection

If a chemical had more than one dose-response data set, we selected the BMDs and BMDLs as follows: (i) the lowest BMD (without warnings, if available) and the BMDL associated with it, and (ii) the lowest BMDL (if different from the previous BMDL). These were selected regardless of endpoint/effect.

Data analysis

We examined the features of the overall resulting data set, including the range and distribution of the batch-processed BMD and BMDL values. BMDs and BMDLs calculated using the method described here were compared to BMDLs and other PODs, particularly NOAELs, for the same chemicals as reported in previous human health assessments, using several linear regression methods to calculate Pearson (R^2 values) and Spearman (ρ values) correlations. Tests for significance were calculated using a two-tailed unpaired t-test; p -values <0.05 were considered significant. Chi-squared test for trend in proportions was used to test for significance in trends; p -values <0.05 were considered significant. Statistical analysis and graphical outputs were produced by Microsoft Excel, R Statistical Package (Version 2.15), GraphPad Prism (La Jolla, CA) software, and the Health Assessment Workspace Collaborative (<https://hawcproject.org>) (Shapiro 2013).

Results

Curation of chemicals and data

We identified 1,260 chemicals with at least one EPA- or California EPA-derived toxicity value. Mixtures, chemicals missing structural information, and inorganic, organometallic, and duplicate structures were removed during curation (n=374). We collected dose-response data for 352 of the remaining 886 chemicals with toxicity values, yielding 880 dose-response data sets. We prioritized data collection according to public availability of the information (Supplemental Material, Figure S1).

BMD modeling

Of the 880 dose-response data sets available for analysis, we successfully (termed Viable in BMDS Wizard) modeled 603 according to the pre-specified statistical and other adequacy criteria given in Supplemental Material, Table S1 without any adjustments. Ninety-nine dose-response data sets contained fewer than 3 dose groups (including control) and thus could not be modeled. For 178 dose-response data sets, a first-pass attempt to model with all dose groups failed. When the highest dose group was omitted, we obtained successful models for an additional 66 dose-response data sets while 112 remained unmodelable. In total, 669 dose-response data sets were successfully modeled while 211 dose-response data sets were not (Supplemental Material, Figure S2). The modeled data sets covered 255 chemicals, whereas dose-response data sets for a remaining 97 chemicals did not pass model fit and completion tests. Overall, the modeling success rates were 86, 91 and 75% for cancer, dichotomous, and continuous data sets, respectively. The most frequently used model was exponential for continuous data sets and log-logistic for dichotomous data sets. See Supplemental Material,

Figure S3 for additional information on models used, including a characterization of models used by the number of dose groups.

We also evaluated the model-fit warnings associated with successful models (271 out of 669 or 40.5% successful data sets had at least one warning), and found that the majority (64%) of these concerned extrapolating more than 3 times below the lowest non-zero dose (median values were 6.4 for BMDL and 5.0 for BMD extrapolations). The next most common (13.2%), but not mutually exclusive, warning was high (>5) BMD/BMDL ratio (Supplemental Material, Figure S4).

Comparison to points-of-departure reported in human health assessments

We made statistical comparisons among previously reported and batch-calculated PODs for the PODs used as the basis for published RfDs (fewer data were available for comparison of PODs for other toxicity values and analyses were designed to be as consistent as possible). The lowest batch-calculated $BMD_{10/1SD}$ and $BMDL_{10/1SD}$ were compared with BMDLs from the same data set used for PODs in previous human health assessments. We found these untransformed values to be significantly and highly linearly correlated (R^2 of 0.95 and 0.83, respectively, $n=42$) (Figures 2A,B). More than 88% of values were within one order of magnitude of the BMDLs used in past assessments, and the mean values were not significantly different (Supplemental Material, Figure S5). We noted two outliers (both were included in the correlation analysis): dichloromethane and trichloroethylene (marked “a” and “b”, respectively, in Figures 2A,B).

These same batch-calculated $BMD_{10/1SD}$ and $BMDL_{10/1SD}$ were also compared with NOAELs from the same data set used as PODs in previous human health assessments, and after log-transformation to account for skewness were found to be significantly linearly correlated (R^2 of

0.66 for both, n=75) (Figures 2C,D; see Supplemental Material, Table S2 for the list). The comparison was further made with LOAELs used previously as PODs, or all previous PODs aggregated together, with significant linear correlation after log transformation (LOAELs: R^2 of 0.78 and 0.63, respectively, n=20; PODs: R^2 of 0.62 and 0.59, respectively, n=138) (Supplemental Material, Figure S6).

Comparison to NOAELs reported in human health assessments

We calculated the ratios of batch-calculated $BMD_{S_{10}/1SD}$ and $BMDL_{S_{10}/1SD}$ to oral NOAELs reported in previous health assessments (Figures 3A,B; n=75) (there was an insufficient number of inhalation NOAELs for statistical comparison), respectively. The median ratio of $BMD_{S_{10}/1SD}:NOAEL$ was 1.96, with a 5th to 95th percentile range of 0.24 to 56.9. The median ratio of $BMDL_{S_{10}/1SD}:NOAEL$ was 0.89, with a 5th to 95th percentile range of 0.06 to 23.7. In addition, we compared LOAELs from the studies used to identify the NOAELs used in the previous health assessments when available, and found a median ratio of 3.81 with a 5th to 95th percentile range of 1.87 to 10.7 (Figure 3C, n=68).

