

Statistical approaches to forcefield calibration and prediction uncertainty in molecular simulation

Fabien Cailliez, Pascal Pernot

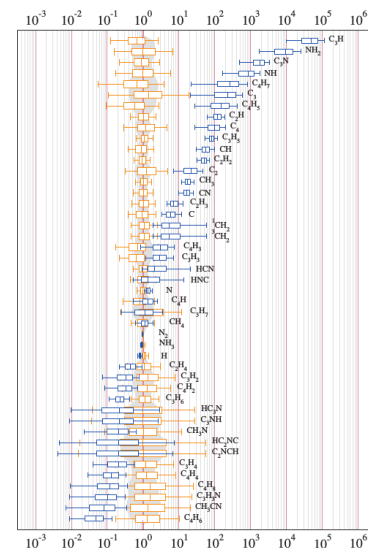
Laboratoire de Chimie Physique, CNRS UMR8000

Who are we?

- Laboratoire de Chimie Physique d'Orsay:
 - RISMAS (Réactivité des Ions, Spectrométrie de Masse, Analyse et Spectroscopies)
 - Biophysique
 - TEMiC (Transfert d'Electrons en Milieu Condensé)
 - **ThéoSim (Théorie et Simulation)**



- Théosim group:
 - Quantum Dynamics and Quantum chemistry
 - Molecular simulations of fluids
 - Molecular simulations of biophysical processes
 - **Management of uncertainties in chemical physics models (P. Pernot & F. Cailliez)**

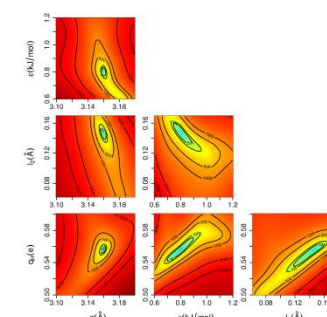
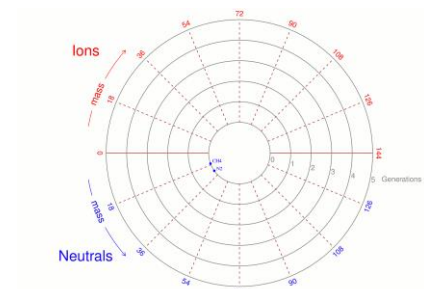
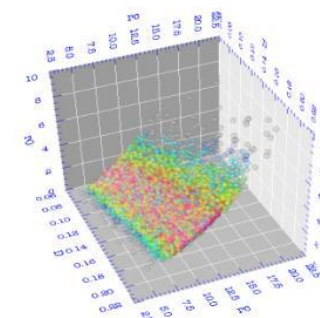


Who are we?

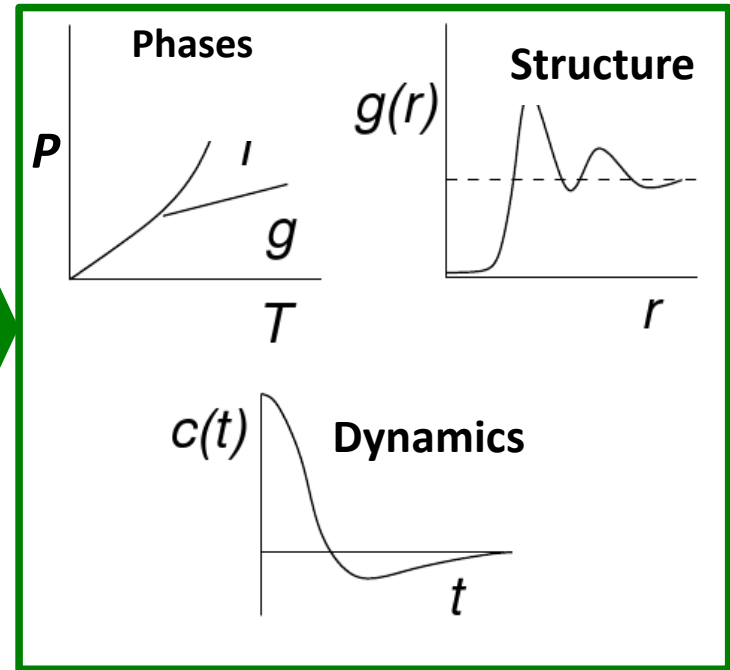
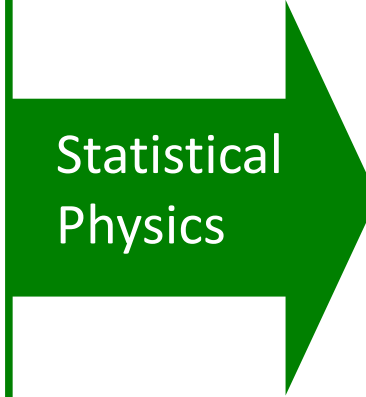
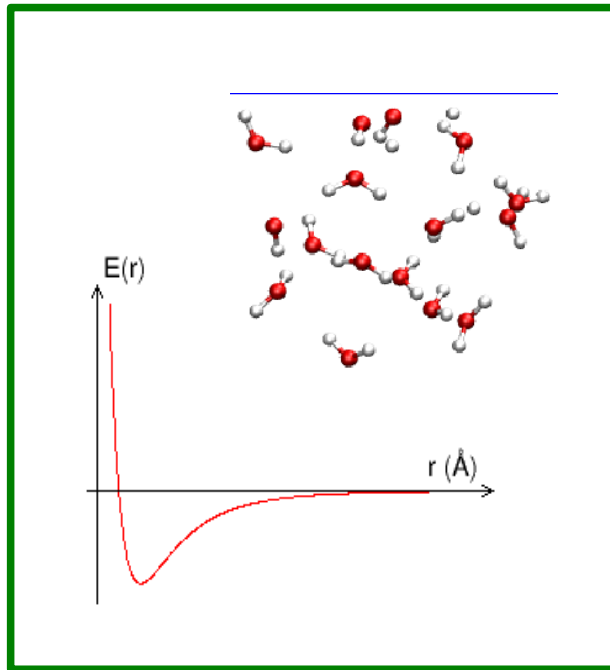
- Bayesian analysis of time-resolved spectroscopic data:
 - Determine chemical reaction schemes
 - Estimation of chemical reaction rates

- Modeling of the chemical complexification in Titan ionosphere:
 - Interpretation of experimental data
 - Use of data analysis to identify the formation of complex species and underlying mechanisms

- Representation and management of uncertainties in physico-chemical models:
 - Statistical treatment of the reactivity in planetary atmospheres
 - Calibration/Prediction uncertainties in theoretical chemistry: scaling factors for vibrational ZPE, DFT parameters, molecular simulation forcefields



Molecular Simulations



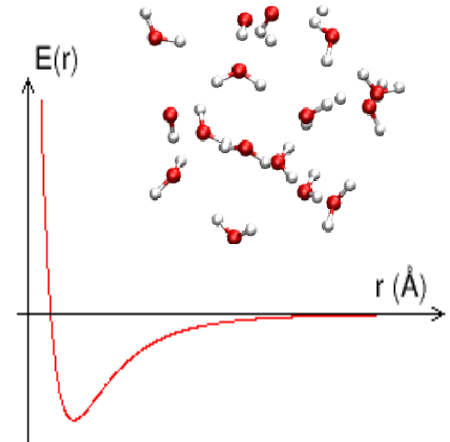
Scale	Laws of evolution	Energetic description
Electronic / Atomic	Schrödinger equation	<i>Ab initio</i> , DFT methods
Atomic / Molecular	Newtonian dynamics	Forcefields
Macromolecular	Langevin dynamics	Coarse-grained forcefields

Atomistic simulations and forcefields

- Principle of molecular simulations:
 - Sampling of representative configurations of the system (MD or MC)
 - Computation of macroscopic properties (laws of statistical physics)

- **Forcefield**: mathematical expression of the interatomic potential as a function of the nuclei positions

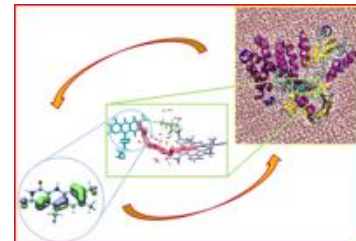
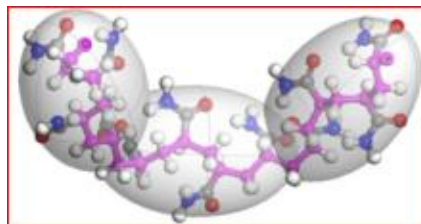
$$E(r_{ij}) = 4\varepsilon_{ij} \left[\left(\frac{\sigma_{ij}}{r_{ij}} \right)^{12} - \left(\frac{\sigma_{ij}}{r_{ij}} \right)^6 \right] + \frac{q_i q_j}{4\pi\varepsilon_0 r_{ij}} + \dots$$



- Characteristics of molecular simulations:
 - **Costly** evaluation of the properties (few hours to few days)
 - Number of (forcefield) parameters can be big
 - Stochastic outputs

Uncertainties in molecular simulations

- Purpose of Molecular Simulations:
 - Qualitative use: understand and interpret the behaviour of molecular systems
 - Quantitative use: predict properties
- Importance of monitoring uncertainties in molecular simulations:
 - In industry: molecular simulation used as a decision tool
→ **Confidence interval for the prediction needed**
 - In academy: development of multi-scale simulations
→ **Transfer of uncertainties along the different scales**



Uncertainties in molecular simulations

- Numerical uncertainties:

- limited sampling
- algorithmic parameters (cutoff radius,...)

Can be monitored and are reduced by the increase in computational power

- Parametric uncertainties:

Forcefield parameters most of the time calibrated over uncertain experimental data

Have received little attention until now

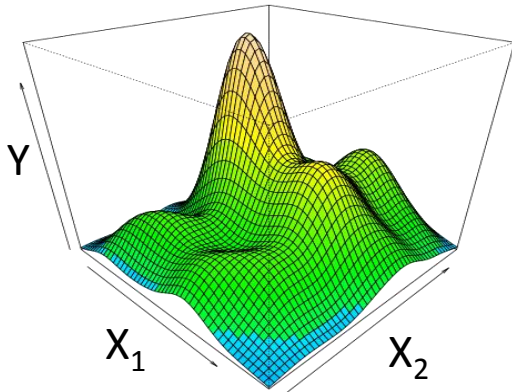
Brief survey of the litterature

- 2008: Cooke & Schmidler, *Biophysical Journal*, 95: 4497–4511
Calibration of dielectric constant in peptides to better reproduce helical folding of peptides
- 2011: Cailliez & Pernot, *J. Chem. Phys.* 134, 054124
Bayesian calibration of a LJ forcefield for Argon and UP in MD simulations
- 2012: Angelikopoulos *et al.*, *J. Chem. Phys.* 137, 144103 (2012)
Bayesian calibration of a LJ forcefield for Argon using GP surrogate models
- 2012: Rizzi *et al.*, *Multiscale model. Simul.* 10(4):1428–1492
Estimation of forcefield parameters of a water model and UP in MD simulations using Polynomial Chaos
- 2013: Rizzi *et al.*, *J. Chem. Phys.* 138, 194105
Statistical calibration of LJ parameters of monoatomic ions
- 2014: Cailliez *et al.*, *J. Comp. Chem.*, 35, 130–149
Estimation of forcefield parameters of a water model and UP in MD simulations using GP surrogate models
- 2016: Wu *et al.*, *Phil. Trans. R. Soc. A* 374: 20150032
Hierarchical modeling to calibrate LJ parameters among heterogeneous data
- 2016 : Pernot & Cailliez, arXiv:1611.04376
Review of statistical calibration/prediction models handling data inconsistency and model inadequacy

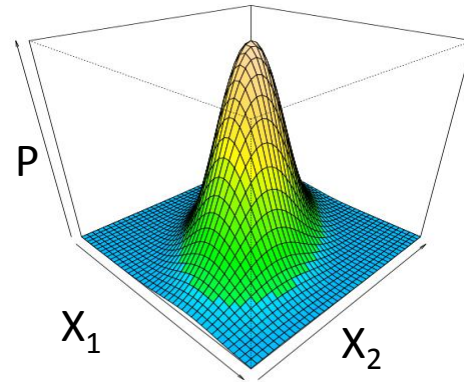
- Bayesian calibration framework
- Statistical calibration of a forcefield for Argon:
 - Comparing numerical and parametric uncertainties
 - Transferability between properties
- Calibration of a water forcefield:
 - How to decrease the computational burden?
 - Use of surrogate models and Efficient Global Optimization strategies
- How to deal with model inadequacy?
- Conclusions

Statistical calibration and uncertainty propagation

Calibration data:
 $Y = Y(X_1, X_2)$



PDF:
 $P(\Omega|\{Y_{i,exp}\})$



$\theta = \{X_1, X_2\}$



Bayesian calibration

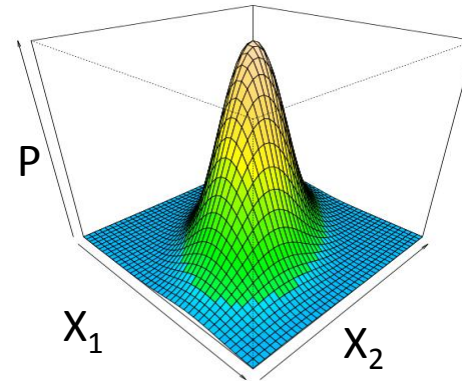
$$P(\theta|\{Y_{i,exp}\}) \propto \underbrace{P(\{Y_{i,exp}\}|\theta)}_{\text{Likelihood}} \times \underbrace{P(\theta)}_{\text{Prior}}$$

