



# Quantitative variability in repeat dose toxicity studies: Implications for scientific confidence in NAMs

Katie Paul Friedman, PhD

December 17, 2019

State of the Science on Development and Use of New Approach Methods (NAMs) for Chemical Safety Testing

*The views expressed in this presentation are those of the authors and do not necessarily reflect the views or policies of the U.S. EPA*



# What is needed to understand the acceptability of NAMs for risk assessment?

- In US, Section 4(h) in the Lautenberg amendment to TSCA:
  - “...Administrator shall reduce and replace, to the extent practicable and scientifically justified...the use of vertebrate animals in the testing of chemical substances or mixtures...”
  - New approach methods (NAMs) need to provide “information of equivalent or better scientific quality and relevance...” than the traditional animal models
- “Directive to Prioritize Efforts to Reduce Animal Testing” memorandum signed by Administrator Andrew Wheeler on September 10, 2019
  - “1. Validation to ensure that NAMs are equivalent to or better than the animal tests replaced.”

**How do we define expectations of *in silico*, *in chemico*, and *in vitro* models for predicting repeat-dose toxicity?**

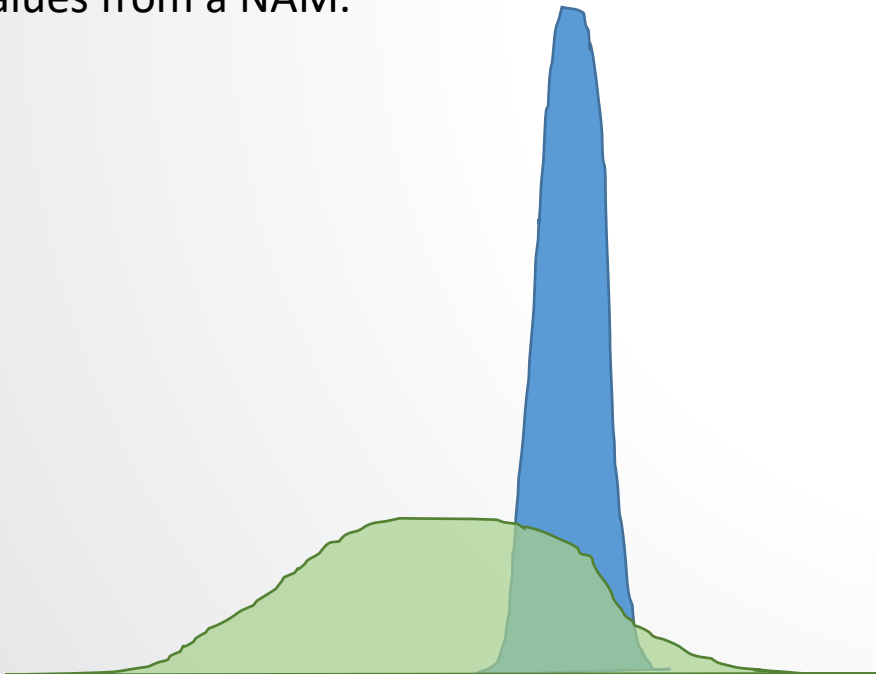
*In silico*, *in chemico*, and *in vitro* models cannot predict *in vivo* systemic effect values with greater accuracy than those animal models reproduce themselves.



# How do we express variability in traditional animal toxicity tests?

**Quantitative: variance is a measure of how far values are spread from the average.**

We need to know what the “spread” or variability of traditional effect levels (e.g., lowest effect levels, LELs, or lowest observable adverse effect levels, LOAELs) might be to know the range of acceptable or “good” values from a NAM.



**Qualitative: We need to know if a specific effect is always observed or not.**

		“Truth” (traditional toxicology)	
		Negative	Positive
Predicted (NAM)	Negative	True negative	False negative
	Positive	False positive	True positive