Batch-calculated BMD/Ls permit comparisons among adverse effects and chemicals

We selected nitroguanidine (CASRN 556-88-7) as an example chemical to illustrate how the standardized BMD approach can be used to calculate “batch-calculated candidate reference values” among multiple adverse health effects. Several dose-response data sets were available for nitroguanidine, including body weight changes, maternal toxicity, and non-neoplastic histopathological changes. In the original human health assessment, all of these endpoints were used to select a single NOAEL and derive an RfD. The collection of batch-calculated $BMDL_{S_{10}/1SD}$ was arrayed and compared to the NOAEL (Figure 4A) (U.S. EPA 1993).

Uncertainty factors for interspecies uncertainty ($UF_A=10$), intraspecies variability ($UF_H=10$), sub-chronic to chronic extrapolation ($UF_S=10$), and database incompleteness ($UF_D=3$) were applied in the original assessment to derive a reference dose of 0.1 mg/kg/day. “Batch-calculated candidate RfDs” based on batch-calculated BMDLs and the same uncertainty factors are presented in Figure 4A. The same type of analysis was conducted for di(2-ethylhexyl)adipate (CASRN 103-23-1) and pentachlorophenol (CASRN 87-86-5) (Supplemental Material, Figure S7).

We also used BMDs to illustrate comparisons across chemicals, as they reflect central estimates of the dose associated with a standardized level of benchmark response based only on the mathematical representation of the response (continuous or dichotomous). We ranked multiple chemicals according to their calculated $BMDs_{10/1SD}$ (i.e., relative potency) in Figure 4B.

Study design features as a factor in BMD modeling success

Because about a quarter of the dose-response data sets could not be successfully modeled using the BMD approach (i.e., Unusable or Questionable according to BMDS Wizard), we reviewed study design characteristics that may be associated with success or failure of modeling. Dose-response data sets that were not modeled successfully failed for a variety of reasons, including poorly modeled variance, goodness of fit p -test values <0.05 , or a lack of confidence in calculated values, such as by having a BMDL higher than highest dose or a BMD/BMDL ratio >20 (Supplemental Material, Figure S8).

We found that there is a significant ($p<0.05$) difference in the number of dose groups of successful dose-response data sets vs. unsuccessful dose-response data sets (Supplemental Material, Figure S9). Upon further examination, we observed a significant ($p<0.01$) trend of

increasing viability of models with increasing numbers of dose groups (Figure 5A). We found that the number of animals per dose group is statistically significantly associated with BMD modeling success ($p < 0.001$) (Figure 5B). Successful models had lower numbers of animals per dose group than unsuccessful models, across all dose-response data types (i.e., dichotomous, dichotomous cancer, continuous). There was no correlation between the number of dose groups and number of animals per dose group (data not shown). The spacing between the dose level of dose group 2 and dose group 3 was not associated with BMD modeling outcome (data not shown).

Discussion

We evaluated the efficacy and reliability of a standardized BMD approach, compared it to chemical-specific BMD modeling, and identified lessons learned for future application of BMD modeling in human health assessments. Our analysis indicates that a standardized approach can be successfully applied to a large number of chemicals and data sets. We limited our analysis to the dose-response data sets from which PODs were identified in past assessments, but which were not necessarily chosen with BMD modeling in mind. It is likely that this approach would be even more successful if applied to data sets specifically chosen for BMD modeling (e.g., those with sufficient dose groups and dose-response trends) (Davis et al. 2011).

We compared batch-calculated BMD/Ls based on a standardized, guidance-driven choice of benchmark responses and models to BMD/Ls based on chemical-specific decisions made by different assessors and at different times. Batch-calculated BMD/Ls were significantly correlated with BMDLs derived one chemical at-a-time. Approximately 20% of the batch-calculated values used a different BMR from the BMR used in the original assessment (Supplemental Material,

Table S3). Two outliers were dichloromethane and trichloroethylene and the difference was largely due to use of PBPK model-based dosimetry in the original assessments. The PODs for these two chemicals already reflected a conversion from animal to human equivalent dose and an adjustment for human toxicokinetic variability (U.S. EPA 2011b, U.S. EPA 2011c). For trichloroethylene, an additional difference was the use of a 10% extra risk in the batch-calculated modeling as opposed to a 1% extra risk in the assessment (U.S. EPA 2011c).

Because our analysis uniquely included BMD, BMDL, and NOAEL values for 75 chemicals, we evaluated the relationship between batch-calculated BMDs and BMDLs and NOAELs selected during the course of a human health assessment. NOAELs are thought to approximate the dose that represents a 1 to 5% BMR (Allen et al. 1994a). However, we show that BMDs based on a 10% or 1SD BMR are similar to NOAELs (Figure 3B) (U.S. EPA 2012a). Similarly, Sand et al (2011) found that the median upper bound on extra risk at the NOAEL was approximately 10% using 786 NTP cancer data sets.

Our analysis also highlights the utility of BMD modeling and batch-processed candidate reference value calculations in evaluating the entirety of a database on a specific chemical. While we only used data from the critical study evaluated in the original human health assessment, we demonstrate that BMDLs can be calculated in a standardized way to facilitate comparison among multiple health effects and multiple studies at a fixed BMR, consistent with the advice from the National Academies (National Research Council 2009). This approach also aids identification of outlier evidence or studies if some calculations are orders of magnitude higher or lower than the balance of the data. Thus, this approach can increase objectivity in evaluating multiple studies, enhance transparency, and improve communication with assessors, peer-reviewers, and the general public.