Independent measurements and Gaussian hypothesis for residuals:

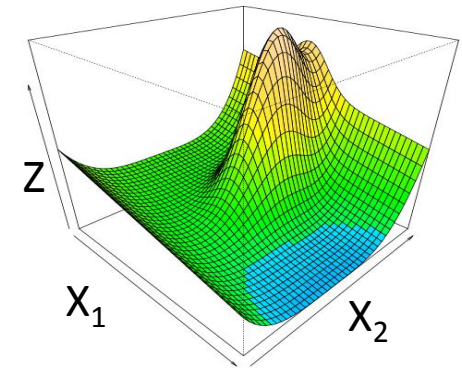
$$P(\{Y_{i,exp}\}|\theta) = \prod_i \left(\frac{1}{u_i \sqrt{2\pi}} \exp\left(-\frac{(Y_i - Y_{i,exp})^2}{2u_i^2}\right) \right)$$

$$u_i^2 = u_{i,exp}^2 + u_{i,mod}^2$$

PDF:
 $P(\theta|\{Y_{i,iexp}\})$



Results of a simulation:
 $Z = Z(X_1, X_2)$



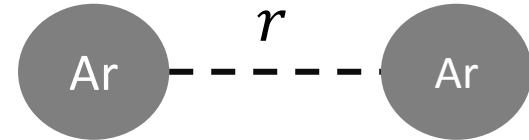
Uncertainty propagation

Monte Carlo sampling of the PDF and estimation of the model:

$$\bar{Z} = \frac{1}{N} \sum_{i=1}^N Z(\theta_i) \quad u_Z^2 = \text{var}(Z_i)$$

A simple test-case: forcefield for Argon

- Two-parameters Lennard-Jones forcefield for Argon:



- Statistical calibration:

- *Uniform prior* : $P(\sigma, \varepsilon) = \text{Cte}$
- Experimental data for calibration:

2nd virial coefficient B from 150 to 450K

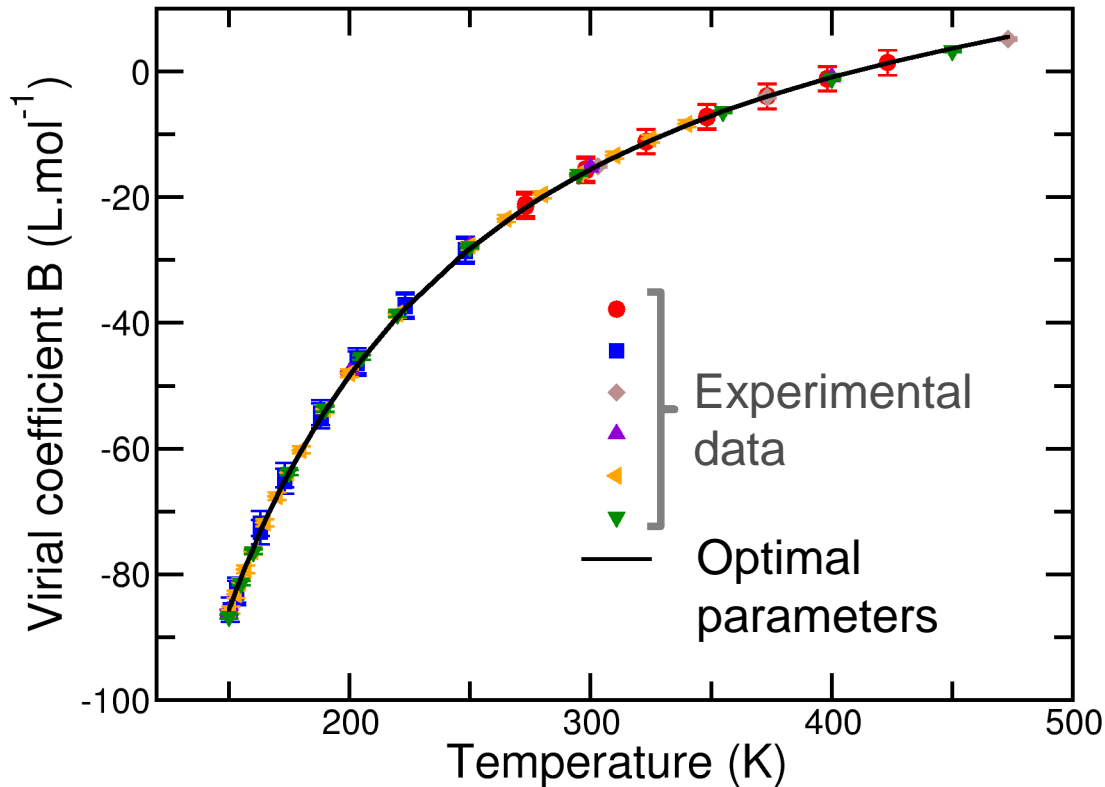
(lit.: $\sigma = 3.405\text{\AA}$; $\varepsilon = 119\text{K}$)

$$E(r) = 4\varepsilon \left[\left(\frac{\sigma}{r}\right)^{12} - \left(\frac{\sigma}{r}\right)^6 \right]$$

- Specificities of this calibration:

- Only two parameters
- Analytical expression linking B to the parameters σ and ε
 → Analytical PDFs and $u_{i,mod} = 0$

Calibration on 2nd virial coefficient



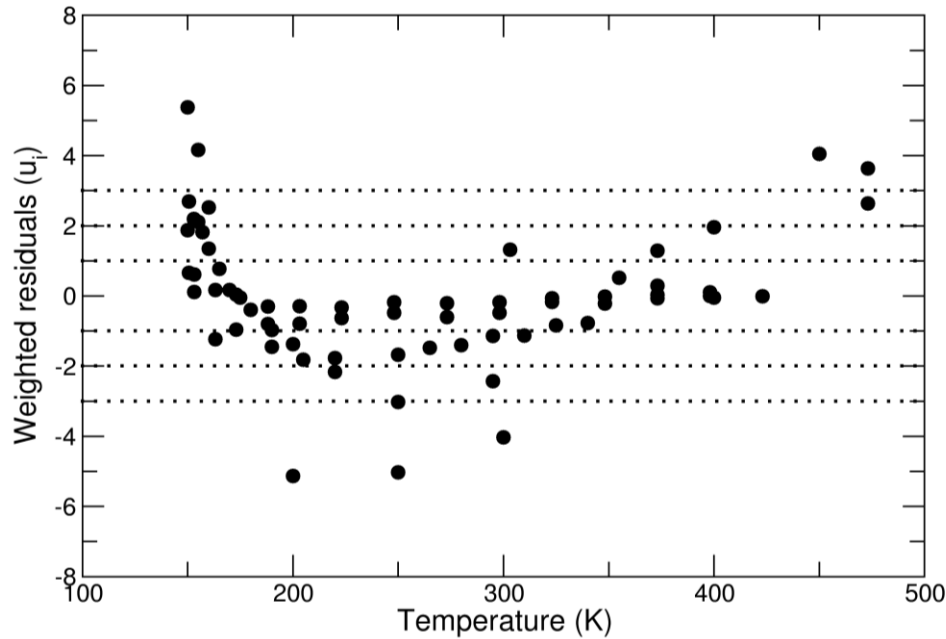
$$\sigma = 3.422 \pm 0.002 \text{ \AA}$$

$$\varepsilon = 119.61 \pm 0.09 \text{ K}$$

(lit.: $\sigma = 3.405 \text{ \AA}$; $\varepsilon = 119 \text{ K}$)

- Various sources of measurements
- At first sight: successfull calibration and small uncertainties

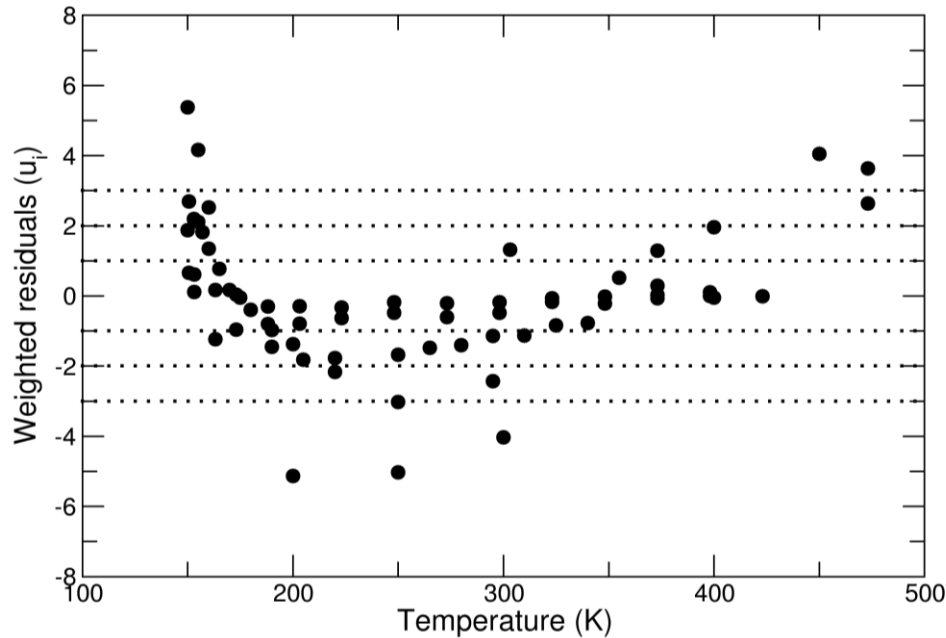
Calibration on 2nd virial coefficient



$$w_i = \frac{B(\bar{\sigma}, \bar{\varepsilon}, T_i) - B_{i,exp}}{u_i}$$

- Gaussian hypothesis for the residues violated
- Possible origins of the problem:
 - Inadequacy of the model
 - Inconsistency between some experimental data
 - Underestimated experimental uncertainties

Calibration on 2nd virial coefficient

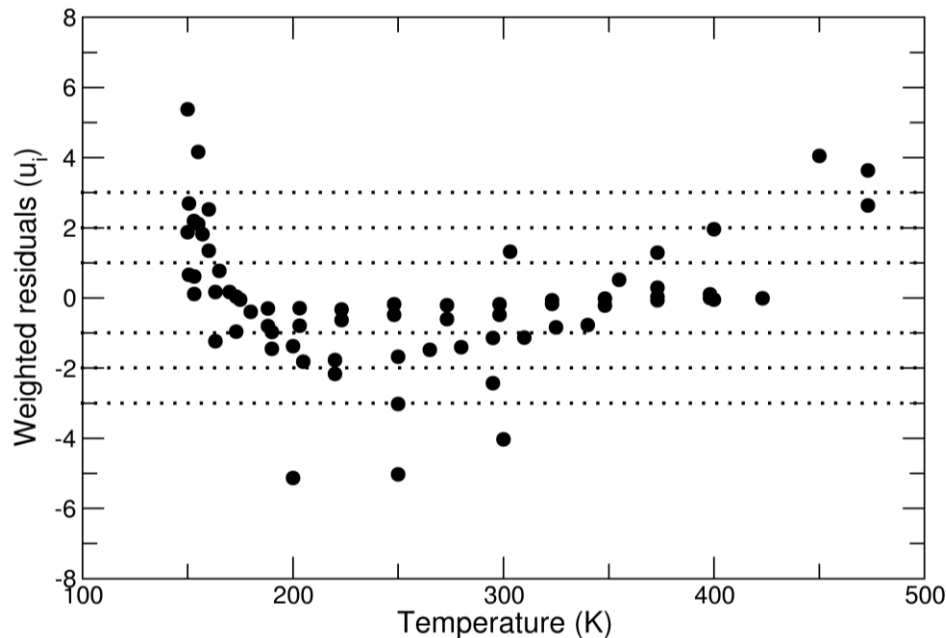


$$w_i = \frac{B(\bar{\sigma}, \bar{\varepsilon}, T_i) - B_{i,exp}}{u_i}$$

- Gaussian hypothesis for the residues violated
- Possible origins of the problem:
 - **Inadequacy of the model** $u_i'^2 = u_{i,exp}^2 + u_{i,m}^2$
 - Inconsistency between some experimental data
 - Underestimated experimental uncertainties

Additive uncertainties on model predictions
 $u_m = \mathcal{N}(0, s^2)$
 → untransferable for prediction to other types of data

Calibration on 2nd virial coefficient



$$w_i = \frac{B(\bar{\sigma}, \bar{\varepsilon}, T_i) - B_{i,exp}}{u_i}$$

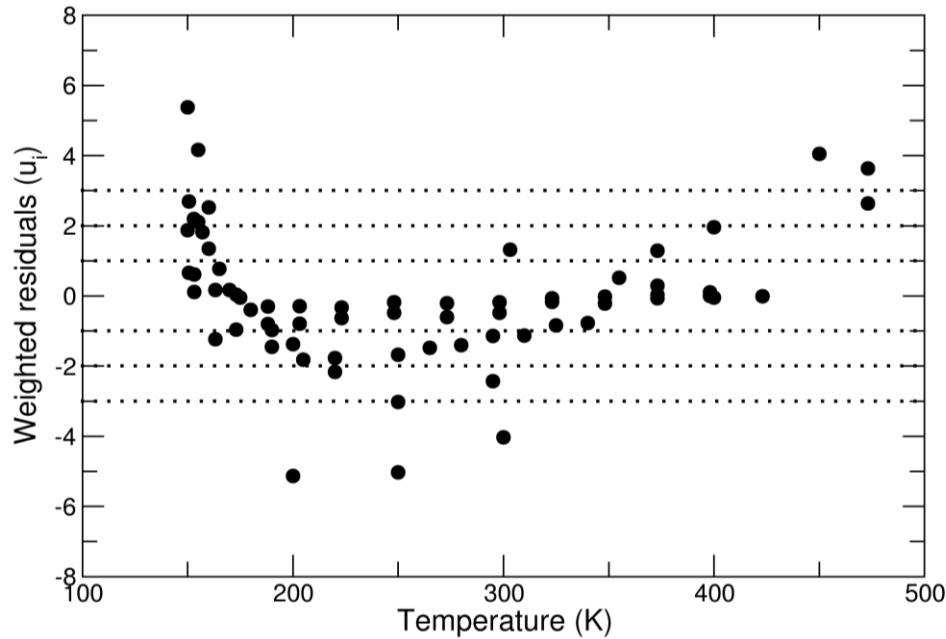
Fixed Laboratory Effects Model (FLEM):

- Each experimental set bears an unknown but constant bias.
- Modification of the data by adding a term (to be calibrated) that depends on the experimental data set.