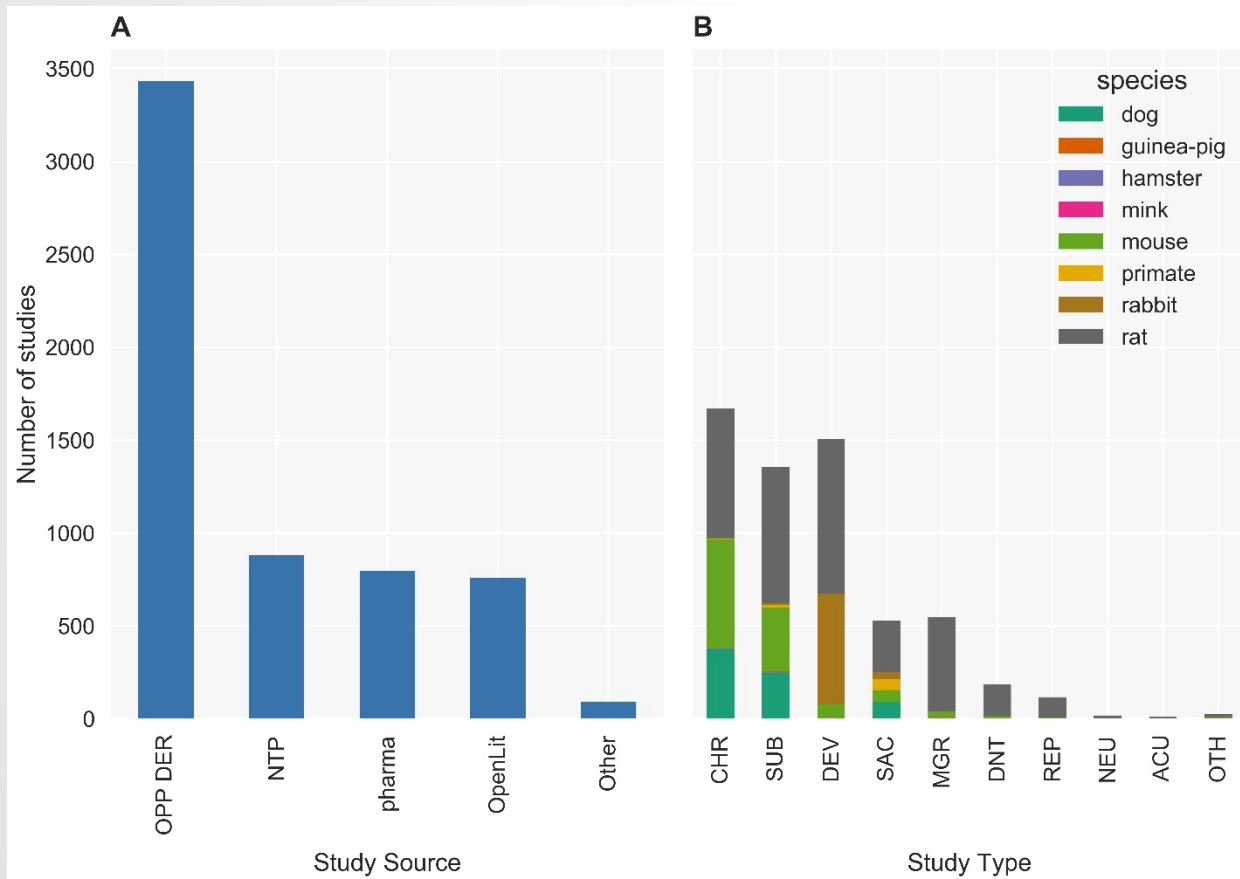


# Research questions for understanding this variability

3 main questions	What is the range of possible systemic effect values (mg/kg/day) in replicate studies?	What is the maximal accuracy of a model that attempts to predict a systemic effect values for an unknown chemical?	What is the probability that an effect in adult animals will be observed in replicate studies?
Statistical approach to the question	<ul style="list-style-type: none"><li>Residual root mean square error (RMSE) is an estimate of variance in the same units as the systemic effect values.</li><li>The RMSE can also be used to define a minimum prediction interval, or estimate range, for a model.</li></ul>	<ul style="list-style-type: none"><li>The mean square error (MSE) is used to approximate the unexplained variance (not explained by study descriptors).</li><li>This unexplained variance limits the R-squared on a new model.</li></ul>	<ul style="list-style-type: none"><li>Understand the reproducibility of treatment-related changes in specific endpoint targets (e.g., any effect on liver).</li></ul>



ToxRefDB v2.0 is a source for a dataset to address these questions of quantitative variability.



**Figure 1. Number of studies by study type and species in ToxRefDB v2.0.** The study designs include chronic (CHR), sub-chronic (SUB), developmental (DEV), subacute (SAC), multigeneration reproductive (MGR), developmental neurotoxicity (DNT), reproductive (REP), neurotoxicity (NEU), acute (ACU), and other (OTH) for numerous species, but mostly for rat, mouse, rabbit, and dog.

ToxRefDB v2.0 contains relevant study data to evaluate variability in traditional data for >1000 chemicals and >5000 studies.

Figure from Watford S, Pham LL, Wignall J, Shin R, Martin MT, Paul Friedman K. 2019. "ToxRefDB version 2.0: Improved utility for predictive and retrospective toxicology analyses." *Reproductive Toxicology*; 89: 145-158.

<https://doi.org/10.1016/j.reprotox.2019.07.012>



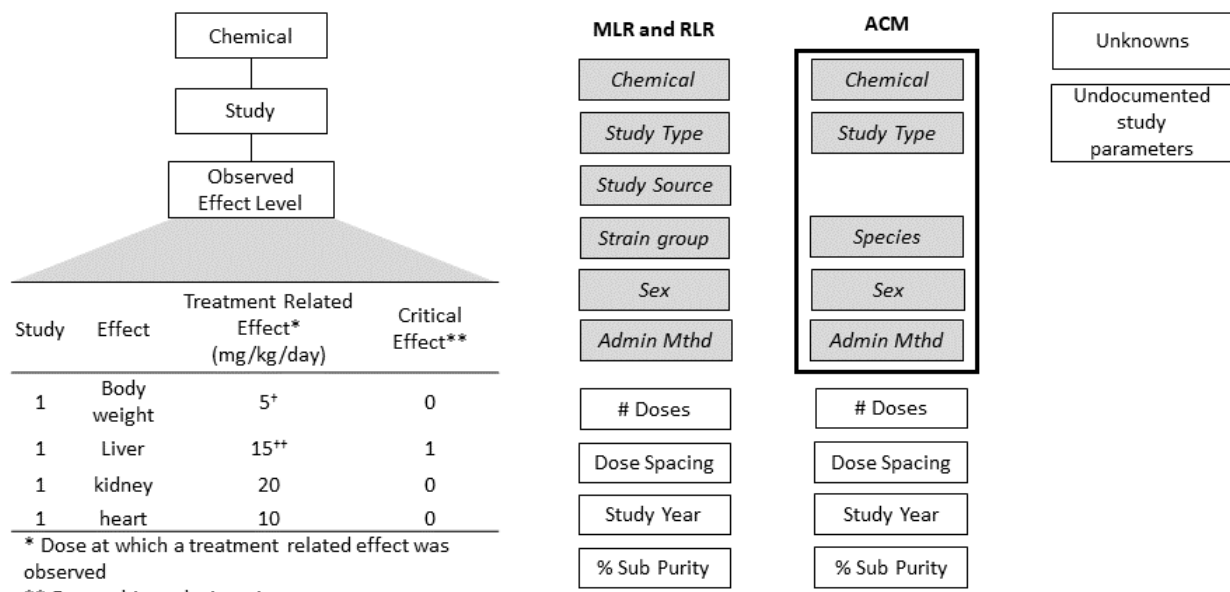
Based on the study descriptors in ToxRefDB v2.0, we developed statistical models of the variance in quantitative systemic effect level values.

