

Validation of Prediction Uncertainty in Computational Chemistry

Pascal PERNOT

Institut de Chimie Physique,
CNRS & Univ. Paris-Saclay



CECAM (2022-06-22)

- 1 Uncertainty estimation in computational chemistry
- 2 Prediction uncertainty validation framework
- 3 Examples from the CC literature
- 4 Conclusion

CC-UQ validation practices

- **Qualitative appreciation** of conformity between PU and errors amplitude
 - visual check of normality of z -scores histogram¹
 - visual estimation of local coverage of 95% prediction interval (PI)²
- **Statistical estimation**
 - coverage of 95% prediction intervals³
 - correlation of uncertainty and absolute errors⁴
 - *calibration/sharpness* (for CC-applied ML methods)⁵

We need a consistent and shared validation framework !

¹Mortensen *et al.* (2005) *Phys. Rev. Lett.* (<https://tinyurl.com/mvwk3fff>)

²Bakowies and von Lilienfeld (2021) *JCTC* (<https://tinyurl.com/ms3dy7yv>)

³Pernot *et al.* (2015) *J. Chem. Phys.* (<https://tinyurl.com/3c9aw28r>), Proppe & Kircher (2022) *ChemPhysChem* (<https://tinyurl.com/yckxvjkk>)

⁴Zheng *et al.* (2022) *J. Phys. Chem. Lett.* (<https://tinyurl.com/ccefk79z>)

⁵Tran *et al.* (2020) *Mach. Learn.: Sci. Technol.* (<https://tinyurl.com/2p849fs6>), Scalia *et al.* (2020) *J. Chem. Inf. Model.* (<https://tinyurl.com/yc7rn7dp>)

Uncertainty vs error⁶

In order to estimate a measurement uncertainty it is assumed that the result of a measurement has been corrected for all recognized significant systematic effects and that every effort has been made to identify such effects.

systematic error

component of measurement error that remains constant or varies in a **predictable** manner

random error

component of measurement error that varies in an **unpredictable** manner

uncertainty

non-negative parameter characterizing the **dispersion** of the quantity values being attributed to a measurand

- often estimated by a standard deviation u , or the half-width of a probability interval U_p

⁶Guide to the expression of uncertainty in measurement (GUM), JCGM 100:2008, International Vocabulary of Metrology (VIM), JCGM 200:2012

Error sources in Computational Chemistry⁹

● Numerical errors

- finite arithmetics, stochastic algorithms. . .
- mostly *random errors*; assumed to be well controlled⁷, except for numerical chaos⁸

● Parametric uncertainty

- semi-empirical methods, statistical corrections. . .
- *random errors*; decrease with size of calibration set

● Model errors

- level-of-theory errors (density functional approximation, correlation level, force-field expression. . .), representation errors (basis-sets, grids). . .
- mostly *systematic errors*; often the dominant error source;
no reason to be normally distributed

⁷Irikura *et al.* (2004) *Metrologia* **41**:369

⁸Feher and Williams (2012) *J. Chem. Inf. Model.* **52**:3200-3212

⁹Lejaeghere (2020) *Uncertainty Quantification in Multiscale Materials Modeling*, pp. 41–46

CC-UQ outputs

- **Prediction distributions or representative samples**
 - available for some methods (Stochastic methods, Statistical models, Bayesian Ensembles. . .)
- Most UQ studies in the CC literature provide statistical summaries:
 - **expanded uncertainties** (U_p , typically for $p = 0.95$)¹⁰
 - **standard uncertainties** (u)

Note: no prediction interval without distribution hypothesis

¹⁰Ruscic (2014) *Int. J. Quantum Chem.* **114**:1097

Validation data sets

Let us consider a typical validation set

- V_i : predicted value at point $i \in 1 : M$
- u_{V_i} : uncertainty on V_i (*model* uncertainty)
- R_i : reference value
- u_{R_i} : uncertainty on R_i (*data* uncertainty)