- Gaussian hypothesis for the residues violated
- Possible origins of the problem:
 - Inadequacy of the model
 - Inconsistency between some experimental data
 - Underestimated experimental uncertainties

$$B'_{i,exp}(\sigma, \varepsilon) = B_{i,exp}(\sigma, \varepsilon) + \lambda_S(i)$$

Calibration on 2nd virial coefficient



$$w_i = \frac{B(\bar{\sigma}, \bar{\varepsilon}, T_i) - B_{i,exp}}{u_i}$$

Random Laboratory Effects Model (RLEM):

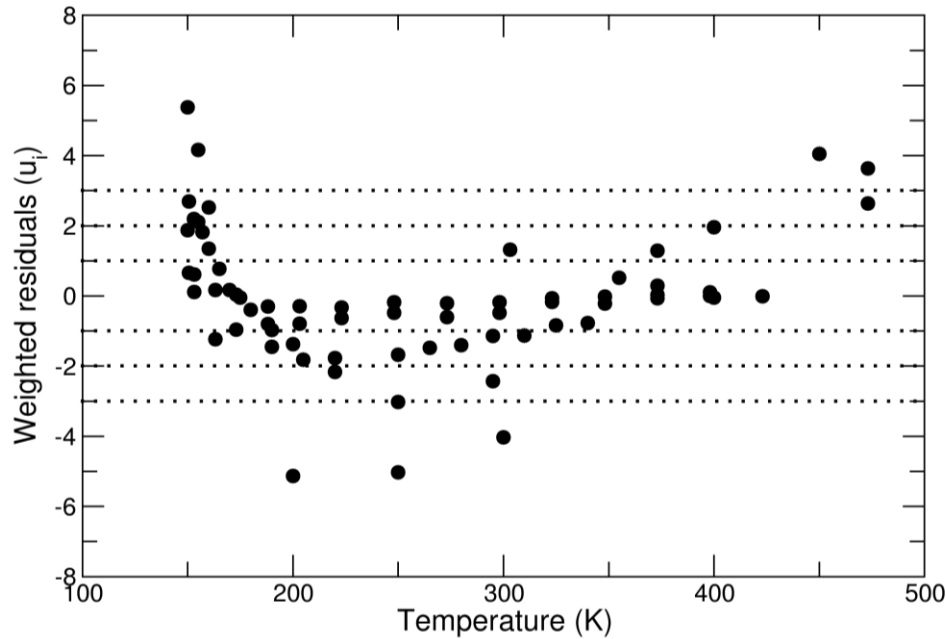
Each experimental data bears a supplementary unknown bias

$$u_{add} = \mathcal{N}(0, s^2)$$

- Gaussian hypothesis for the residues violated
- Possible origins of the problem:
 - Inadequacy of the model
 - **Inconsistency between some experimental data**
 - Underestimated experimental uncertainties

$$u_i'^2 = u_{i,exp}^2 + u_{i,add}^2$$

Calibration on 2nd virial coefficient



$$w_i = \frac{B(\bar{\sigma}, \bar{\varepsilon}, T_i) - B_{i,exp}}{u_{i,exp}}$$

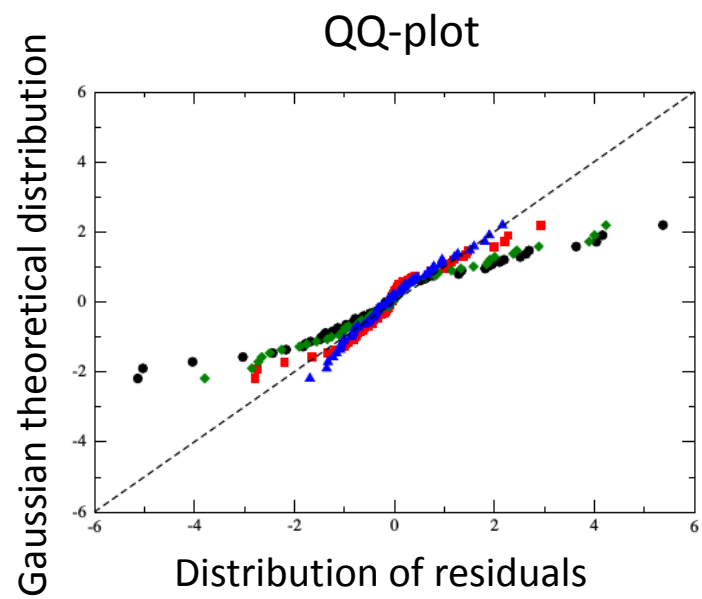
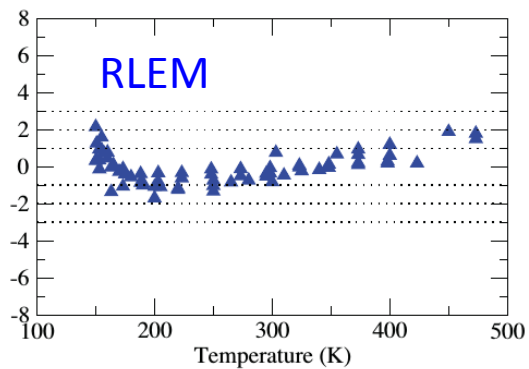
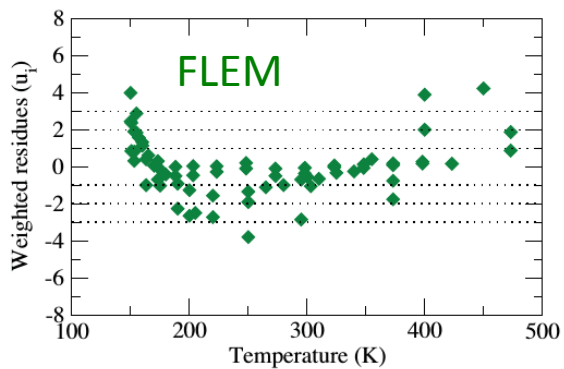
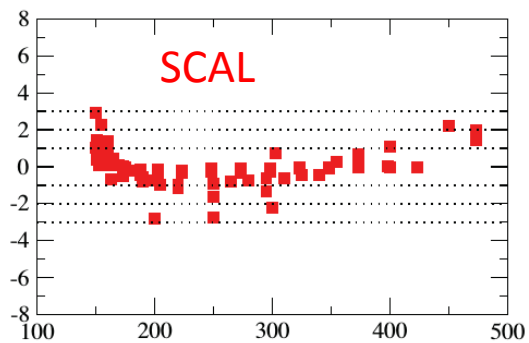
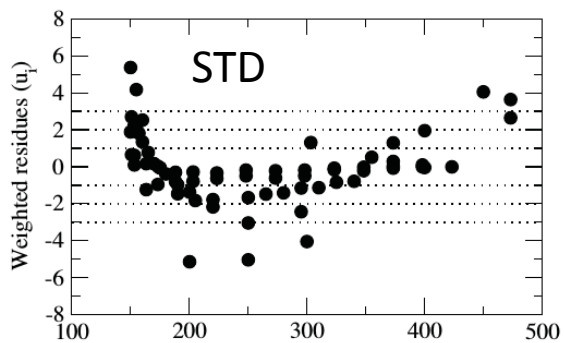
SCAL:

Scaling of experimental uncertainties by an a priori unknown factor

- Gaussian hypothesis for the residues violated
- Possible origins of the problem:
 - Inadequacy of the model
 - Inconsistency between some experimental data

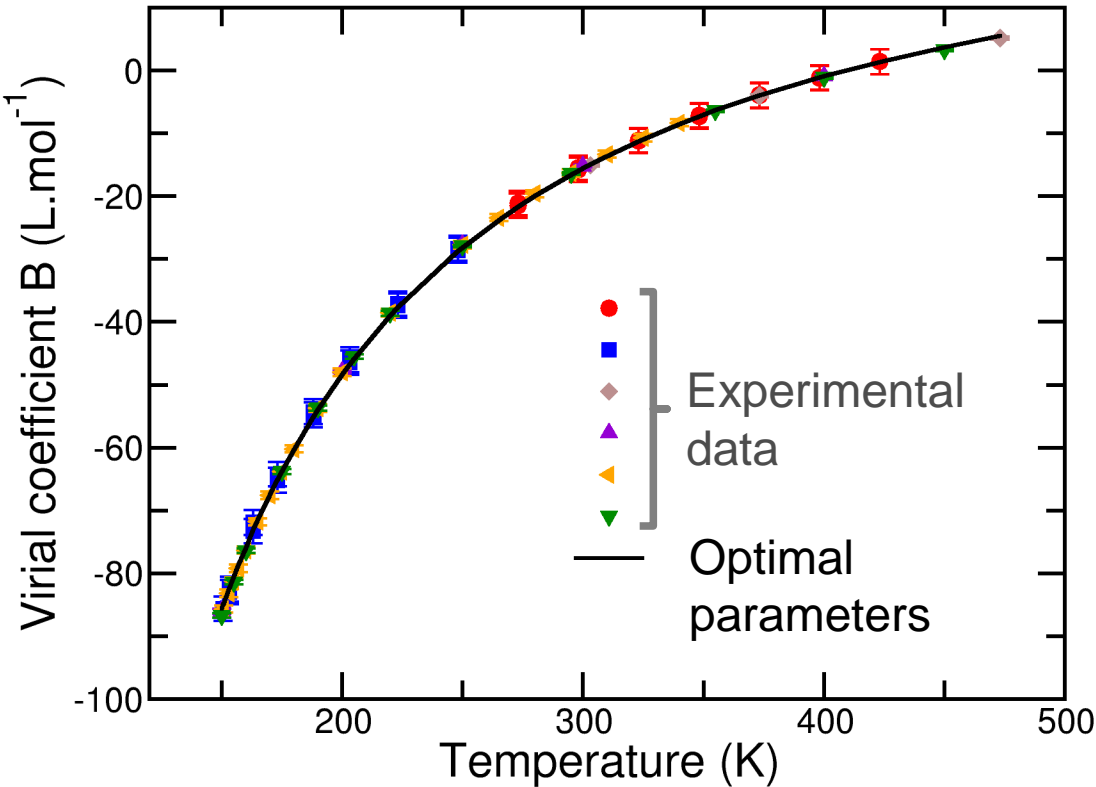
■ **Underestimated experimental uncertainties:** $u'_i = s \times u_{i,exp}$

Various calibration schemes



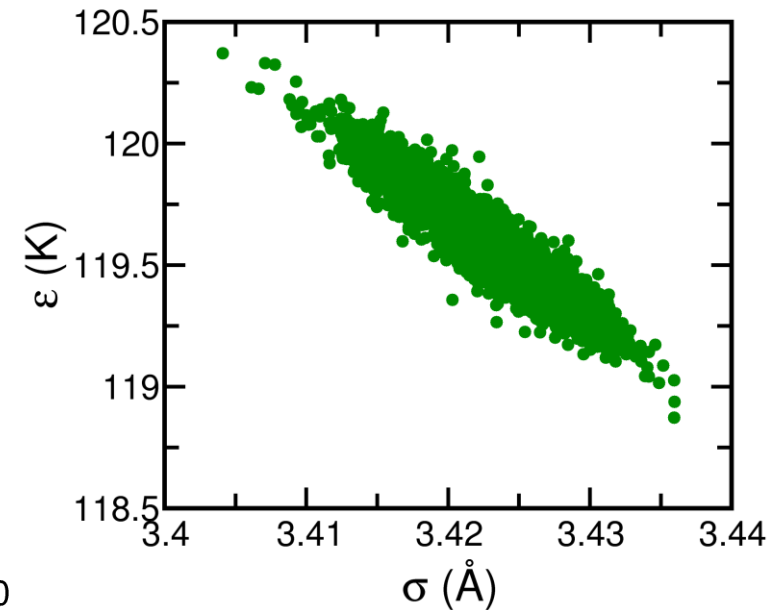
- SCAL and RLEM models achieve statistical consistency in the residuals
- Long-range correlation remains at different temperatures (limitation of the LJ interaction model)