**Total variance**

**Approximated by mean square error**

**Using two approaches:**

$$\text{Observed Variance (LEL or LOAELs)} = \text{Variance Explained by Study Parameters} + \text{Unexplained Variance}$$



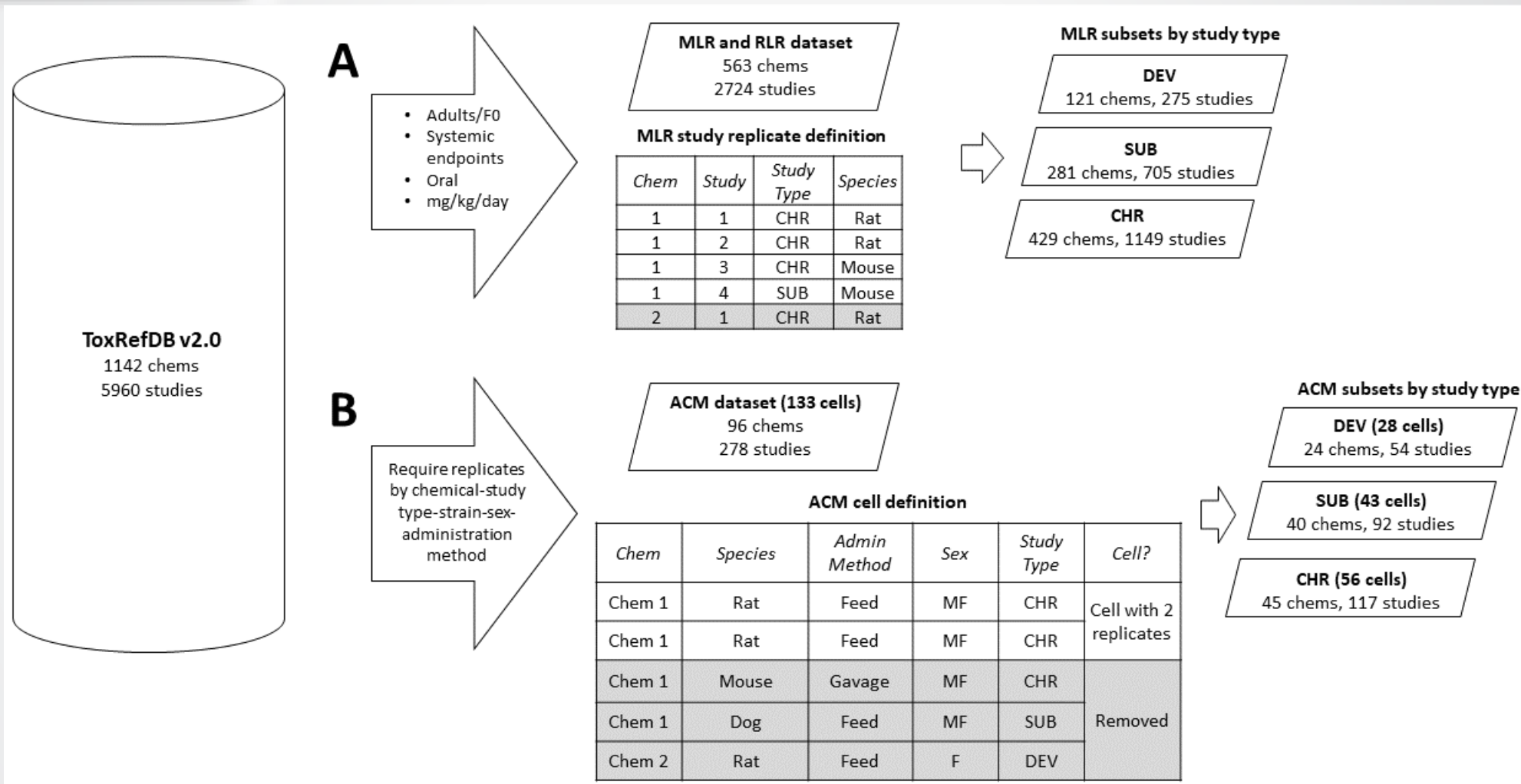
\* Dose at which a treatment related effect was observed  
 \*\* Expert driven designation  
 † Observed effect level used in LEL dataset  
 \*\* Observed effect level used in LOAEL dataset

	Multilinear regression (MLR, RLR)	Augmented cell means (ACM)
Aggregation level	Chemical	Chemical-Study Type-Species-Sex-Admin Method combination
Replicate definition stringency	Not stringent	Stringent
N	Maximized; ↓ impact of outliers/database error rate	Small; may bias variance estimate
Study descriptors	Contribute independently to variance	Accounts for possible interactions among descriptors

**Figure 2. Statistical model of the variance.** LEL = lowest effect level; LOAEL = lowest observable adverse effect level. The LEL is the lowest treatment-related effect observed for a given chemical in a study, and the LOAEL is defined by expert review as coinciding with the critical effect dose level from a given study. Multiple studies for a given chemical yield multiple LELs and LOAELs for computation of variance. MLR = multilinear regression; RLR = robust linear regression; ACM = augmented cell means; Adm. Method = administration method; % Sub Purity = % substance purity used in the study. The gray shaded study descriptor boxes are categorical variables, and the white study descriptor boxes are continuous variables. The box around five categorical study descriptors for the ACM indicates these were concatenated to a factor to define study replicates.