We posit that a comparable approach can be applied in other contexts. For example, high-throughput *in vitro* testing is producing vast amounts of data, consisting of hundreds of dose-response data sets on thousands of chemicals. However, it is unrealistic to expect that individual evaluation of concentration-response relationships in each data set would be commensurate with timely and efficient analyses of these data. Calculation of BMD-like values from *in vitro* data has been suggested (National Research Council 2009), and our approach can be applied to increase efficiency and transparency in processing such large data sets. Sand et al. (2012) provide a comprehensive review of the considerations for selecting appropriate standardized BMRs when performing concentration-response analysis of *in vitro* data. Consistent selection and application of BMRs and a standardized decision logic yields values that enable comparisons across chemicals (Sirenko et al. 2013) and may inform further testing using a process that is relevant for and familiar to risk managers and decision makers.

Additionally, consistently derived BMDs that represent the same biological response can provide valuable quantitative information for other analyses. For example, they can be used to evaluate the potential for quantitative structure-activity relationship modeling. If a chemical structure is found to be predictive of a chemical's BMD, this would allow decision-makers to evaluate a chemical's potential hazard to human health even if animal or human data on that chemical are lacking.

Our analysis also informs advancement of a unified dose-response modeling framework that is applied consistently to cancer and non-cancer effects proposed by NRC (National Research Council 2009). The exact nature and implementation of this framework has yet to be determined. For dichotomous endpoints, current EPA BMDS guidance specifies a smaller and more constrained set of models for cancer than for non-cancer endpoints (U.S. EPA 2012a). This is a

potential area for harmonization as health assessments move towards unifying the cancer and non-cancer assessments that could be readily explored by the batch-processing approach explored herein.

Finally, results of our analysis also give insight into study design attributes that increase the potential for BMD modeling success. We observed that successful dose-response data sets tend to have higher numbers of dose groups with fewer animals in each dose group. This result is in accord with Slob et al. (2005) who found, using simulated data, that a higher number of dose groups will help to define the shape of the dose-response and may minimize the risk for unfavorable dose placement. This may be due to several factors. First, as the number of animals in each dose group increases, flexibility in slight deviations between the statistical model's shape and the true underlying dose-response function decreases. Second, for dichotomous models, there may be sources of variation beyond the binomial statistics assumed by BMDS. In either case, a statistically poor model fit is more likely with more animals per dose group, all other things being equal. This may arise because the test for lack of fit has more power and is more likely to reject the model fit when group sample size is high. Nonetheless, this finding does not imply that fewer animals per dose group is preferable overall. Modeling success needs to be balanced against having enough statistical power to detect a response (Melnick et al. 2008). Because the majority of warnings found in otherwise successful models are due to extrapolation more than 3 times below the lowest non-zero dose, it is likely that those data sets did not have adequate data to support the BMRs used, and such a warning would not have occurred if a higher BMR has been selected. Additionally, the models may not account for non-biological sources of variation (e.g., group effects) and are dependent on a biological or statistical dose-response trend (Sand et al. 2008). Consideration of these factors together with a more detailed evaluation of the

characteristics of dose-response data sets associated with BMD modeling success might illuminate additional useful trends that can inform future study design.

We acknowledge several limitations. Because we did not conduct chemical-by-chemical evaluation, the BMR was not adjusted based on data source or effect severity. A higher or lower BMR may be warranted based on the study type (e.g., epidemiological vs. experimental animal) or severity of the biological response (e.g., developmental malformations vs. organ-specific histopathological changes). However, it is likely that a fixed BMR would be appropriate still (i.e., to enable comparisons among chemicals with the same critical effect and observed severity) in contexts using a standardized BMD process. Additionally, BMD models might fit the data mathematically, but may not inform plausibility of the biological response (Davis et al. 2011). Statistical evaluation was limited to model-fit criteria and did not include other considerations, such as evaluation of the model fit in the low dose region. Also, cutoffs were fixed in an automated manner according to the decision logic, resulting in less flexibility in assessing model viability than if each cutoff were independently adjusted. These issues can be addressed by a chemical-by-chemical or model-by-model analysis, if necessary.

Furthermore, when using BMD modeling for the purpose of deriving a chemical-specific POD, EPA guidance recommends an evaluation of the pertinent literature to first identify the most appropriate study(ies) for analysis, based on hazard identification, the type of data, and study design (U.S. EPA 2012a). However, our analysis was based on studies that were not necessarily selected for their amenability to BMD modeling. Thus, for a given chemical, it was possible that the dose-response data were unavailable due to inadequate reporting (e.g., original data not provided or only represented graphically in primary literature, group means reported without standard deviations, no control group reported). This highlights the importance of presenting the

raw data used to identify the POD in assessment summaries (such as the online IRIS Summaries).