Results of the SCAL calibration



- Small scaling factor for uncertainties (<2)
- Small uncertainties on the parameters
- Strong correlation between σ and ε

Markov chain over the PDF

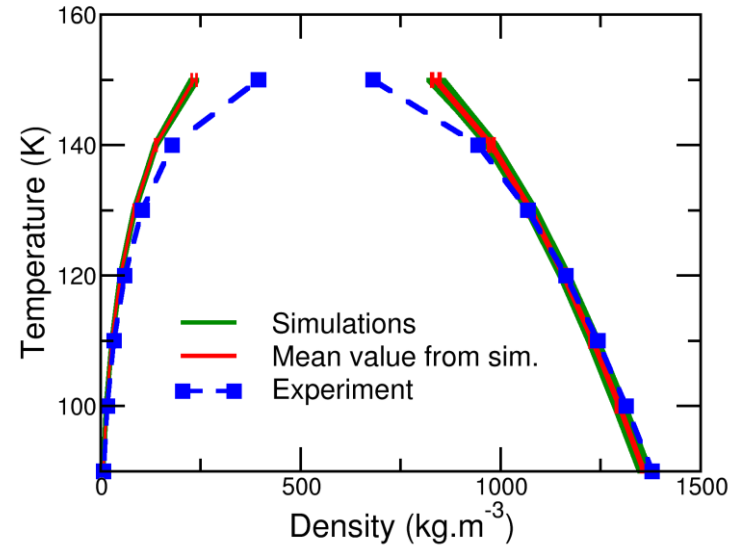
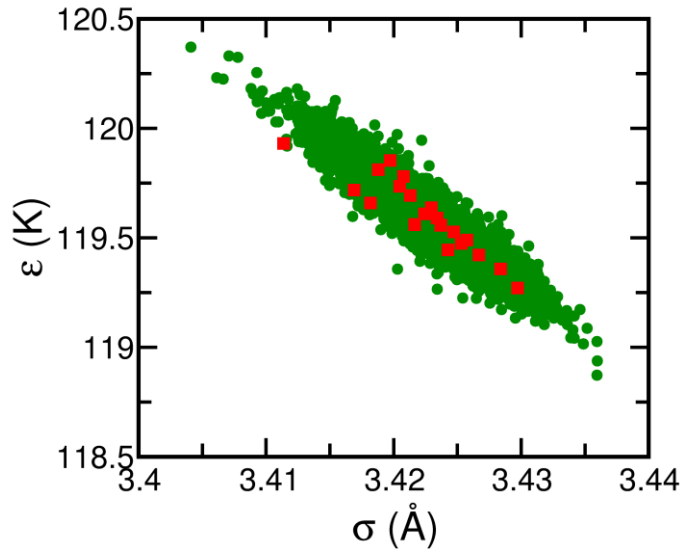


$$\sigma = 3.422 \pm 0.004 \text{ \AA}$$

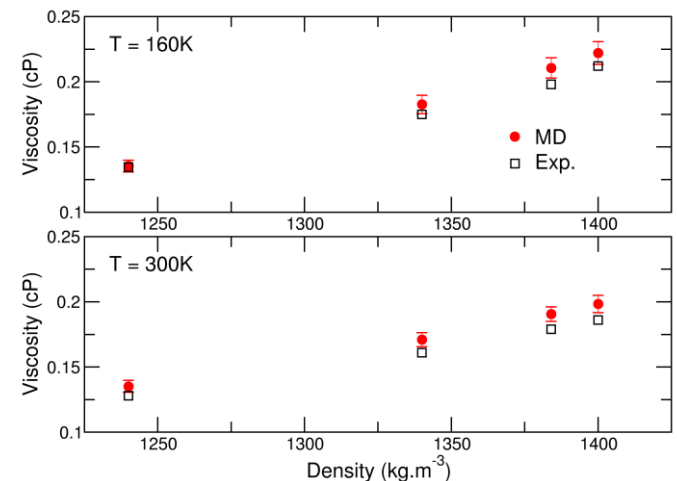
$$\varepsilon = 119.61 \pm 0.18 \text{ K}$$

Values from literature:
 $\sigma = 3.405 \text{ \AA}$, $\varepsilon = 119.8 \text{ K}$

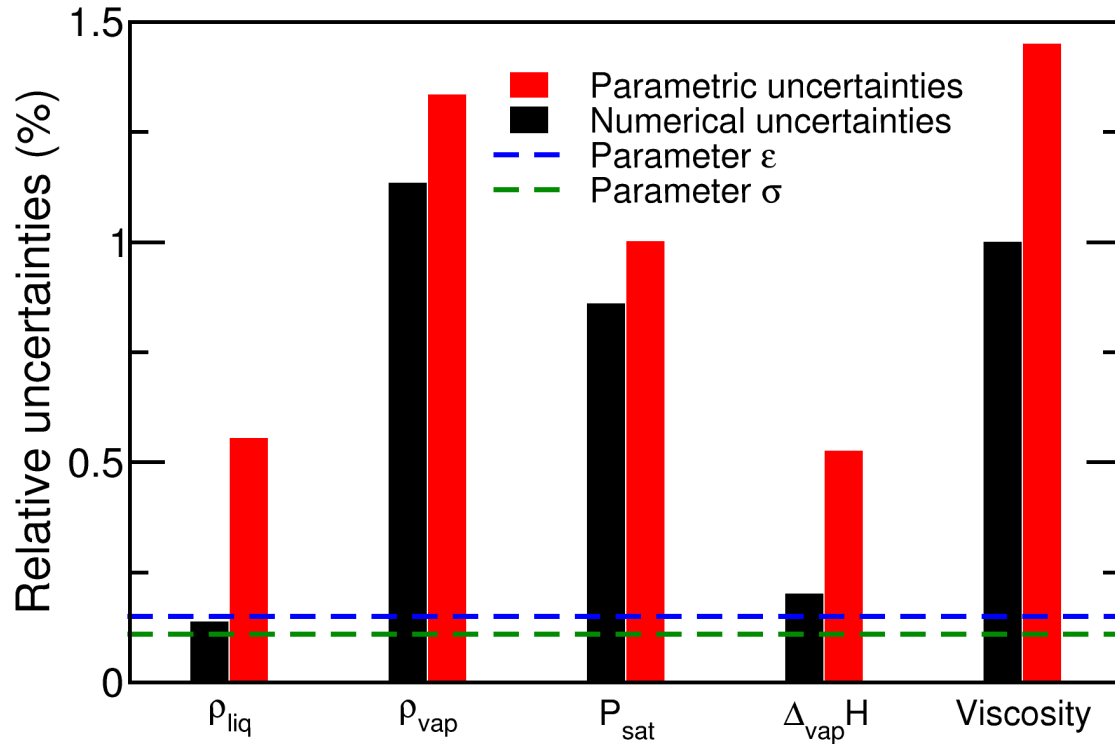
Uncertainty propagation in molecular simulation



- Sample (LHS) of the PDF of (σ, ε)
- Computation of L/V phase diagrams and liquid viscosities
- Parametric uncertainties remain small and do not allow to reconcile computed and experimental values



Numerical/parametric uncertainties



- Parameters uncertainties amplified by molecular simulation
- **Parametric uncertainties** bigger than numerical uncertainties

A first conclusion

- What we have learnt from this LJ system:
 - An operative methodology for statistical calibration and uncertainty propagation
 - Parametric uncertainties small but greater than numerical uncertainties
 - Taking into account parametric uncertainty is not sufficient to have quantitative transferability to other properties

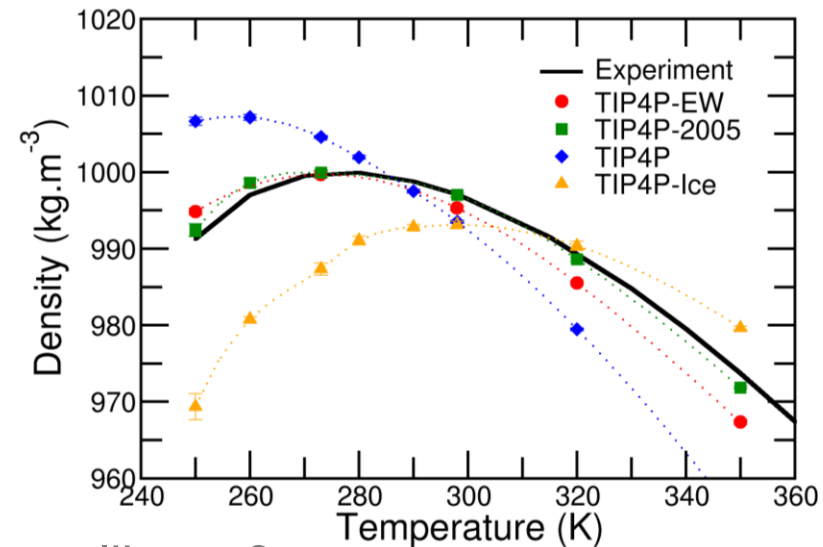
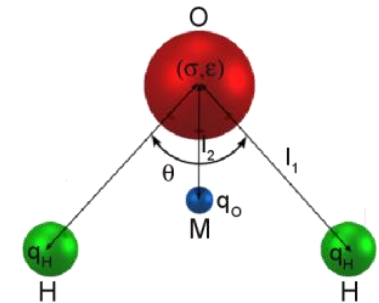
Cailliez et Pernot, *J. Chem. Phys.* **134**, 054124 (2011)

- What to do next:
 - Increase the complexity of the forcefield?
 - Is the method tractable when calibration data requires molecular simulation to be evaluated?
 - How to deal with model inadequacy

TIP4P water forcefield: strategy of calibration

- TIP4P water forcefield:

- Various parameter sets available in the literature
- 4 parameters: $\sigma, \epsilon, q_H, l_2$
- Calibration data: liquid water density at 4 temperatures from 253K to 350K
- Molecular simulations needed to compute the calibration data



- Strategy of calibration:

- Reducing the number of parameters to calibrate?
Global sensitivity analysis
- Avoiding the use of expensive molecular simulations:
Use of surrogate models

Surrogate modeling

- Objective: Replace the costly computer model \mathcal{M} (molecular simulation) by a cheap estimator $\tilde{\mathcal{M}}$ (surrogate model)
- Use of a response function $f: Y_i(\theta_i) = f(\theta_i) + R_i$ } Error made by the response function
- Gaussian Process (GP) or Kriging surrogate models:
 - Stochastic function Z to describe the $\{R_i\}$
 - if two points θ_a and θ_b are close in the parameter space, R_a and R_b should be close too

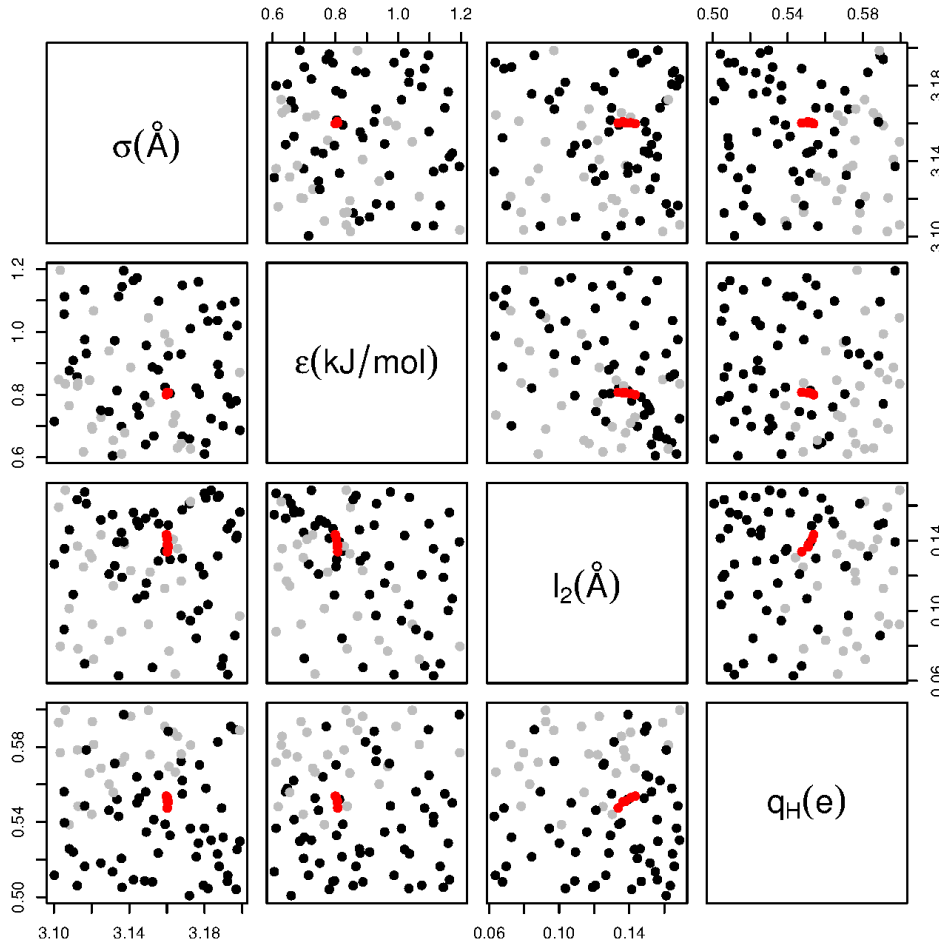
$$\tilde{\mathcal{M}}(\theta) = \underbrace{f(\theta)}_{\text{OK: constant function}} + Z(\theta)$$

$$E[Z(\theta)] = 0$$

$$\text{cov}(Z(\theta_1), Z(\theta_2)) = \sigma^2 R(\theta_1, \theta_2)$$

$$R(\theta_1, \theta_2) = R(\theta_1 - \theta_2) = \exp\left(-\sum_{i=1}^p \frac{(\theta_{1i} - \theta_{2i})^2}{2l_i^2}\right)$$