Pham LL, Watford S, Pradeep P, Martin MT, Thomas RS, Judson RS, Setzer RW, Paul Friedman K. *in Agency review*. “Variability in in vivo studies: Defining the upper limit of performance for predictions of systemic effect levels”

# Our workflow for evaluating variance in repeat dose toxicity information



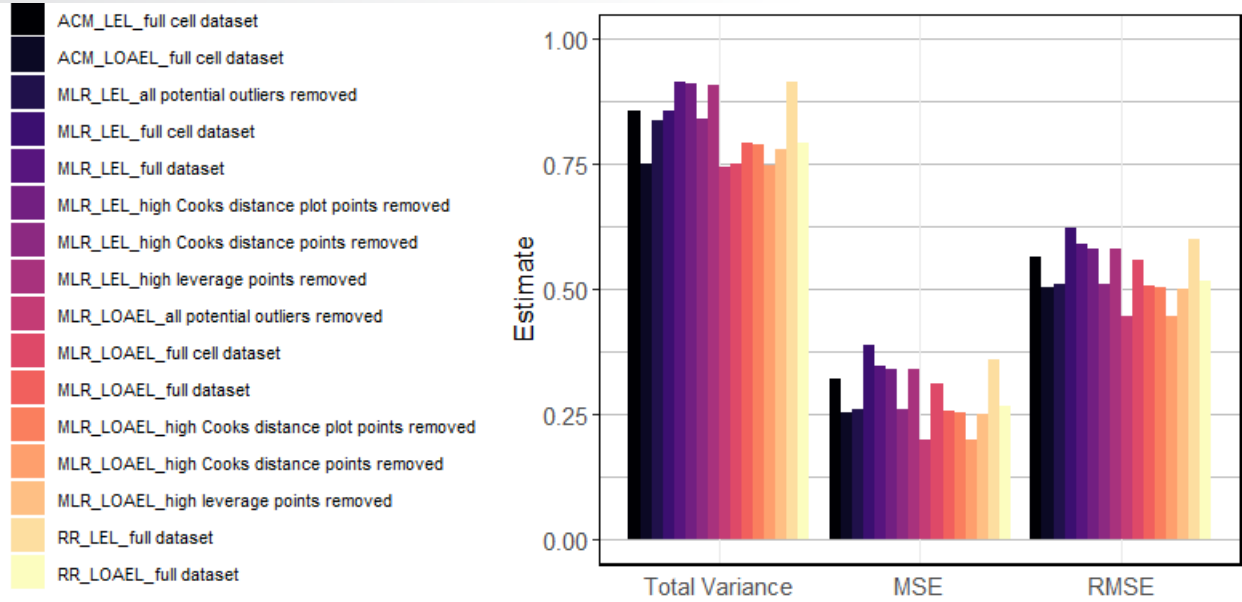
**Figure 1. Variance estimation workflow.**

CHR = chronic; DEV = developmental (adults only); SUB = subchronic; cells are defined by the factor of all categorical variables; MF = males and females; F = females; MLR = multilinear regression; RLR = robust linear regression; ACM = augmented cell means.

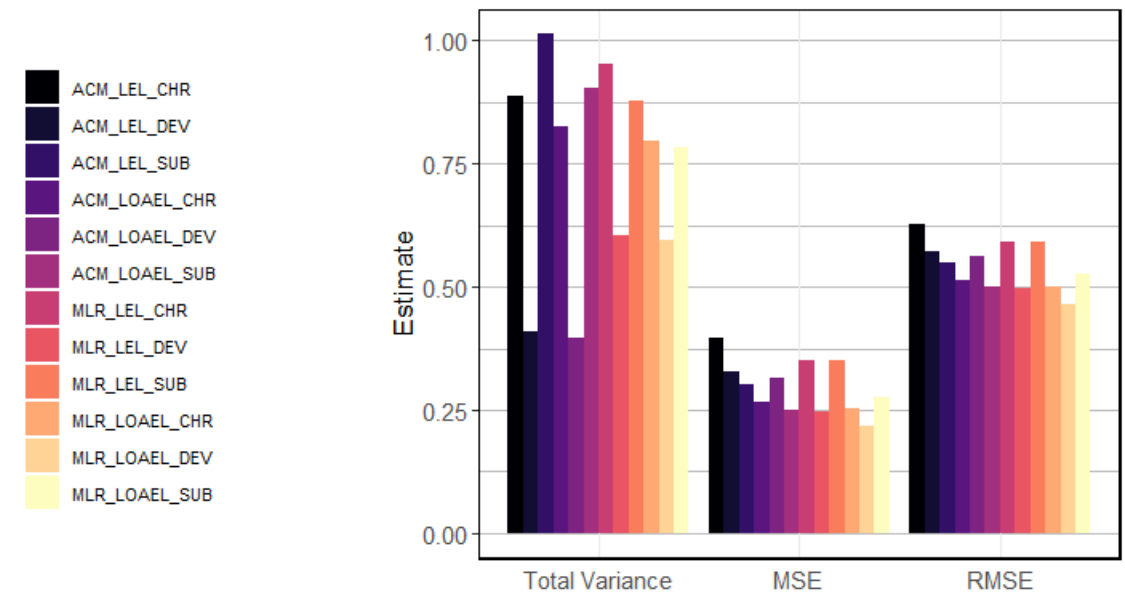


# 28 models to approximate total variance, unexplained variance (MSE), and then the spread of the residuals from the statistical models (RMSE)