Conclusions

We demonstrate that a standardized BMD modeling approach can be used to derive $BMD/L_{S_{10}/1SD}$ that are significantly and highly correlated with BMDLs derived one chemical at a time. The median ratio of $BMD_{S_{10}/1SD}$ to NOAEL was less than 2, while $BMDL_{S_{10}/1SD}$ values were generally even lower than NOAELs. Deriving BMD/Ls in a consistent way across chemicals and endpoints gives values that represent the same response level and which are therefore useful in various decision-making contexts, such as identifying a candidate reference value, or determining relative potency of chemicals. Such a standardized approach can also be applied to data sets when speed and efficiency are priorities (e.g., *in vitro* assays). Ultimately, we demonstrated that a standardized approach, which makes BMD modeling transparent and easy to reproduce, is feasible and thus may be considered for wider use in certain decision contexts and types of assessments. In specific cases, expert judgment will still be needed in evaluations of alternative BMRs based on the study type or severity of biological response. Such judgment will assure that the standardized BMD modeling yields an accurate reflection of the underlying biology.

References

- Allen BC, Kavlock RJ, Kimmel CA, Faustman EM. 1994a. Dose-response assessment for developmental toxicity. II. Comparison of generic benchmark dose estimates with no observed adverse effect levels. *Fundam Appl Toxicol* 23:487-495.
- Allen BC, Kavlock RJ, Kimmel CA, Faustman EM. 1994b. Dose-response assessment for developmental toxicity. III. Statistical models. *Fundam Appl Toxicol* 23:496-509.
- ATSDR (Agency for Toxic Substances and Disease Registry). 2013. Minimal Risk Levels (MRLs). Available: <http://www.atsdr.cdc.gov/mrls/index.asp> [accessed 1 November 2013].
- Barnes DG, Dourson M. 1988. Reference dose (RfD): Description and use in health risk assessments. *Regul Toxicol Pharmacol* 8:471-486.
- Berthold MR, Cebron N, Dill F, Gabriel TR, Kötter T, Meinel T, et al. 2007. KNIME: The Konstanz information miner. Springer.
- Crump KS. 1984. A new method for determining allowable daily intakes. *Fundam Appl Toxicol* 4:854-871.
- Crump, KS. 1995. Calculation of benchmark doses from continuous data. *Risk Anal* 15:79-89.
- Davis JA, Gift JS, Zhao QJ. 2011. Introduction to benchmark dose methods and U.S. EPA's benchmark dose software (BMDS) version 2.1.1. *Toxicol Appl Pharmacol* 254:181-191.
- Fourches D, Muratov E, Tropsha A. 2010. Trust, but verify: On the importance of chemical structure curation in cheminformatics and qsar modeling research. *J Chem Inf Model* 50:1189-1204.
- ICF International. 2012. BMDS Wizard: Installation and user's guide. Durham, NC.
- Melnick RL, Thayer KA, Bucher JR. 2008. Conflicting views on chemical carcinogenesis arising from the design and evaluation of rodent carcinogenicity studies. *Environ Health Perspect* 116:130-135.
- National Research Council. 1983. Risk assessment in the federal government: Managing the process. Washington, D.C.: National Academy Press.
- National Research Council. 2009. Science and decisions: Advancing risk assessment. Washington, D.C.: National Academy Press.

- OEHHA (Office of Environmental Health Hazard Assessment). 2009. California EPA OEHHA Cancer Potency Values as of July 21, 2009. Available: <http://www.oehha.ca.gov/risk/pdf/tcdb072109alpha.pdf> and http://www.oehha.ca.gov/air/hot_spots/2009/AppendixA.pdf [accessed 1 November 2013].
- Sand S, Portier CJ, Krewski D. 2011. A signal-to-noise crossover dose as the point of departure for health risk assessment. *Environ Health Perspect* 119:1766-1774.
- Sand S, Ringblom J, Hakansson H, Oberg M. 2012. The point of transition on the dose-effect curve as a reference point in the evaluation of in vitro toxicity data. *J Appl Toxicol* 32(10):843-849.
- Sand S, Victorin K, Filipsson AF. 2008. The current state of knowledge on the use of the benchmark dose concept in risk assessment. *J Appl Toxicol* 28:405-421.
- Shapiro A. 2013. Health assessment workspace collaborative (HAWC). Available: <https://hawcproject.org/> [accessed November 1, 2013].
- Slob W, Moerbeek M, Rauniomaa E, Piersma AH. 2005. A statistical evaluation of toxicity study designs for the estimation of the benchmark dose in continuous endpoints. *Toxicol Sci* 84(1):167-185.
- Sirenko O, Cromwell EF, Crittenden C, Wignall JA, Wright FA, Rusyn I. 2013. Assessment of beating parameters in human induced pluripotent stem cells enables quantitative in vitro screening for cardiotoxicity. *Toxicol Appl Pharmacol* 273:500-507.
- Travis KZ, Pate I, Welsh ZK. 2005. The role of the benchmark dose in a regulatory context. *Regul Toxicol Pharmacol* 43:280-291.
- U.S. EPA. 1993. Integrated risk information system (IRIS): Online substance file for nitroguanidine (casrn 556-88-7). Available: <http://www.epa.gov/iris/subst/0402.htm> [accessed November 1, 2013].
- U.S. EPA. 2011a. Health Effects Assessment Summary Tables (HEAST) for Superfund. Available: <http://epa-heat.ornl.gov/> [accessed 1 November 2013].
- U.S. EPA. 2011b. Toxicological Review of Dichloromethane (Methylene Chloride) (CAS No. 75-09-2): In Support of Summary Information on the Integrated Risk Information System (IRIS). EPA/635/R-10/003F. Washington, DC