Building of the surrogate models



- Initial sampling of the parameter space:
 - Maximin LHS
 - 84 parameter sets (D1)

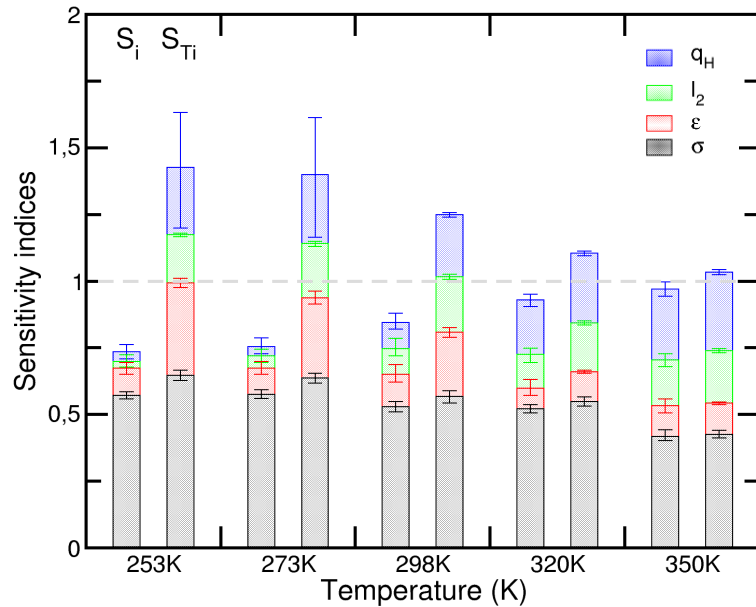
- One surrogate model for each property

- Leave-one-out predictivity coefficients:

$Q_2 \geq 92\%$

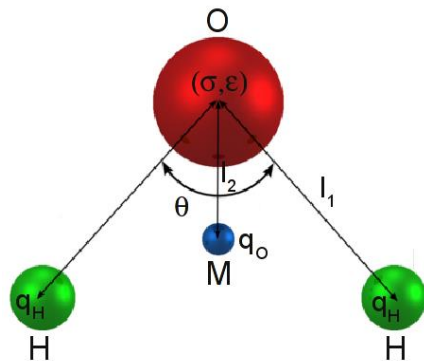
- Some parameter sets lead to badly converged simulations:
 - Subsample of 57 parameter sets (D2)

Global sensitivity analysis on density



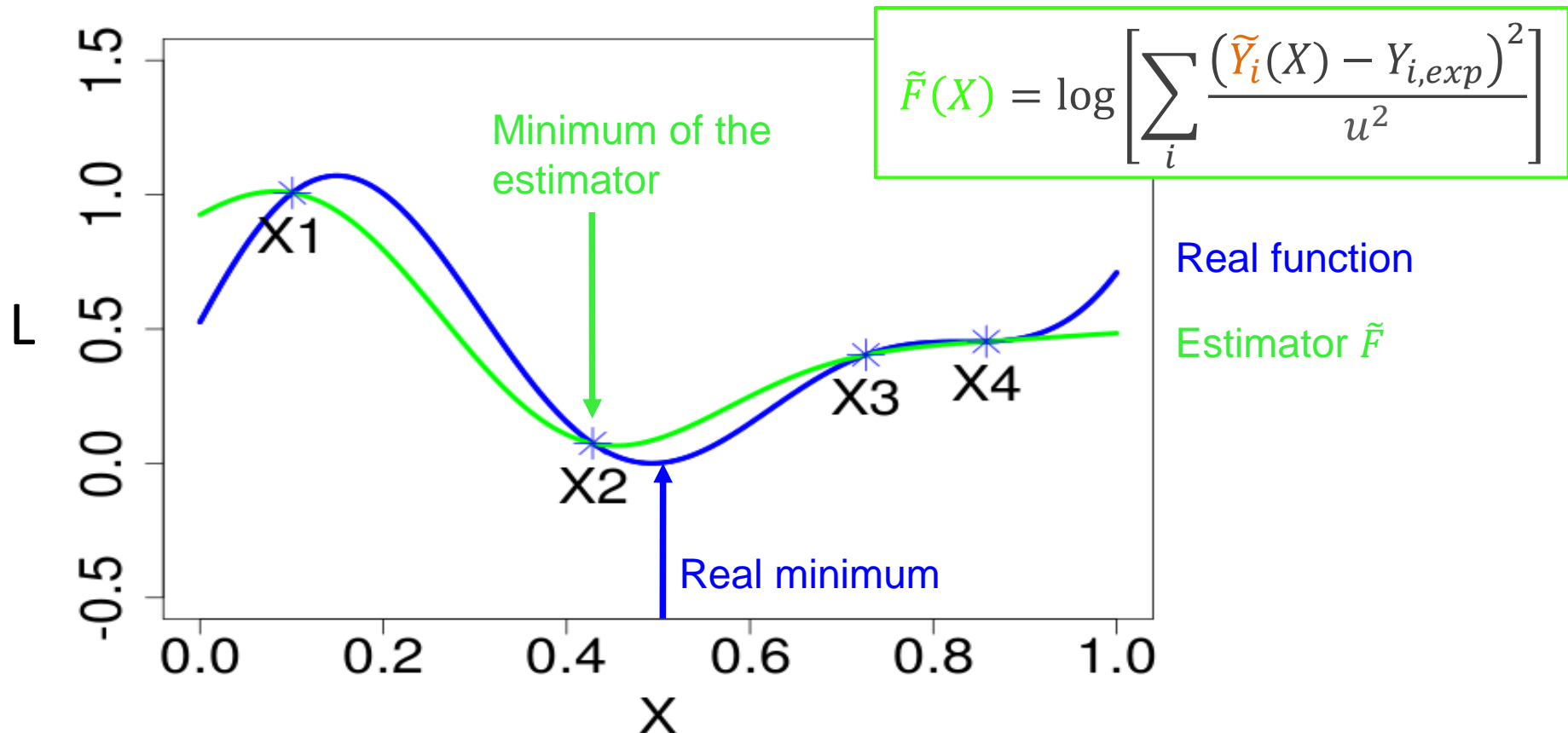
- Analysis of Sobol sensitivity indices:
 - σ and q_H are the more important parameters
 - All parameters have non-negligible effect
 - Interactions between parameters

- Consequences:
 - No reduction of the dimension of the parameter space possible
 - One-at-a-time calibration inefficient



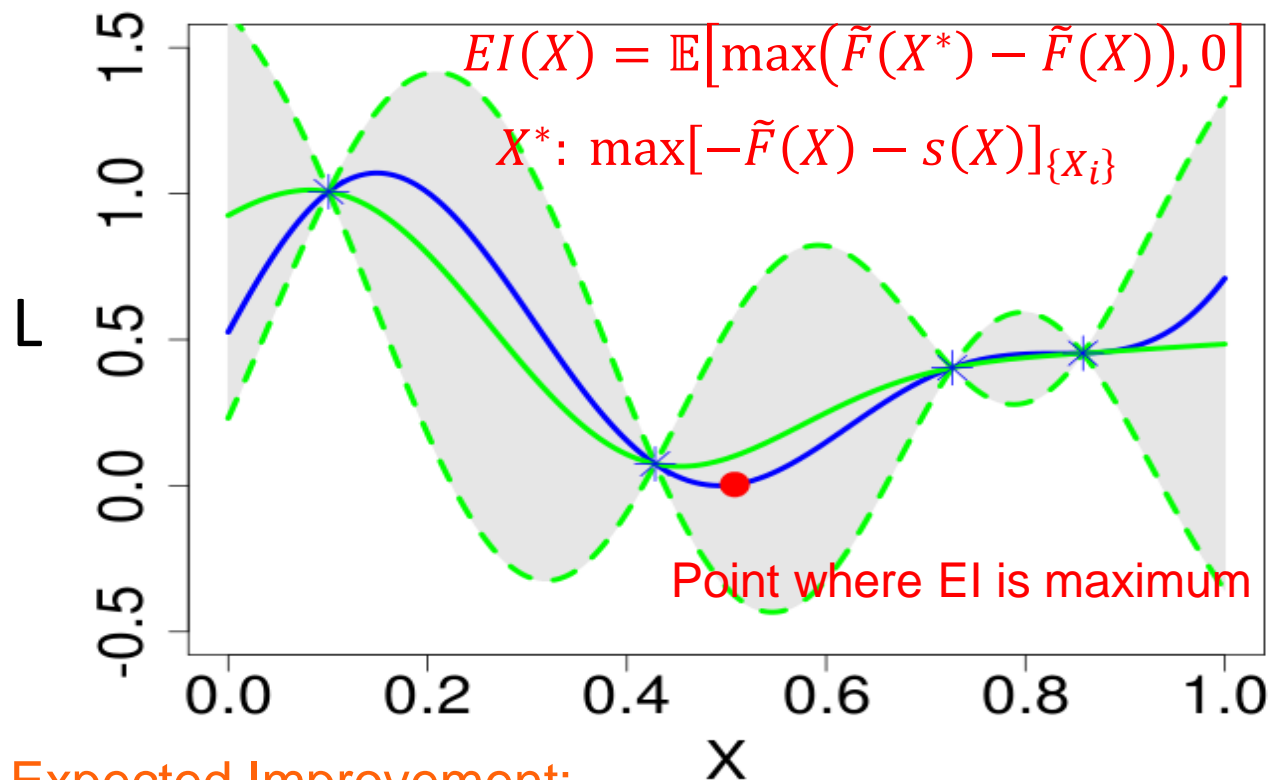
Saltelli et al., 2010, *Comp Phys Comm*, **181**: 259–270

Calibration with GP surrogate models



- The minimum region of the estimator does not necessarily reproduce accurately the real minimum
- Iterative improvement of the estimator of the PDF:
Use of “Efficient Global Optimization” (EGO) algorithms

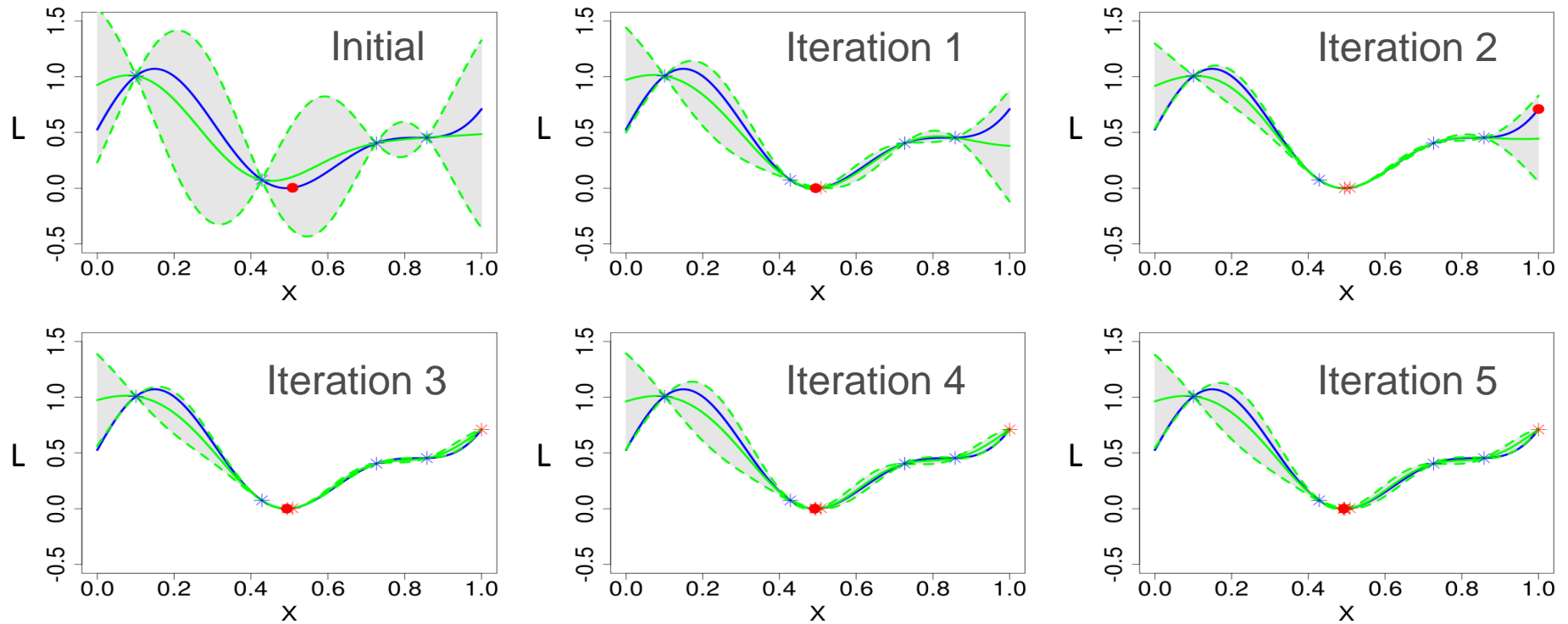
Efficient Global Optimization



Huang et al., 2006,
J. Glob. Opt., **34**: 441

- **Expected Improvement:**
 use of the uncertainty prediction $s(X)$ of the surrogate model \tilde{F}
- Two technical difficulties:
 - GP optimised on stochastic data
 - \tilde{F} not a GP: EI computed numerically