Validation is based on

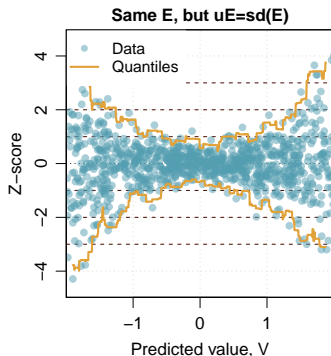
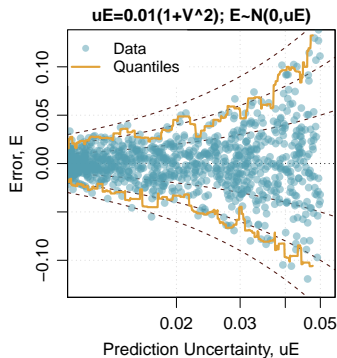
- $E_i = R_i - V_i$: error / prediction error
- for standard uncertainties $u_{E_i} = \sqrt{u_{V_i}^2 + u_{R_i}^2}$ (*prediction* uncertainty)

Prediction uncertainty quantifies the dispersion of pred. errors

Visual appreciation

Do the errors scale with uncertainty ?

- 1 if $u_E \neq c^{te}$, plot E vs u_E and guiding lines $y = k * x$
- 2 if $u_E = c^{te}$, plot E/u_E (z-score) vs V and guiding lines $y = k$
- 3 as helper, plot *running quantiles* (CI₉₅)



The Calibration/Sharpness framework¹²

Calibration a method is considered to be *calibrated* if the confidence of predictions matches the probability of being correct for all confidence levels

Sharpness the concentration of a predictive distribution in absolute terms. Conditional to calibration

Pb: sharpness is a property of the forecast alone and does not involve the observations.

→ useful in benchmarking, not in validation. . .

Tightness¹¹ a method is considered to be *tight* if it is calibrated for any relevant subgroup of the validation data (small-scale calibration)

¹¹Pernot (2022) (<https://arxiv.org/abs/2204.13477>)

¹²Gneiting et al. (2007) *Stat. Meth. B* (<https://tinyurl.com/2p8nr3ab>)

Limits of average calibration

- Average calibration does not guarantee tightness
 - in benchmarking, *sharpness* is used to select tighter forecasts
- Stronger calibration modes have been introduced:
 - *group* calibration¹³ where calibration is assessed on relevant subgroups of the validation dataset
 - *adversarial* group calibration¹⁴ where calibration is assessed on any random group of useful size
 - *perfect* calibration¹⁵
 - I propose to use **local** calibration, a variant of group calibration, where the validation set is split into contiguous areas of a chosen coordinate (predicted value, prediction uncertainty. . .)

¹³Chung *et al.* (2021) arXiv:2109.10254; H bert-Johnson (2017) arXiv:1711.08513

¹⁴Zhao (2020) arXiv:2006.10288

¹⁵Levi *et al.* (2020) (<http://arxiv.org/abs/1905.11659>)

Different validation approaches¹⁸

- **Interval-based**¹⁶

$$\lim_{M \rightarrow \infty} \frac{1}{M} \sum_{i=1}^M \mathbf{1}(E_i \in I_{E_i,p}) = p, \forall p \in [0, 1]$$

where $I_{E_i,p}$ is a prediction interval at probability level p

- **Variance-based**¹⁷

$$\text{Var}(E | u_E^2 = \sigma^2) = \sigma^2, \forall \sigma^2$$

which *does not* operate at the same level as interval-based validation

- **Note:** *ranking-based* methods (correlation between u_E and $|E|$; confidence curves) cannot validate calibration/tightness, but can invalidate tightness. . .