- U.S. EPA. 2011c. Toxicological Review of Trichloroethylene (CAS No. 79-01-6): In Support of Summary Information on the Integrated Risk Information System (IRIS). EPA/635/R-09/011F. Washington, DC.
- U.S. EPA. 2012a. Benchmark dose technical guidance. Washington, DC: Risk Assessment Forum, US EPA.
- U.S. EPA. 2012b. Step 2 - dose-response assessment. Available:
http://www.epa.gov/risk_assessment/dose-response.htm [accessed April 1, 2013].
- U.S. EPA. 2013a. US Environmental Protection Agency Integrated Risk Information System (IRIS) Homepage. Available: <http://www.epa.gov/iris/> [accessed 1 November 2013].
- U.S. EPA. 2013b. US Environmental Protection Agency Office of Pesticide Programs: Human Health Benchmarks for Pesticides. Available:
<http://iaspub.epa.gov/apex/pesticides/f?p=HHBP:home> [accessed 1 November 2013].
- U.S. EPA. 2013c. US Environmental Protection Agency Provisional Peer Reviewed Toxicity Values for Superfund (PPRTV). PPRTV Assessments Electronic Library. Available:
<http://hhpprtv.ornl.gov/>. [accessed 1 November 2013].
- U.S. EPA. 2013d. US Environmental Protection Agency Region 9: Regional Screening Levels (Formerly PRGs). Available: <http://www.epa.gov/region9/superfund/prg/> [accessed 1 November 2013].
- UNC. 2013. Standardized benchmark dose calculation database. Available:
<http://comptox.unc.edu/bmddata.php> [accessed November 1, 2013].
- Weininger D, Weininger A, Weininger JL. 1989. Smiles. 2. Algorithm for generation of unique smiles notation. *J Chem Inf Comput Sci* 29:97-101.

Table 1. Summary of BMRs and models used in BMDS, according to dose-response type.

Dose-response type	Dichotomous	Continuous	Dichotomous-cancer
Benchmark response	10% extra risk	Change in the mean equal to 1 control group SD ^a	10% extra risk
Models used to calculate BMDs and BMDLs ^b	Gamma, Dichotomous-Hill, Logistic, LogLogistic, Probit, LogProbit, Weibull, and Multistage ^c	Exponential 2, Exponential 3, Exponential 4, Exponential 5, Hill, Power, Polynomial ^c , and Linear (both constant and modeled variance models for each model above)	Cancer multistage 1 st -order through n-1 order where n is the number of dose groups
Distribution assumption	Binomial	Normal	Binomial

^aThis control group SD is the modeled SD. ^bModels selected based on defaults in BMDS and preferences of EPA IRIS program (U.S. EPA, 2012a). ^cOf order $n-1$ where n is the number of dose groups for each data set modeled.

Figure Legends

Figure 1. Schematic of BMDS Wizard workflow, adapted with permission from (ICF International 2012).

Figure 2. Correlations of batch-calculated BMDs and BMDLs with BMDLs (A,B) and NOAELs (C,D) reported in human health risk assessments. R^2 values represent squared Pearson correlations. ρ values represent Spearman correlations. Dotted line represents the regression line through the origin. Solid line represents the best-fit line. “a” denotes dichloromethane values; “b” denotes trichloroethylene values.

Figure 3. Histograms of log-transformed ratios of batch-calculated BMDs to NOAELs (A), BMDLs to NOAELs (B), and LOAELs to NOAELs (C). Y-axis: frequency counts; X-axis: magnitude of the ratio; red dotted lines: 5th and 95th percentiles of the distribution; red arrows: median values.

Figure 4. Array of batch-calculated BMDLs for the critical effects observed in studies of nitroguanidine as compared to IRIS NOAEL and RfD (A), and array of batch-calculated BMDs for selected chemicals compared to RfDs and PODs reported in human health assessments (B). Yellow circles: batch-calculated BMDs and BMDLs; orange circles: RfDs based on batch-calculated BMDLs; colored bars: uncertainty factors; blue squares: human health assessment PODs; gray squares: human health assessment RfDs. ^aReduced body weight gain. ^bRetarded ossification of pubis. ^cFewer than 3 sternebrae ossified. ^dFewer than 3 caudal vertebra ossified. ^eReduced weight gain in female rats. ^fReduced weight gain in female rats. ^gRetarded ossification of pubis. ^hFewer than 3 caudal vertebra ossified. ⁱFewer than 3 sternebrae ossified. ^jReduced body weight gain. ^kMaternal toxicity. ^lRenal lesions (glomerulosclerosis). ^mDecreased

delayed hypersensitivity response. ⁿRenal tubule regeneration. ^oIncreased splenic weight.

^pRenal cytomegaly. ^qNest-like infolds of the nasal respiratory epithelium. ^rChronic irritation.

^sLung adenoma or carcinoma (combined). ^tHemosiderin deposition in the liver. ^uIncreased

mortality. ^vLung and kidney histopathology. ^wReduced offspring body weight.

Figure 5. Relationship of Viable BMD models to (A) the number of dose groups, (B) number of animals in each dose group. ** indicates significant for trend ($p < 0.01$); *** indicates a significant difference between group means ($p < 0.001$).

Figure 1.