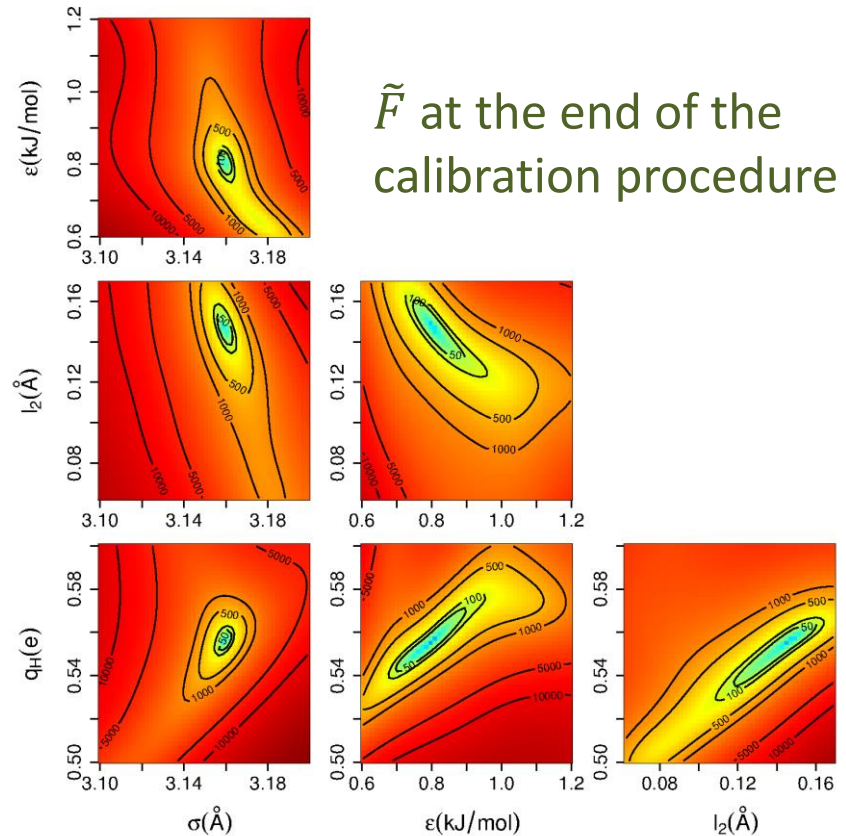
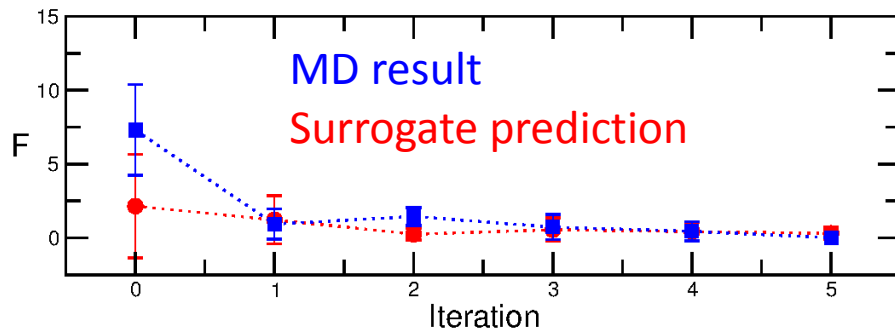
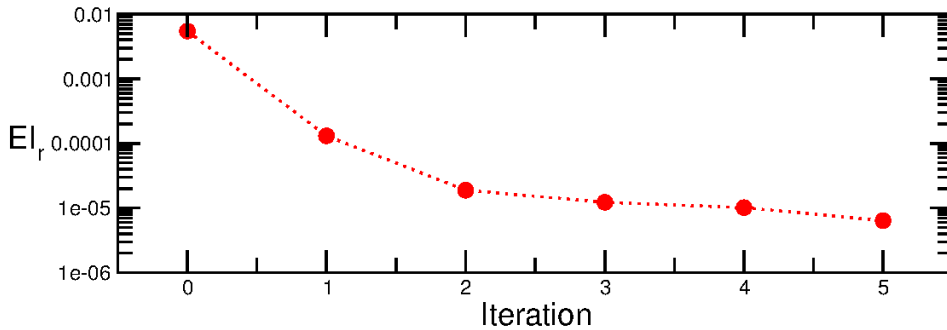
Efficient Global Optimization



- Iterative procedure
- End of the procedure:

$$EI_r = \frac{EI}{\max(\tilde{F}) - \min(\tilde{F})} < 10^{-5}$$

EGO convergence



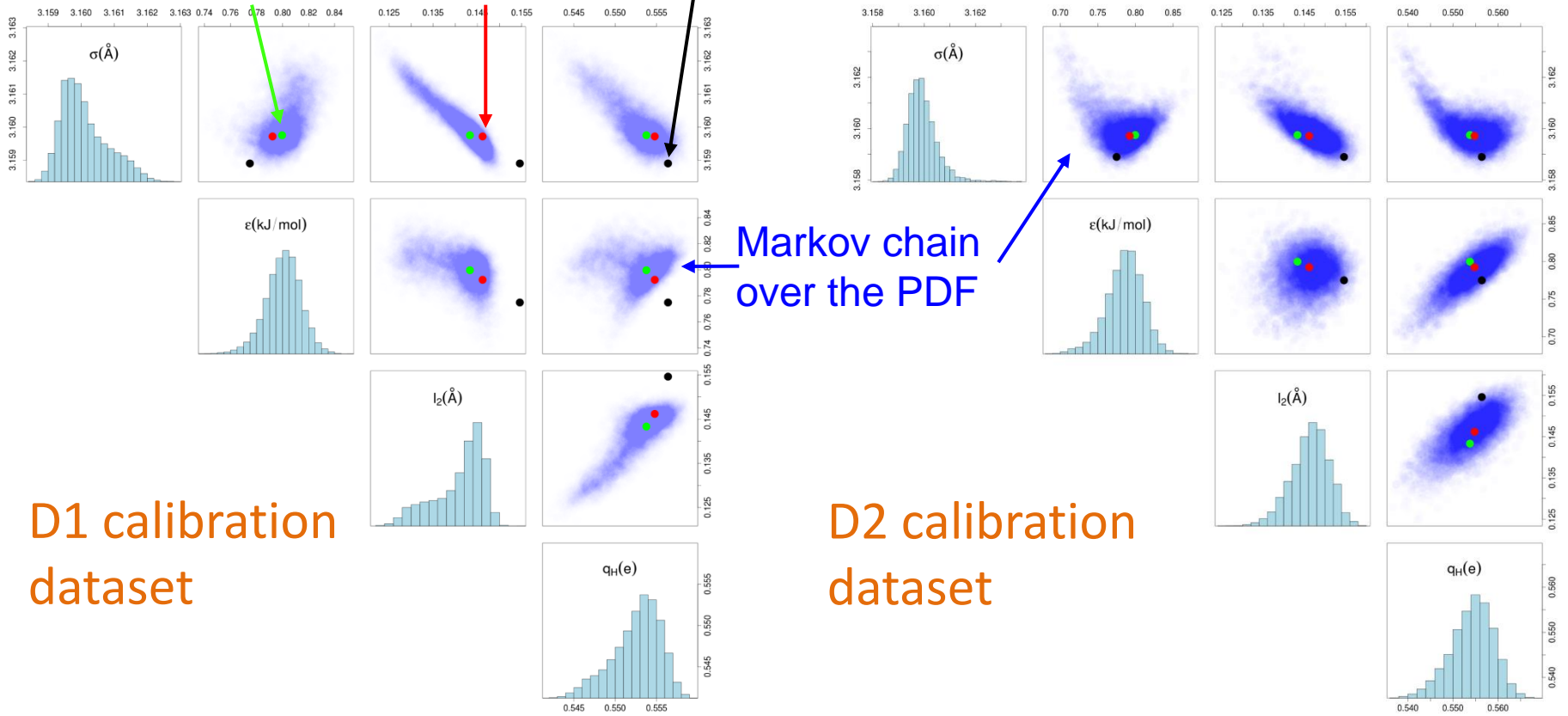
- Rapid convergence of the EGO
- Rapid improvement of the prediction around the optimum
- Single well defined optimum of the score function

Results of the calibration

Opt. EGO-D2

Opt. EGO-D1

TIP4P-2005

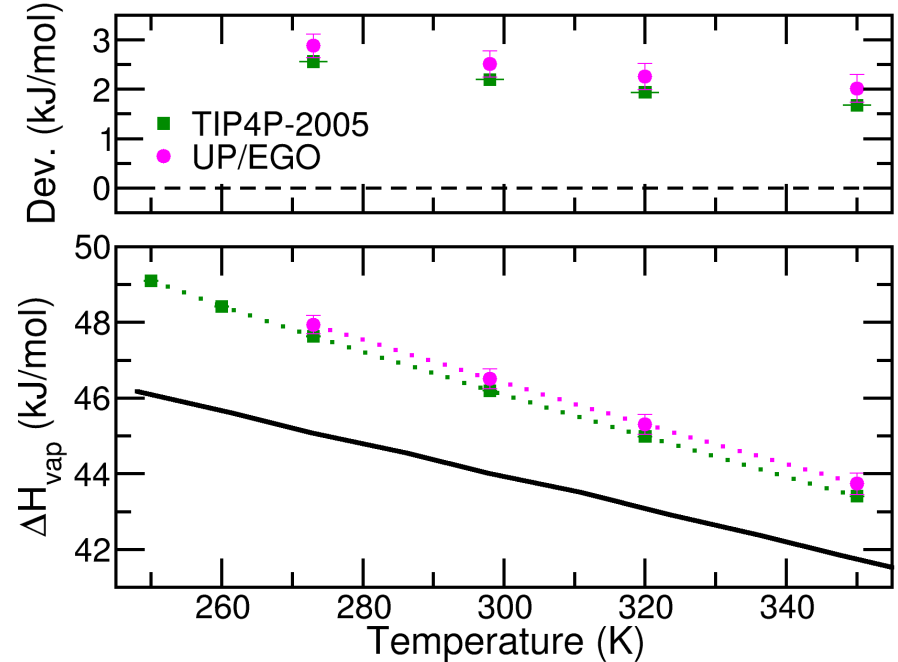
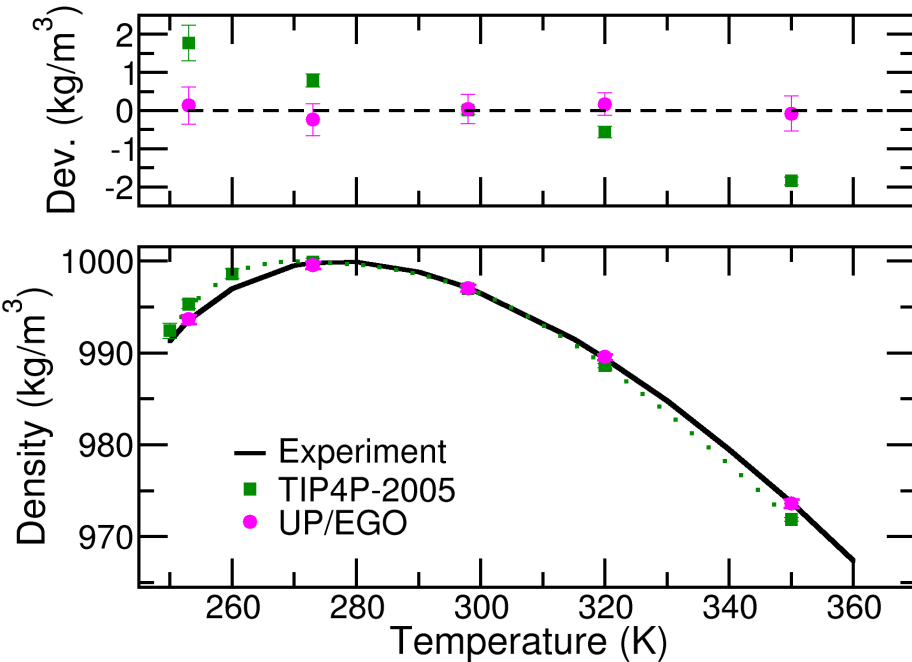


D1 calibration dataset

D2 calibration dataset

- D1 calibration dataset: ≈ 90 parameter sets used
- D1 and D2 calibration: similar results – low sensitivity to badly converged simulations