### Statistical models for LELs and LOAELs for the full dataset



### Statistical models for LELs and LOAELs for datasets subset by study type



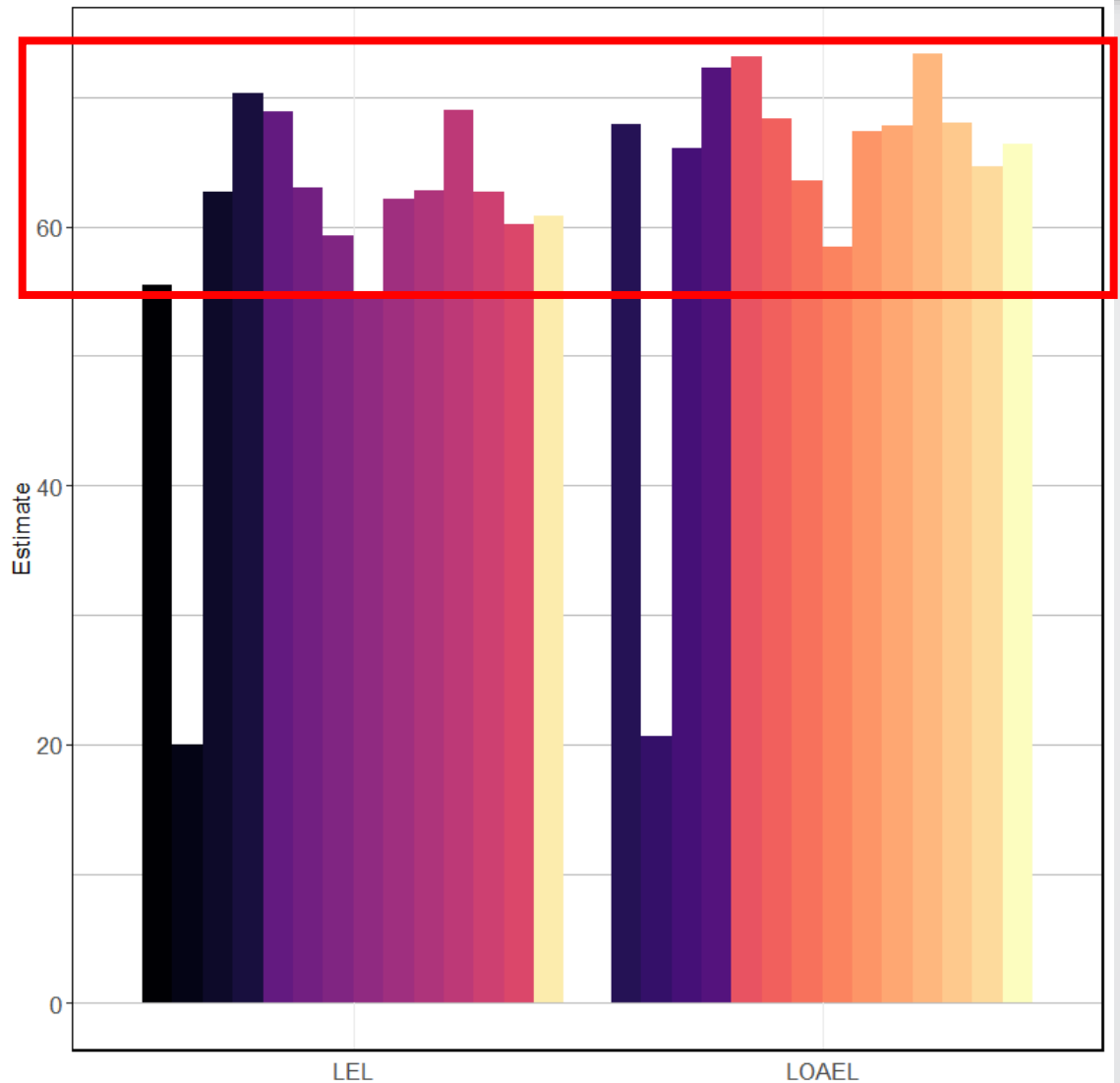
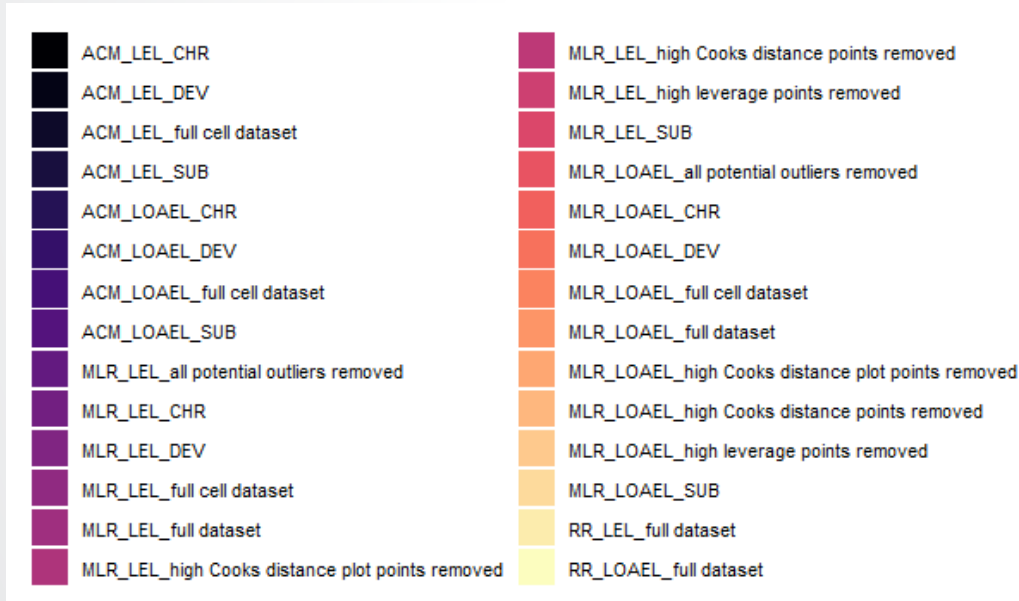
- Total variance in systemic toxicity effect values likely approaches 0.75-1 (units of  $(\log_{10}\text{-mg/kg/day})^2$ )
- MSE (unexplained variance) is 0.2 – 0.4 (units of  $(\log_{10}\text{-mg/kg/day})^2$ )
- RMSE is 0.45-0.60  $\log_{10}\text{-mg/kg/day}$
- RMSE is used to define a 95% minimum prediction interval (i.e., based on the standard deviation or spread of the residuals)





# Percent explained variance is also stable across statistical models.

- The % explained variance (amount explained by study descriptors) likely approaches 55-73%.
- This means that the  $R^2$  on some new, predictive model would approach 0.55 to 0.73 as an upper bound on accuracy.