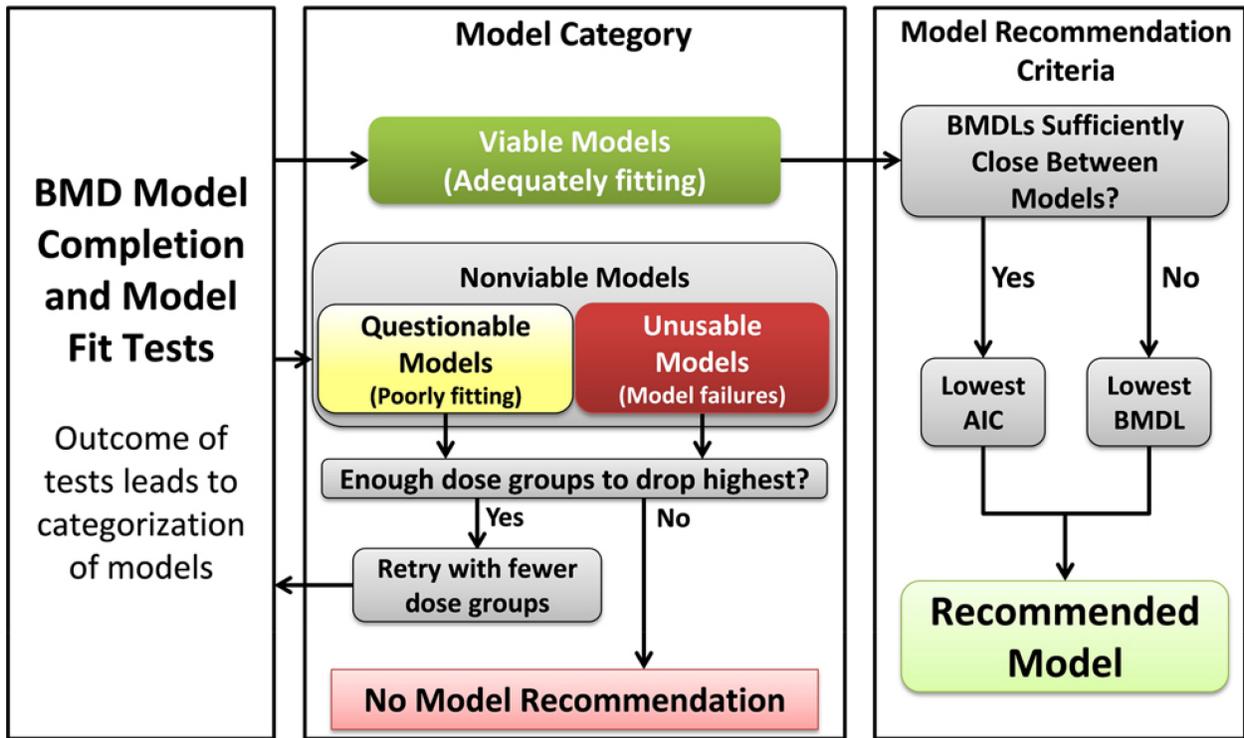


Figure 1

Figure 2.