Parametric uncertainties - TIP4P forcefield



- Uncertainty propagation using kriging surrogate models for density and vaporization enthalpy
- Parametric uncertainties bigger than numerical uncertainties

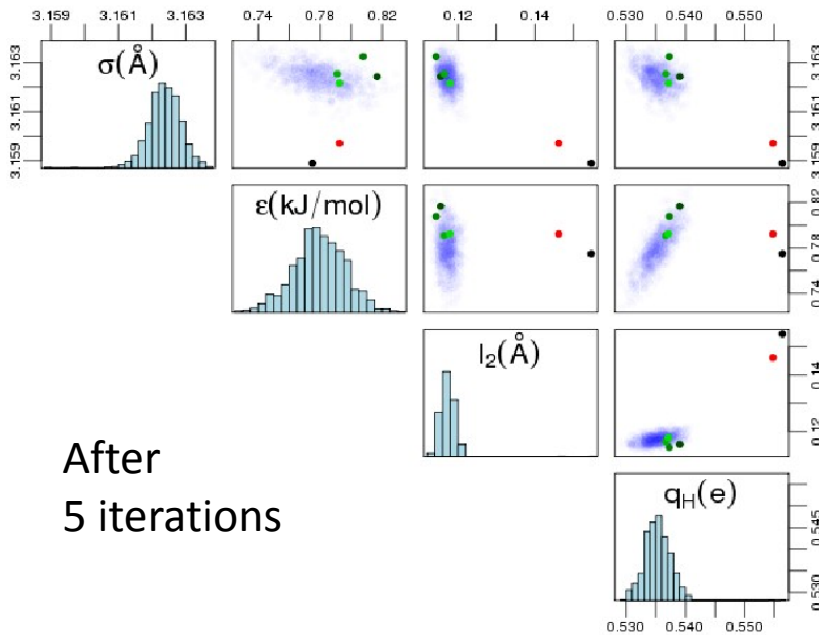
A second conclusion

- Similar conclusions as for the argon case regarding parametric uncertainties:
 - At least as big as numerical uncertainties
 - Taking into account parametric uncertainty is not sufficient to have quantitative transferability to other properties
- Use of surrogate models:
 - Extensive exploration of parameter space at lower cost
 - Global sensitivity analysis (reduce parameter space dimension)
 - Global optimisation of the parameters possible
 - Cailliez, Bourasseau, and Pernot, *J. Comp. Chem.* **35**: 130-149 (2014)
- Limitations and unresolved issues:
 - Reducing the cost of the optimization procedure
 - How to deal with model inadequacy?

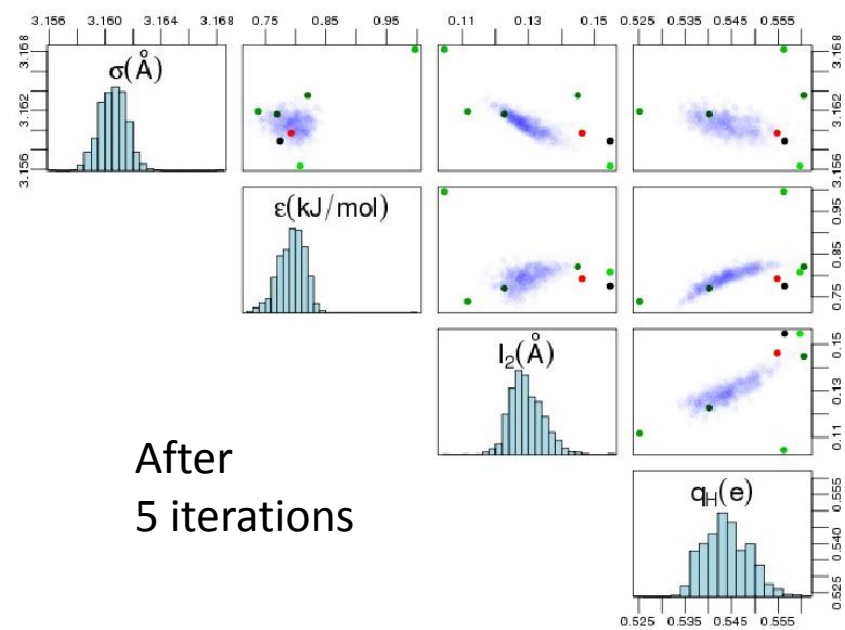
Reduction of the initial sample design

$$EI(X) = \mathbb{E}[\max(\tilde{F}(X^*) - \tilde{F}(X), 0)]$$

$$AEI(X) = EI(X) \left(1 - \frac{\tau^2}{\sqrt{s^2(X) + \tau^2}} \right)$$



After
5 iterations

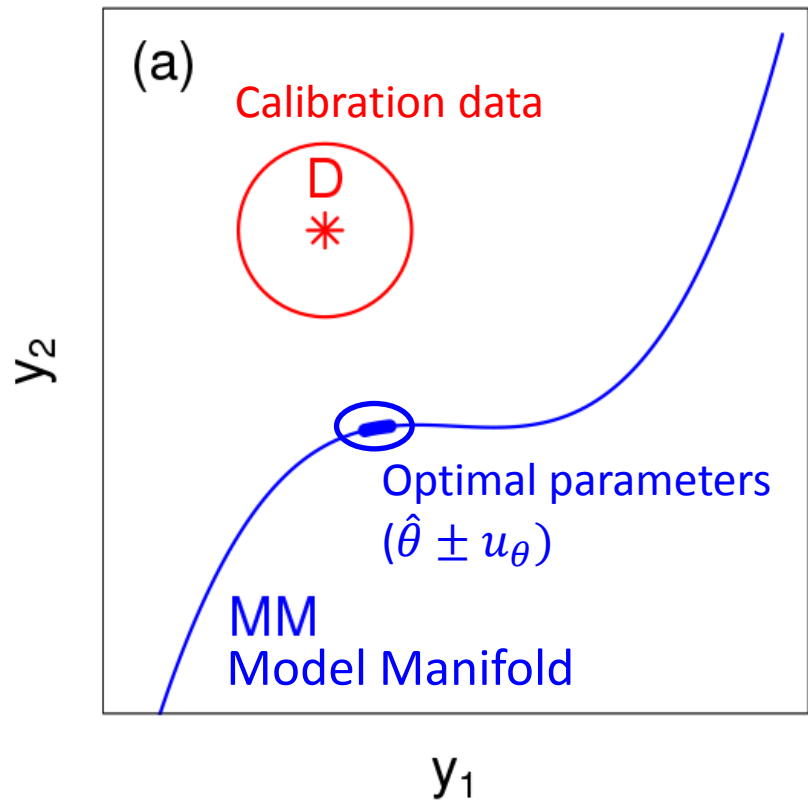


After
5 iterations

- Sparse 12-points initial sample
- Predictivity coefficient of the surrogate models still satisfying
- EGO with EI gets stucked in a wrong region: need to increase exploratory behaviour → AEI

Huang *et al.*, 2006, *J. Glob. Opt.*, **34**: 441

The issue of model inadequacy



(Deterministic) model

$$y_i = \mathcal{M}(x_i, \theta) + R_i + e_i$$

Calibration data

Residual

Experimental uncertainty

- Inadequacy remains at the calibration stage
- Prediction inefficient even taking into account parameters uncertainties

Solving model inadequacy on synthetic data

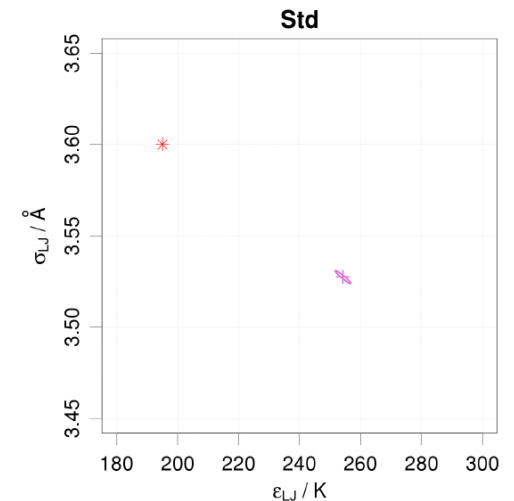
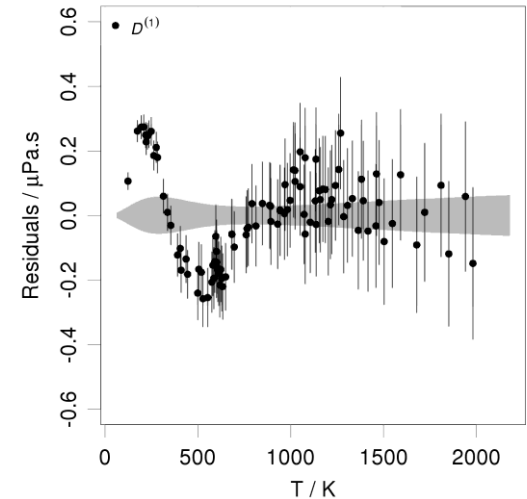
- Krypton described by a Lennard-Jones potential:
 $\theta = \{\sigma = 3.6\text{\AA}; \varepsilon = 195\text{K}\}$
- Gas-phase viscosity: Chapman-Enskog model:

$$\eta = \mathcal{M}(T, \sigma, \varepsilon) = 2.6693 \frac{\sqrt{mT}}{\sigma^2 \Omega} \quad T^* = T/\varepsilon$$

$$\Omega = \frac{A}{(T^*)^B} + \frac{C}{\exp(DT^*)} + \frac{E}{\exp(FT^*)}$$

- Synthetic data:
 - 100 data points for various T , generated with a modified value of C in CE formula
 - Generation of synthetic « experimental » uncertainties

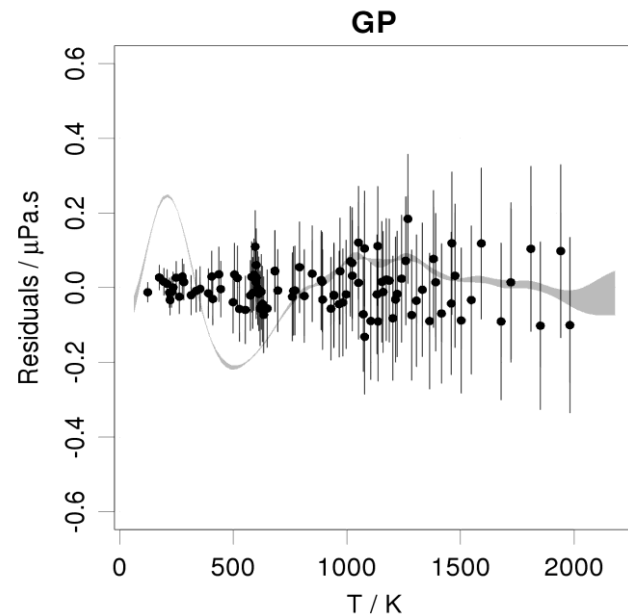
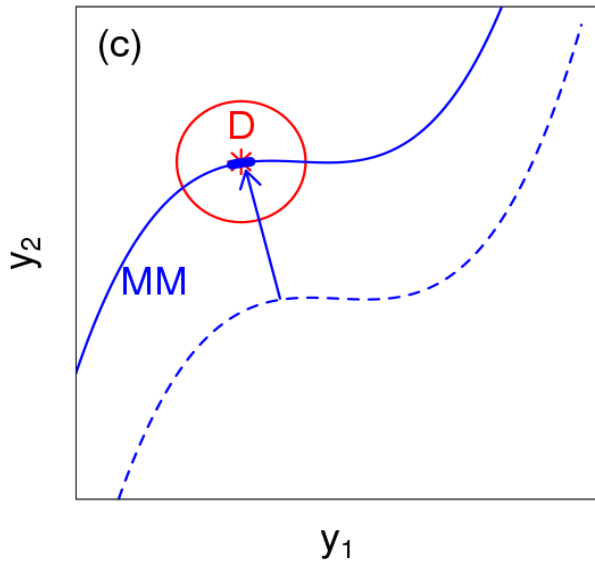
Standard calibration



Correcting the model (GP)