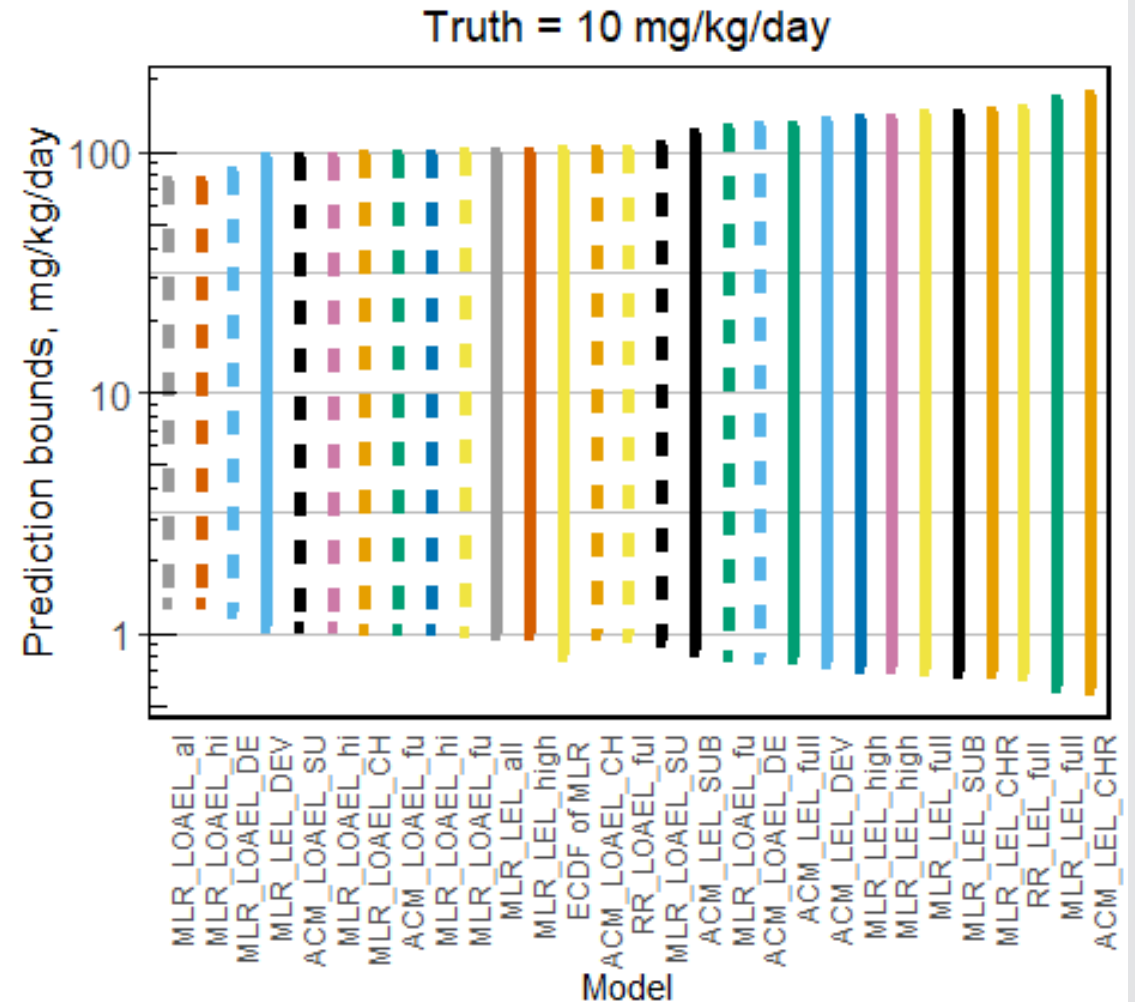
Based on tables in Pham LL, Watford S, Pradeep P, Martin MT, Thomas RS, Judson RS, Setzer RW, Paul Friedman K. *in Agency review*. "Variability in in vivo studies: Defining the upper limit of performance for predictions of systemic effect levels"



## Range of 95% minimum prediction intervals across the modeling approaches, effect levels, and study types is 58-284-fold

If attempting to use a NAM-based predictive model for prediction of a reference systemic effect level value of 10 mg/kg/day, it is likely that given the variability in reference data of this kind, that a model prediction of somewhere between 1 and 100 mg/kg/day would be the greatest amount of accuracy achievable.

- CHR
- DEV
- SUB
- full cell dataset
- full dataset
- all potential outliers removed
- high Cooks distance plot points removed
- high Cooks distance points removed
- high leverage points removed
- LEL
- LOAEL



Based on tables in Pham LL, Watford S, Pradeep P, Martin MT, Judson RS, Setzer RW, Paul Friedman K. *in Agency review*. "Variability in in vivo studies: Defining the upper limit of performance for predictions of systemic effect levels"

- Variability in *in vivo* toxicity studies limits predictive accuracy of NAMs.
- Total variance in systemic effect levels and the fraction explained were quantified.
- Maximal R-squared for a NAM-based predictive model of systemic effect levels may be 55 to 73%; i.e., as much as 1/3 of the variance in these data may not be explainable using study descriptors.
- The estimate of variance (RMSE) in curated LELs and/or LOAELs approaches a 0.5 log<sub>10</sub>-mg/kg/day.
- Estimated minimum prediction intervals for systemic effect levels are likely 58 to 284-fold based on RMSE estimates.
  - This suggests that the current LOAEL-NOAEL uncertainty factor ( $UF_L$ ) covers the estimated one-sided minimum prediction interval.