¹⁶Kuleshov *et al.* (2018) (<http://arxiv.org/abs/1807.00263>)

¹⁷Levi *et al.* (2020) (<http://arxiv.org/abs/1905.11659>)

¹⁸Scalia *et al.* (2020) *J. Chem. Inf. Model.* (<https://tinyurl.com/yc7rn7dp>), Pernot (2022) *J. Chem.*

Interval-based validation

An *expanded uncertainty* is the half-width of a prediction interval

$$I_{E_i,p} = [-U_{E_i,p}, U_{E_i,p}]$$

To validate $U_{E,p}$, one should therefore test

$$p \stackrel{?}{\in} I_{95}(\nu_p, M), \text{ where } \nu_p = \frac{1}{M} \sum_{i=1}^M \mathbf{1}(|E_i| \leq U_{p,i})$$

and ν_p is a PICP (Prediction Interval Coverage Probability)

- $I_{95}(\nu_p, M)$ (Binomial Proportion CI) is estimated by a method avoiding normality hypothesis (Clopper-Pearson, cc-Wilson, Agresti-Coull...)¹⁹

¹⁹Vollset (1993) *Stat. Med.* (<https://tinyurl.com/5dps8u3h>)

Variance-based validation

Let us consider unbiased errors of *unknown* distribution

$$E(E_i) = 0; \text{Var}(E_i) = u_{E_i}^2$$

Then for z-scores $z_i = E_i/u_{E_i}$ one has

$$E(z_i) = 0; \text{Var}(z_i) = 1$$

and for a set of z-scores $Z = \{z_i\}_{i=1}^M$

$$E(Z) = 0; \text{Var}(Z) = 1$$

To validate $\text{Var}(Z)$, one should therefore test

$$1 \in I_{95}^?(\text{Var}(Z), M)$$

- $I_{95}(\text{Var}(Z), M)$ is estimated by an adapted bootstrap method (BC_a, ABC...) ²⁰
to avoid the normality-based textbook method

²⁰Diciccio and Efron (1996) *Stat. Sci.* (<https://tinyurl.com/ssztxy6k>)

How to test tightness ?

- **Proposed approach:** interpret tightness as **local calibration** and use calibration tests on subsets of the validation set wrt V or uE (or any other relevant property)
 - LCP analysis: local PICP estimation and test
 - LZV analysis: local z-scores variance testing
- **Pb:** partitioning reduces sample size (bad for test power)
→ use overlapping/sliding areas for small datasets.
Trends in LCP-LZV curves help diagnostic.
- Link of LZV/ uE with *perfect* calibration (reliability diagram, RD)²¹

$$\text{Var}(E|u_E^2 = \sigma^2) = \sigma^2, \forall \sigma^2$$

but RD does not deal with homoscedastic datasets.

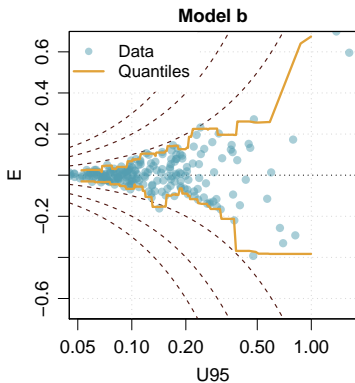
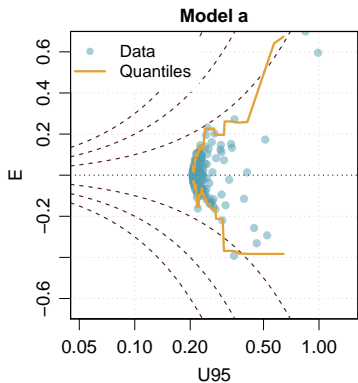
²¹Levi *et al.* (2020) (<http://arxiv.org/abs/1905.11659>)

Overview

Diagnostic	Applicability					Validation	
	q_E	u_E	$U_{E,p}$	Homosc.	Heterosc.	Calibrat.	Tightness
<i>Average</i>							
PIT hist.	✓	✗	✗	✓	✓	✓	✗
Calib. curve	✓	✗	✗	✓	✓	✓	✗
PICP	✓	✗	✓	✓	✓	✓	✗
Var(Z)	✓	✓	✗	✓	✓	✓	✗
Cor($u_E, E $)	✓	✓	✓	✗	✓	✗	✗*
<i>Local</i>							
LCP/LRR	✓	✗	✓	✓	✓	✓†	✓
LZV	✓	✓	✗	✓	✓	✓†	✓
Reliab. diag.	✓	✓	✗	✗	✓	✓†	✓
Confid. curve	✓	✓	✓	✗	✓	✗	✗*