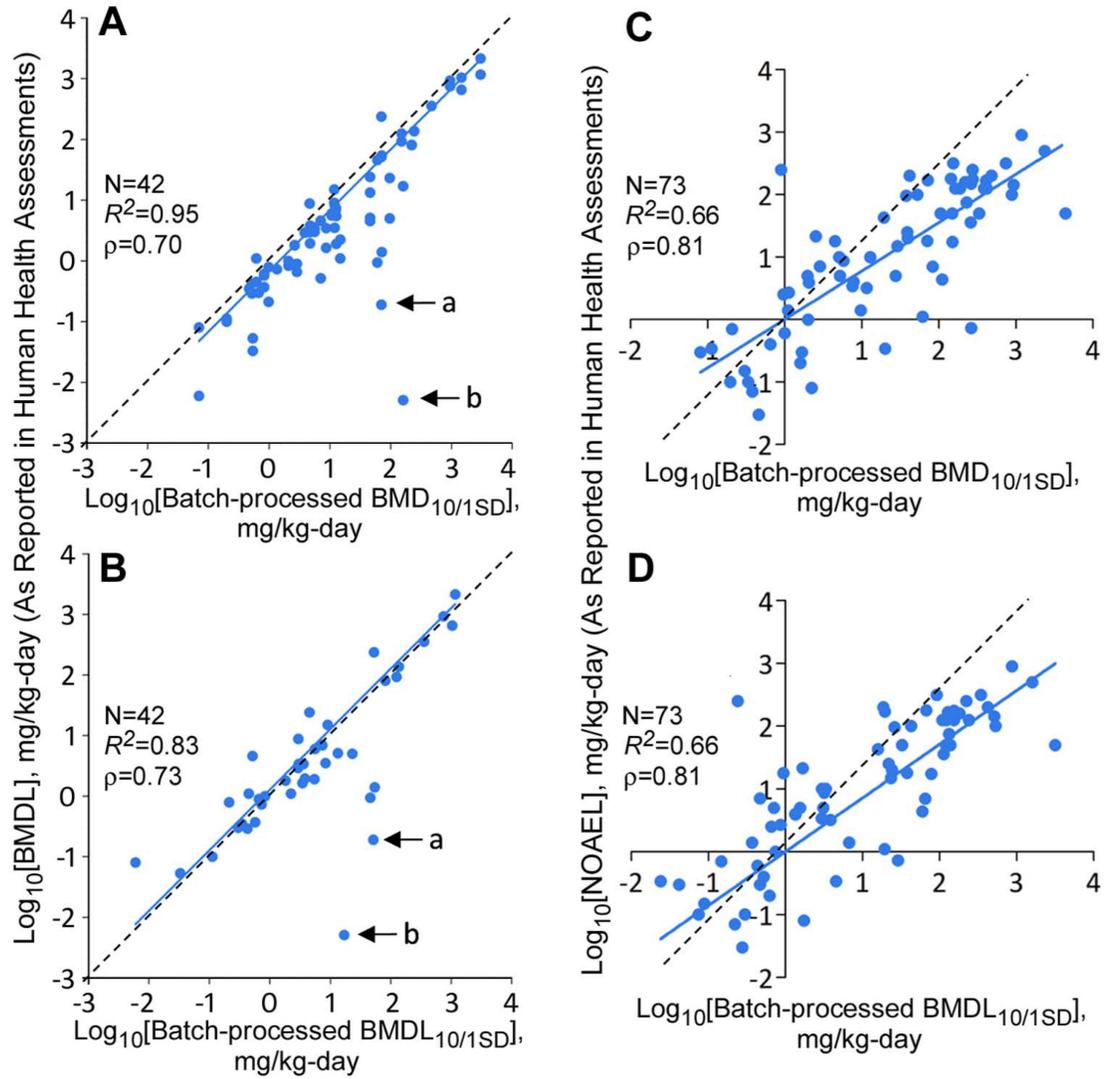


Figure 2

Figure 3.

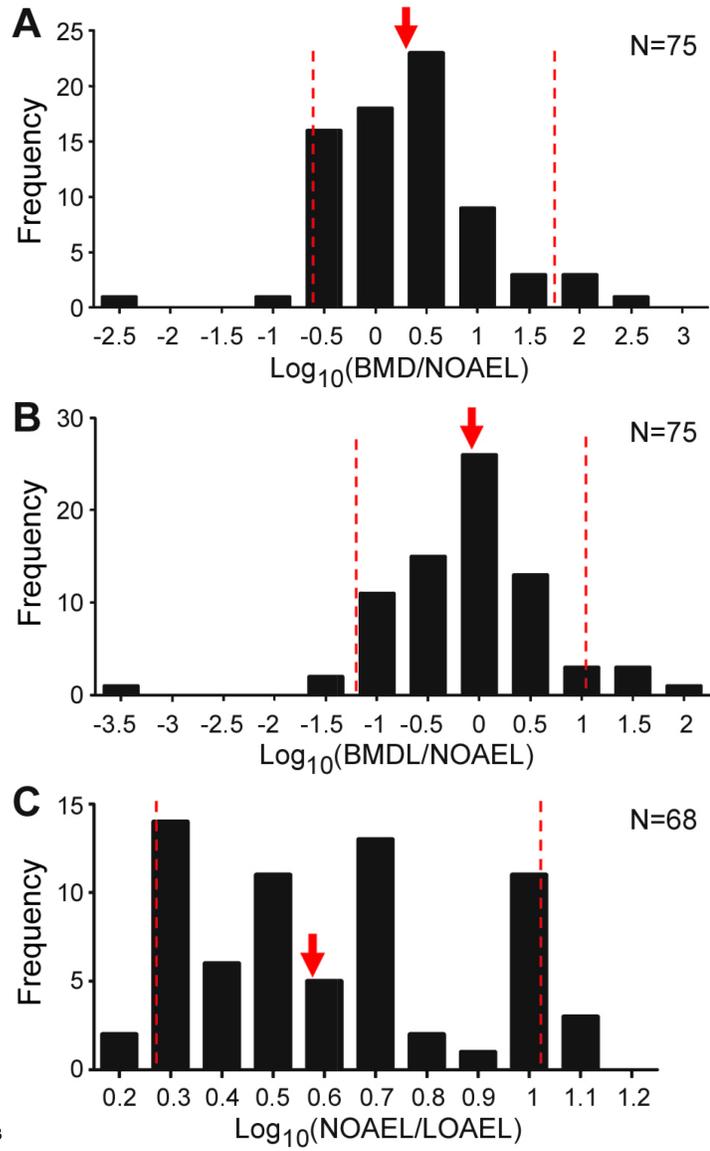


Figure 3

Figure 4.