## How does this compare to previous work in this area?

- Previous QSAR models of subchronic oral rat NOAEL values:  $R^2$  approaches 0.46-0.71, i.e. 46-71% of residual variance could be explained for the reference set (Veselinovic et al. 2016; Toropov et al. 2015; Toropova et al. 2017).
- A multi-linear regression QSAR model of chronic oral rat LOAEL values for approximately 400 chemicals, demonstrated a RMSE of  $0.73 \log_{10}(\text{mg/kg-day})$ , which was similar to the size of the variability in the training data,  $\pm 0.64 \log_{10}(\text{mg/kg-day})$ , suggested that the error in the model approached the error in the reference data from different laboratories (Mazzatorta et al. 2008; Helma et al. 2018).

**Few examples of quantitative variability in this domain to cite, but suggest that similar thresholds of 50-70% explained variance and RMSE of 0.5-0.7 may exist in other larger reference data sets for systemic toxicity in subchronic and chronic animal studies.**

*"We can learn from history, but we can also deceive ourselves when we selectively take evidence from the past to justify what we have already made up our minds to do."*

Dr. Margaret MacMillan, professor and expert at University of Oxford

- **Understanding that a prediction of an animal systemic effect level within  $\pm 1$  log<sub>10</sub>-mg/kg/day fold demonstrates a *very good* NAM is important for acceptance of NAMs for chemical safety assessment.**
- Use of NAM-based information in chemical safety assessment should emphasize NAM value for screening level chemical assessment beyond the accuracy of NAMs for prediction of *in vivo* data from animal models.
- Finally, construction of NAM-based effect level estimates that offer an equivalent level of public health protection as effect levels produced by methods using animals may provide a bridge to major reduction in the use of animals as well as identification of cases in which animals may provide scientific value.



# Thank you for listening

## References

- Congress, U. S., FRANK R. LAUTENBERG CHEMICAL SAFETY FOR THE 21ST CENTURY ACT. In: Congress, (Ed.), H.R.2576, Vol. Public Law 114-182, 2016.
- Dumont, C., et al. (2016). "Analysis of the Local Lymph Node Assay (LLNA) variability for assessing the prediction of skin sensitisation potential and potency of chemicals with non-animal approaches." Toxicol In Vitro 34: 220-228.
- Gold, L. S., et al. (1989). "Interspecies extrapolation in carcinogenesis: prediction between rats and mice." Environ Health Perspect 81: 211-219.
- Gottmann, E., et al., 2001. Data quality in predictive toxicology: Reproducibility of rodent carcinogenicity experiments. Environmental Health Perspectives. 109, 509-514.
- Haseman, J. K. (2000). "Using the NTP database to assess the value of rodent carcinogenicity studies for determining human cancer risk." Drug Metab Rev 32(2): 169-186.
- Mazzatorta, P., et al., 2008. Modeling Oral Rat Chronic Toxicity. Journal of Chemical Information and Modeling. 48, 1949-1954.
- Monticello, T. M., et al. (2017). "Current nonclinical testing paradigm enables safe entry to First-In-Human clinical trials: The IQ consortium nonclinical to clinical translational database." Toxicol Appl Pharmacol 334: 100-109.
- Toropov, A. A., et al., 2015. CORAL: model for no observed adverse effect level (NOAEL). Molecular diversity. 19, 563-75.
- Toropova, A. P., et al., 2017. The application of new HARD-descriptor available from the CORAL software to building up NOAEL models. Food and Chemical Toxicology.
- Toropova, A. P., et al., 2015. QSAR as a random event: a case of NOAEL. Environ Sci Pollut Res Int. 22, 8264-71.
- Veselinović, J. B., et al., 2016. The Monte Carlo technique as a tool to predict LOAEL. European Journal of Medicinal Chemistry. 116, 71-75.
- Wang, B. and G. Gray (2015). "Concordance of Noncarcinogenic Endpoints in Rodent Chemical Bioassays." Risk Anal 35(6): 1154-1166.
- Watford, S., et al., 2019. ToxRefDB version 2.0: Improved utility for predictive and retrospective toxicology analyses. Reprod Toxicol. 89, 145-158.
- Wheeler, A. R., Memorandum: Directive to Prioritize Efforts to Reduce Animal Testing. US Environmental Protection Agency, Washington, D.C., 2019.



**Office of Research and Development  
Center for Computational Toxicology & Exposure (CCTE)  
Bioinformatic and Computational Toxicology Division  
(BCTD)  
Computational Toxicology and Bioinformatics Branch (CTBB)**



# Questions/appendix



## Future work: what about qualitative reproducibility or binary classification?

- Local lymph node assay (LLNA): with same species & vehicle solvent, repeat LLNA were concordant only **78%** of the time, with a 35% chance that a “negative” chemical would test “positive” if the LLNA was repeated (Hoffmann et al., 2018; Dumont et al., 2016).
- Kleinstreuer and colleagues (2014) showed that even in high quality studies for the rodent uterotrophic bioactivity assay, concordance was only achieved **74%** of the time for replicate uterotrophic assays.
- An evaluation of 37 National Toxicology Program repeat dose toxicity studies demonstrated 0-100% concordance, with a median of approximately **70%**, in the non-carcinogenic effects observed between rats and mice, depending on the biological endpoint or tissue measured (Wang & Gray, 2015).
- Concordance among rat and mouse models of carcinogenicity has been shown to range from **57% to 76%** (Gottman et al. 2001; Gold et al., 1989; Haseman 2000).
- The negative prediction value of animal preclinical data may be more informative for (negative) clinical outcome prediction (than positives) (Monticello et al. 2017).