Reactivity Scales²²

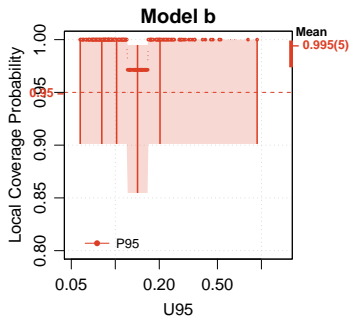
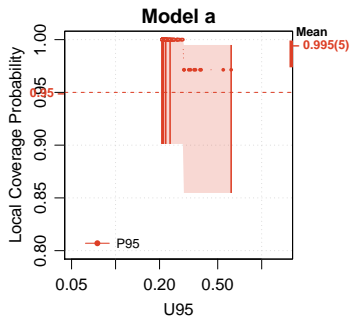
Set of 212 errors for reaction rates by an extended Mayr's reactivity scale and U_{95} values obtained by two UQ methods (a and b).



²²Proppe & Kircher (2022) *ChemPhysChem* (<https://tinyurl.com/yckxvjkk>)

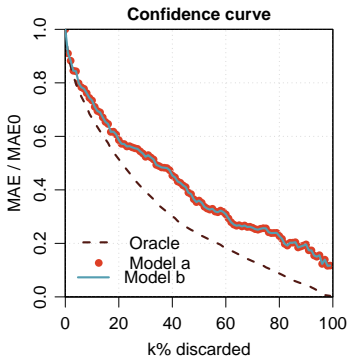
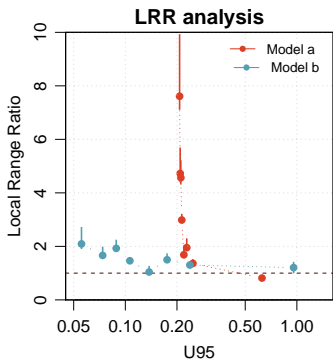
Reactivity Scales²³

$U_{95} \Rightarrow$ interval-based validation



- large local uncertainties because of small groups
- for overestimated uncertainties, the PICP test saturates at 1

²³Proppe & Kircher (2022) *ChemPhysChem* (<https://tinyurl.com/yckxvjkk>)

Reactivity Scales²⁴

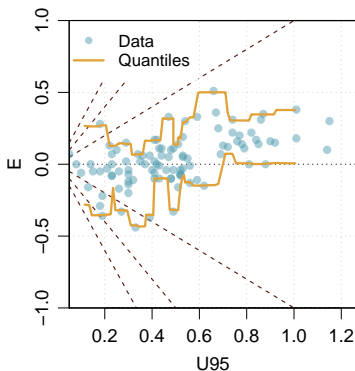
- **Range Ratio** : Mean width of pred. int. / Width of error proba. int.
- even if calibration is rejected, on might reliably use these uncertainties for active learning (Conf. curves)

²⁴Proppe & Kircher (2022) *ChemPhysChem* (<https://tinyurl.com/yckxvjkk>)

ZPE by the ATOMIC-2 composite method²⁵

A-posteriori estimation of U_{95} based on a correlation of errors with the fraction of heteroatoms in a molecule.

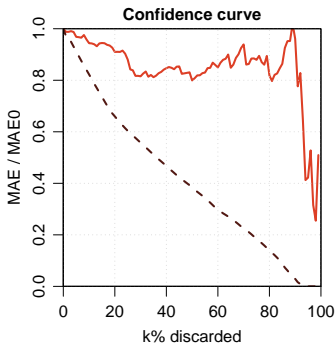
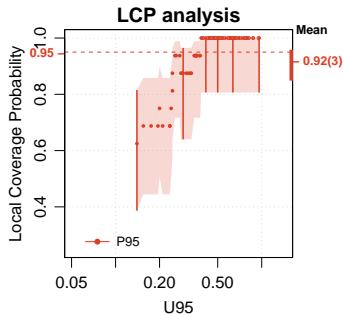
- Small test dataset: $M = 99$
 $\{V_i, R_i, U_{95}, V_i\}$
- Reference data: CCSD(T) (no uncertainty)
- Authors validate by visual appreciation of error bars