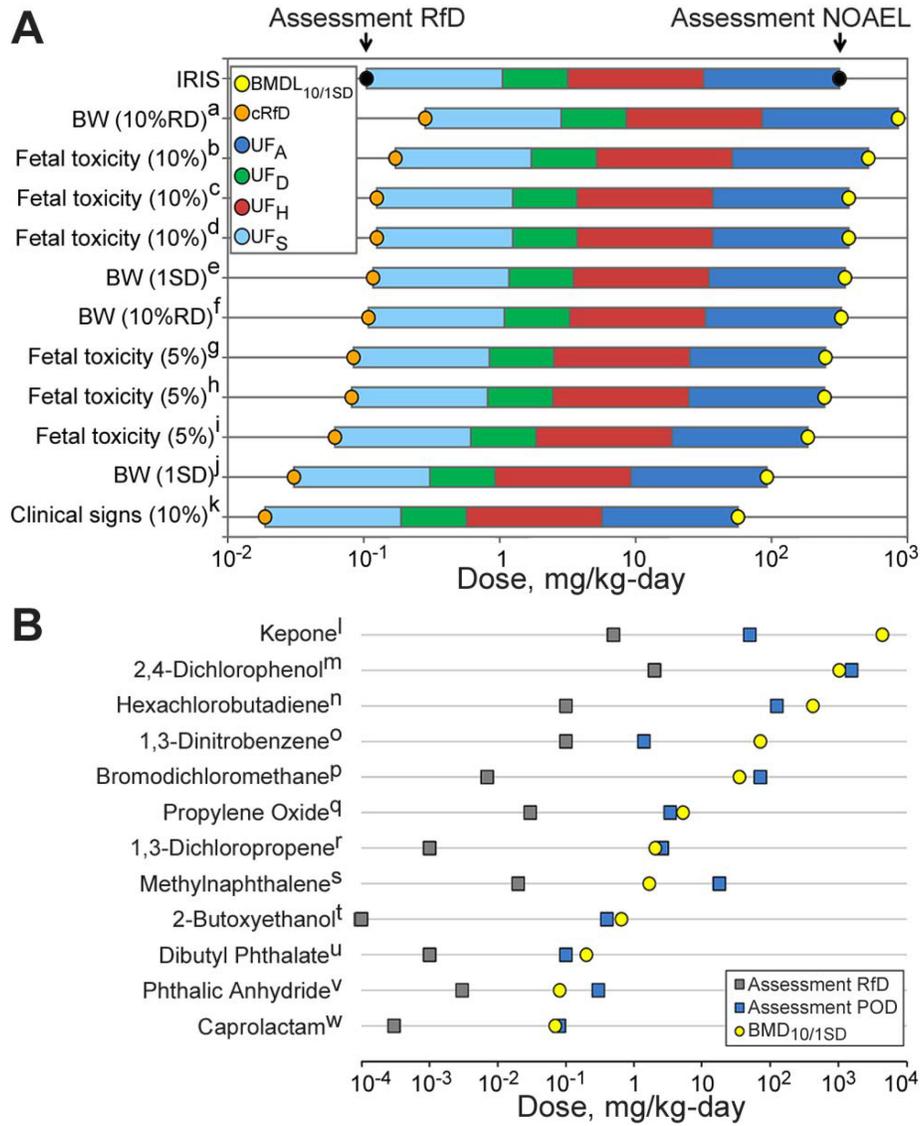


Figure 4

Figure 5.

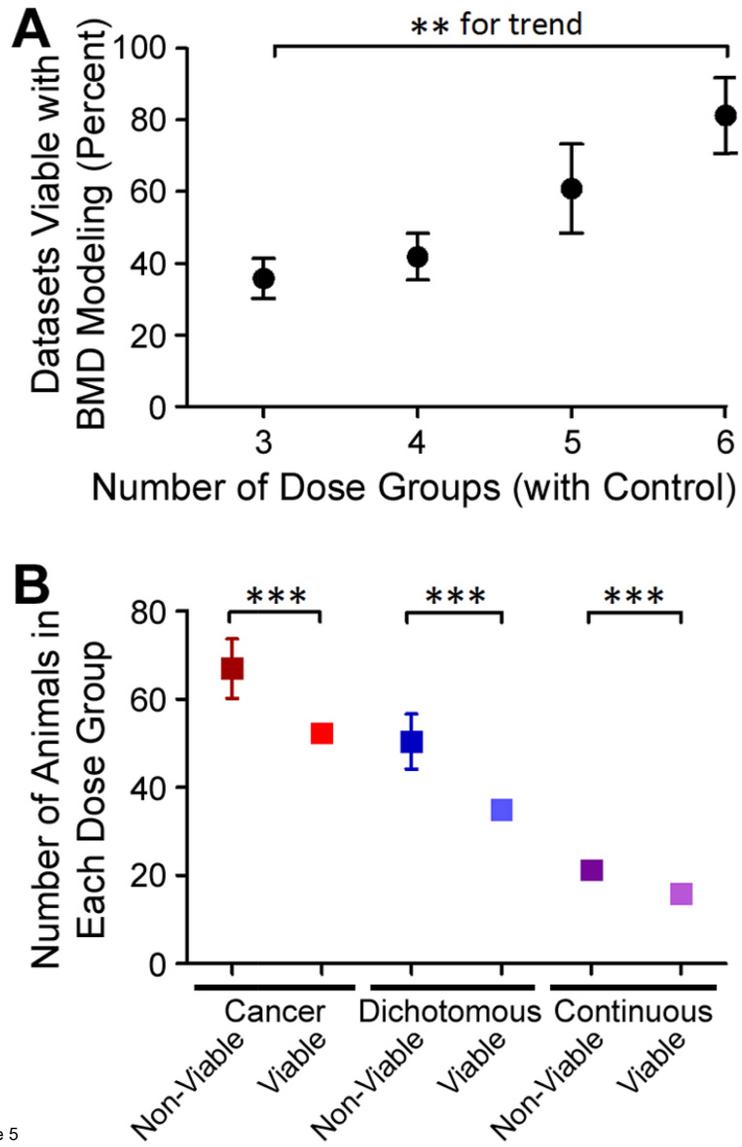


Figure 5