# Is it reasonable to combine study types in this analysis?

- In the MLR regression models, study type is a covariate.
- In the ACM models, study type is used to define each cell.
- The consistency in the observed MSE (0.217 to 0.395 for study type subsets versus 0.200 to 0.387 for full datasets) builds confidence in the estimates resulting from the analyses that have combined multiple study types, especially SUB and CHR data.

Indeed, the difference between the SUB and CHR LOAEL values *on average* is less than the variance in SUB or CHR data.

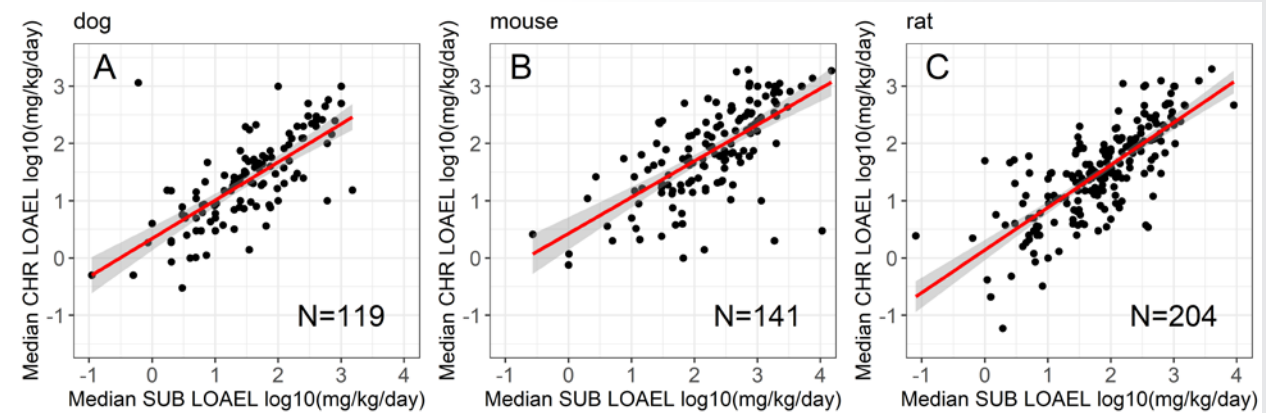


Figure 5. Median chronic versus median subchronic LOAEL values for the ACM datasets.

**For a NAM-based predictive model of animal systemic, repeat dose toxicity effect level, it seems advantageous to train on SUB and CHR reference data in order to provide the best model for prediction.**



# Minimum prediction intervals were estimated using the RMSE from the statistical models and an ECDF.

- ECDF = empirical cumulative distribution function
- The use of RMSE to compute a 95% minimum prediction interval relies on an assumption of normally distributed residuals.
- This assumption was rigorously evaluated using regression diagnostic plots, removal of potential high residual and high leverage points, and by comparison of the ECDF of the MLR full LEL dataset model residuals to a standard normal distribution of the same sample size, mean, and standard deviation.

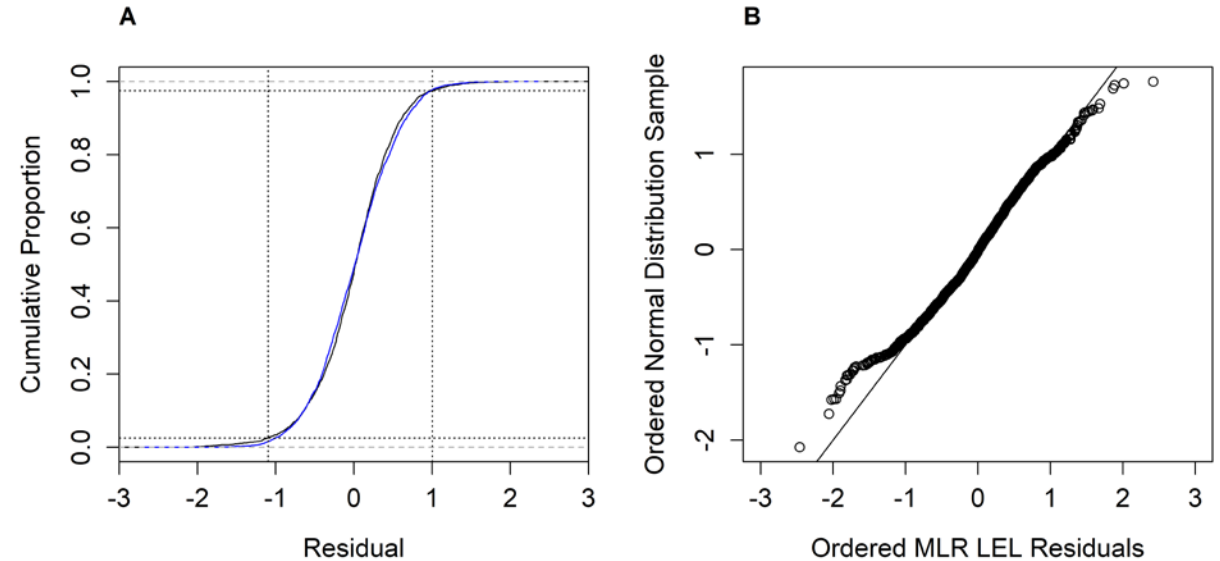


Figure 6. ECDF for MLR LEL model residuals.

*Following evaluation through these different methods, and given that the distribution of residuals from the MLR full dataset model did not appear to deviate substantially from a normal distribution, minimum prediction intervals for a new POD value based on RMSE values from regression models in this work were derived using the assumption of normality.*