²⁵Bakowies and von Lilienfeld (2021) *JCTC* (<https://tinyurl.com/ms3dy7yv>), Pernot (2022) *J Chem Phys* (<https://doi.org/10.1063/5.0084302>)

ZPE by the ATOMIC composite method

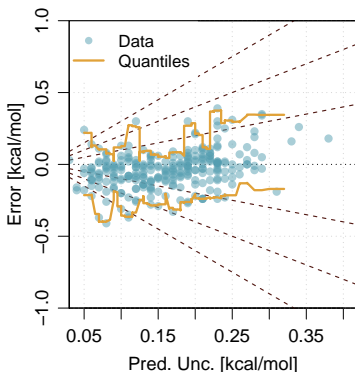
PICP testing at $p = 0.95$



- $\nu_{0.95} = 0.92(3)$: average calibration is OK
- large uncertainty on PICPs, but the **trends are informative**
- from the LCP analysis, one sees a systematic bias:
small PUs are underestimated, large ones are overestimated

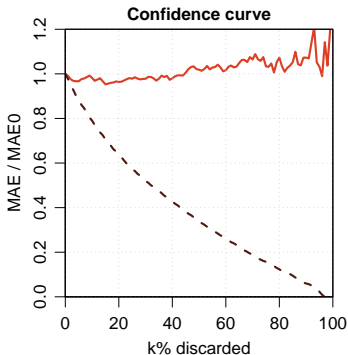
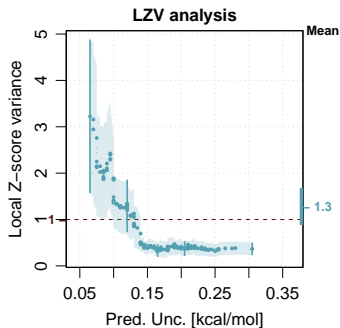
Formation heats by the mBEEF method²⁶

- Bayesian Ensembles method inflates *parametric uncertainty* of exchange-correlation model to cover errors amplitude
 - strong functional constraints: tightness ???
 - does not disambiguate model uncertainty from reference data uncertainty
- Set of $M = 257 \{V_i, R_i, u_{V_i}\}$
- R_i experimental, no uncertainty provided



²⁶Pandey and Jacobsen (2015) *Phys. Rev. B* (<https://tinyurl.com/5dv9spnn>), Pernot (2017) *J. Chem. Phys.* (<https://tinyurl.com/yb6uzwzr>), Pernot and Cailliez (2017) *AIChE J.* (<https://tinyurl.com/2xxcfs2f>)

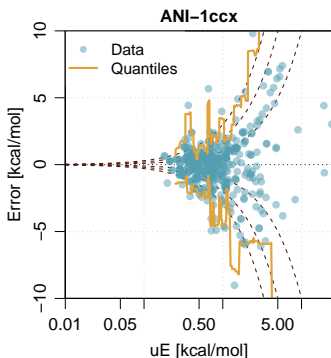
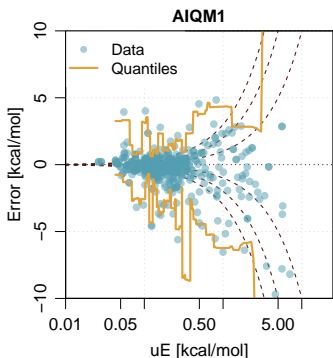
Formation heats by the mBEEF method



- $\text{Var}(Z) = 1.3(2)$, average calibration OK
- the LZV analysis shows that small PUs are underestimated, while large ones are overestimated
- the confidence curve is catastrophic
- these uncertainties should not be trusted