- Add a discrepancy term to correct model errors:

$$y_i = \mathcal{M}(x_i, \theta) + GP(x_i, \theta_K) + e_i$$



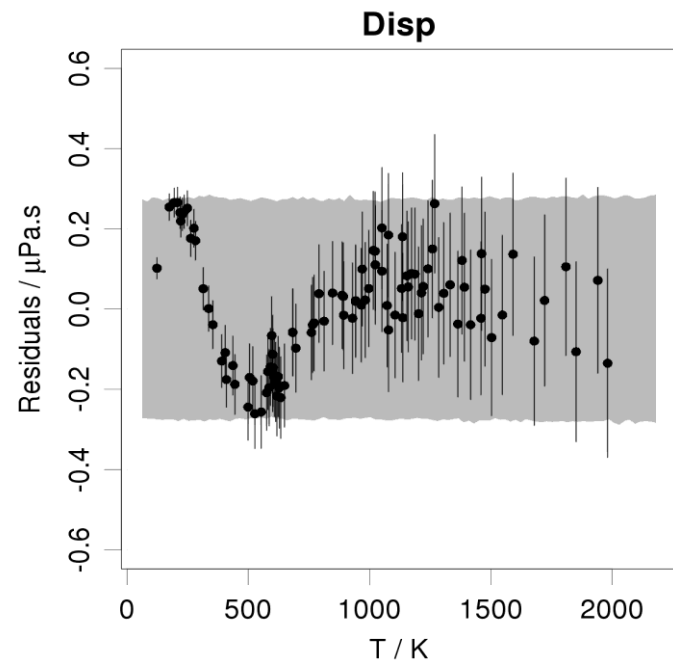
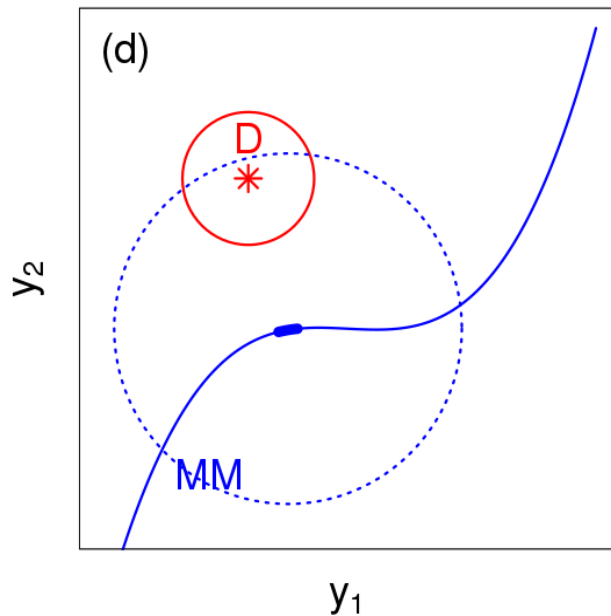
- Non-transferability of the correction to the prediction of another type of data

Kennedy & O'Hagan (2001), J. Roy. Stat. Soc. B, 63: 425-464

Correction at the prediction level (Disp)

- Adding a stochastic term to the model to increase the uncertainty of the prediction:

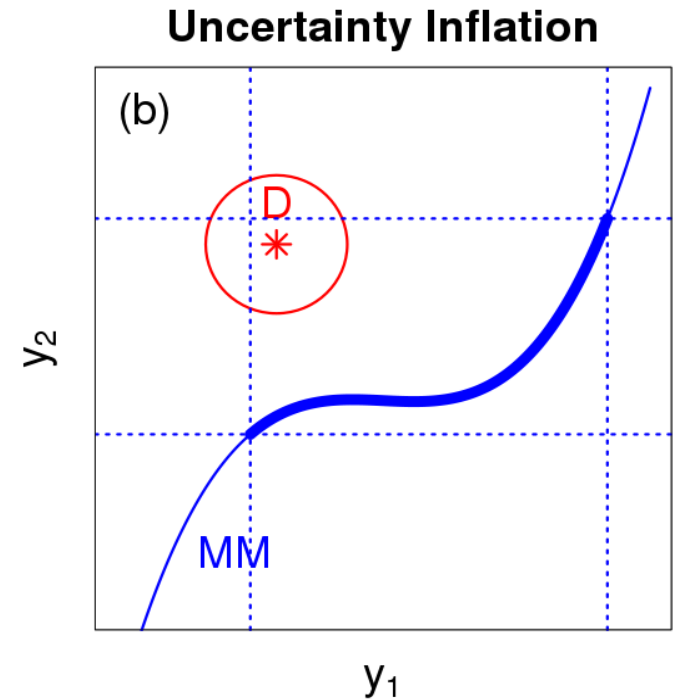
$$y_i = \mathcal{M}(x_i, \theta) + e_D + e_i \quad e_D = \mathcal{N}(0, s^2)$$



- Justified if no trend in the residuals
- Non-transferable to the prediction of another property

Increase parameter uncertainties

- Optimizing the covariance matrix Σ_θ of the parameters
- Variance inflation (VarInf):
 - scaling the covariance matrix obtained from standard calibration: $\Sigma'_\theta = s \times \Sigma_\theta$
- Hierarchical Bayesian framework* (Hier):
 - Divide the dataset D in series D_i
 - Calibrate parameters θ_i for each D_i
 - Find hyperparameters to reproduce the distribution of θ_i : $\theta_i \sim \mathcal{N}(\mu_\theta, \Sigma_\theta)$
- « Direct » stochastic modeling** (ABC):
 - $\mathcal{M}(\theta) \rightarrow \mathcal{M}(\theta, \Sigma_\theta)$
 - Optimize $p(\theta, \Sigma_\theta | D)$

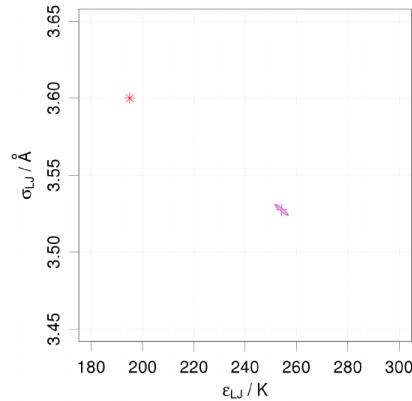
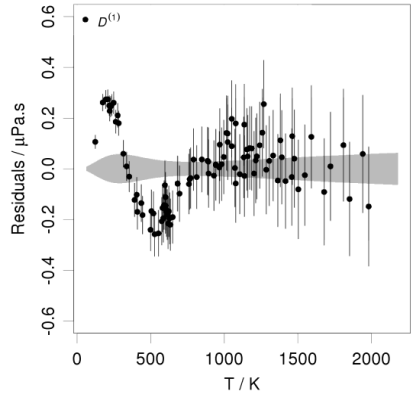


* Wu et al. (2015), *Phil. Trans. R. Soc. A* 374: 20150032

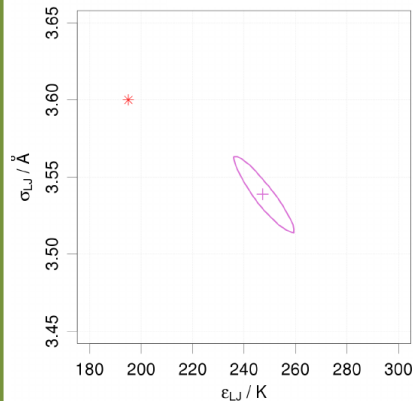
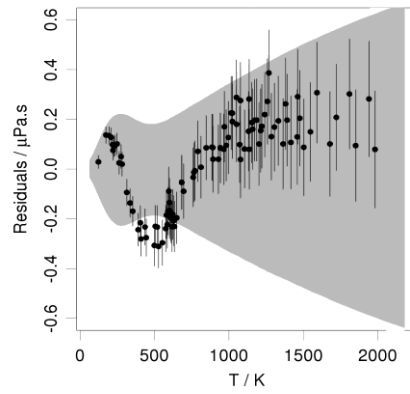
** Csilléry et al., *Trends Ecol Evol.* 2010;25:410–418.

Increase parameter uncertainties

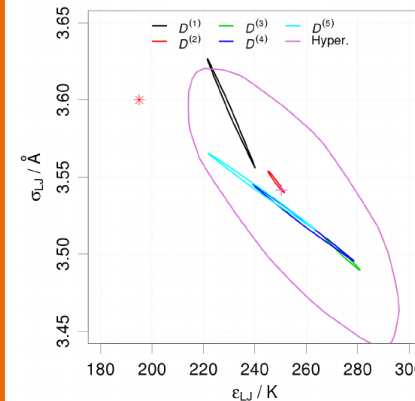
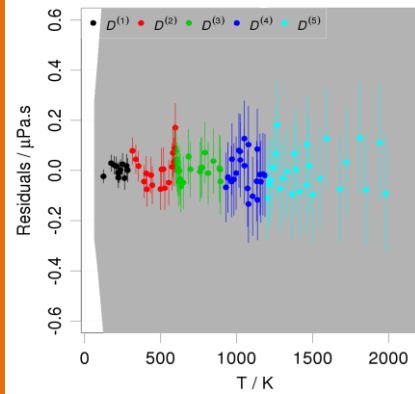
Standard calibration



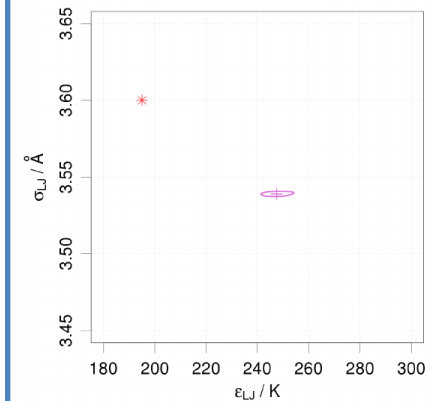
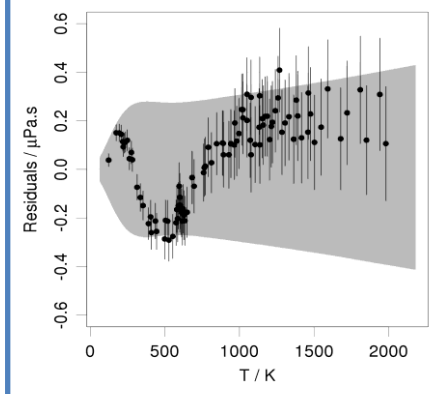
VarInf calibration



Hier calibration



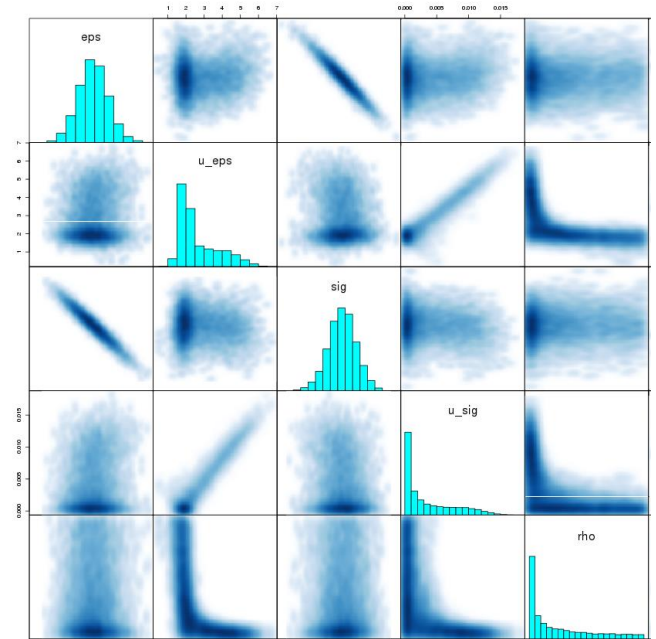
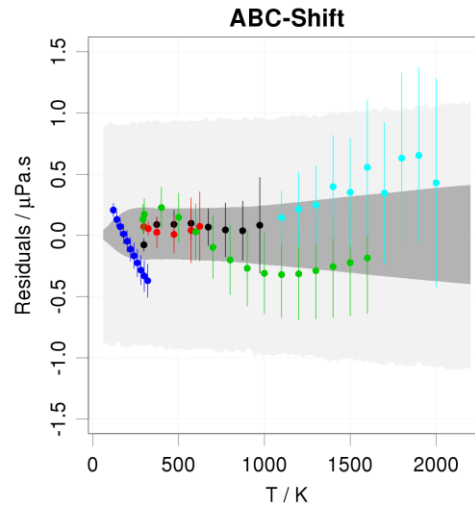
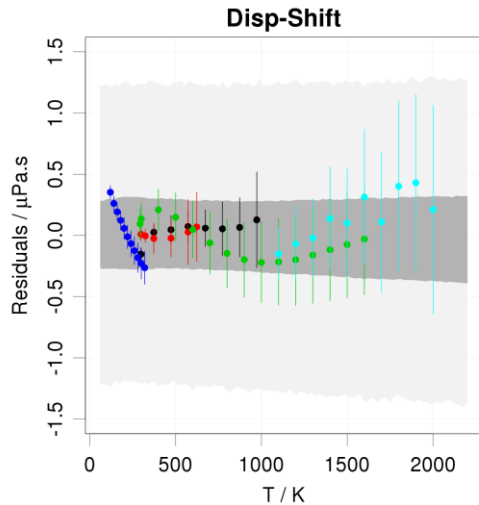
ABC calibration



- VarInf and Hier calibrations: overestimated prediction bands
- ABC calibration: most reasonable option

A real test-case on experimental data

Experimental data for Krypton viscosity



- When transferability to other properties is not an issue, Disp correction is OK

- ABC calibration:

- reasonable but problems of multimodality of the solutions

- Might be improved...

Pernot & Cailliez (2016), arXiv:1611.04376

Concluding remarks

- Study of parameters uncertainties in molecular simulations is still in its infancy
- Bayesian calibration is an adequate framework to determine forcefield parameters and their uncertainties
- Surrogate models and Efficient Global Optimization strategies can be used to alleviate the computational burden of the calibration
- Parametric uncertainties may be the main source of uncertainties in the calculation of fluid properties
- Model inadequacy and transferability to various properties are major issues for quantitative reliable predictions

Acknowledgements

- Laboratoire de Chimie Physique :
 - Pascal Pernot
 - Arnaud Bourasseau
 - Manuel Lopez-Ortiz
 - Jean-Marie Teuler
 - Bernard Rousseau



- Computational facilities :