Query by Committee UQ

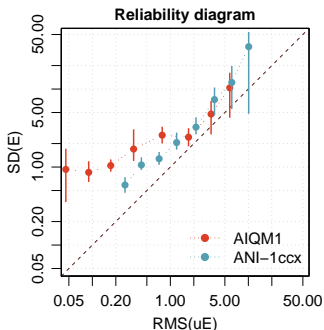
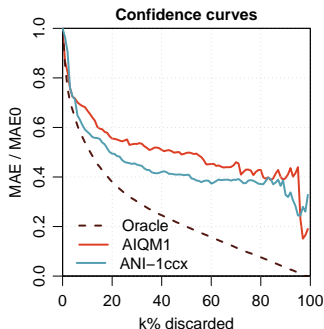
QbC uncertainties ($n = 8$) on formation enthalpies for AIQM1 and ANI-1ccx²⁷. Set of $M = 472 \{E_i, u_{V_i}\}$



²⁷Zheng et al. (2022) *J. Phys. Chem. Lett.* (<https://tinyurl.com/ccefk79z>)

Query by Committee UQ

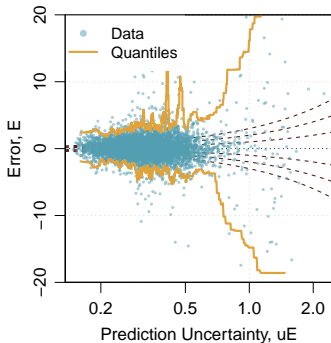
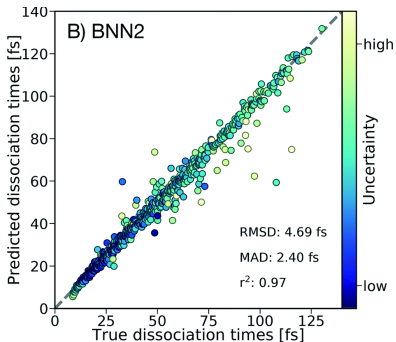
Stat	AIQM1	ANI-1ccx	Target
$Var(Z)$	59	4.3	1.4



- QbC does not provide a prediction uncertainty
- but both methods point reliably to largest errors (good for active learning !)

Bayesian Neural Network

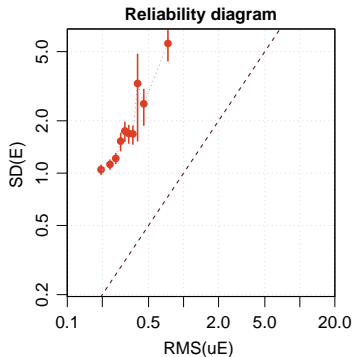
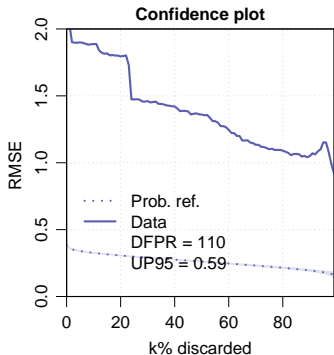
Data issued from a BNN trained to predict a MD potential²⁸ ($M = 5923$)



- The color scale for uncertainty is not a proper tool for validation

²⁸H se *et al.* (2019) *Chem. Sci.* **10**:2298

Bayesian Neural Network



- This BNN uncertainty is NOT a prediction uncertainty ($Var(Z) = 30$) but could still be used for active learning...

Calibration is easy, tightness is tough !

- **Calibration/Tightness : a principled framework for UQ validation**
- **CC-adapted C/T validation methods**
 - *standard PU* : test z -scores variance (LZV analysis), or build RD
 - *expanded PU* : test PICP values (LCP/LRR analysis)
- **CC-UQ methods rarely reach calibration and/or tightness**
 - datasets often too small for solid conclusions
 - should we loosen the validation criteria ?
 - how-much mis-calibration/mis-tightness is acceptable for a given application ? (e.g. calibration is not useful for active learning. . .)

Warmful thanks to...

- **Andreas SAVIN** (LCT, Jussieu)
for enlightening discussions
- **Jonny Proppe** (Univ. Göttingen)
for the PRO2022 dataset
- **Morgane Vacher** (Univ. Nantes)
for the HAS2019 dataset
- and **YOU**
for your attention !