

Large scale Bayesian Uncertainty Quantification in Molecular Dynamics simulations

Panagiotis Angelikopoulos

With

Dr. Stephen Wu Dr. Panagiotis Hadjidoukas



Prof. Costas Papadimitriou
@University of Thessaly



- Prof. Robert Moser (UT Austin)
- Prof. Stacey Finley (USC)



Prof. Petros Koumoutsakos

ETH zürich

Overview

- **Motivation**

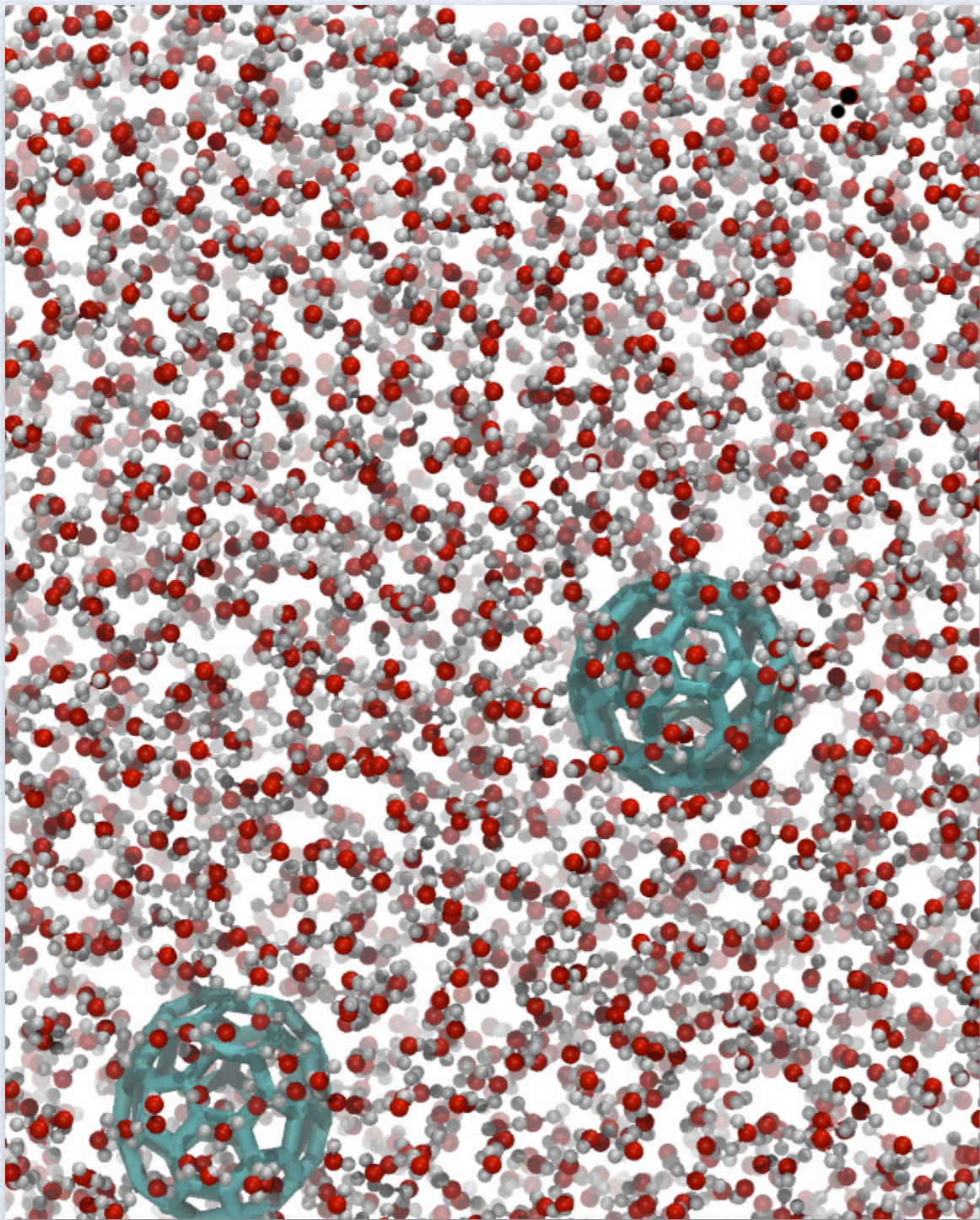
- Molecular Dynamics (MD) - sources of uncertainty
- MD and heterogeneous data fusion

- **Bayesian UQ in MD**

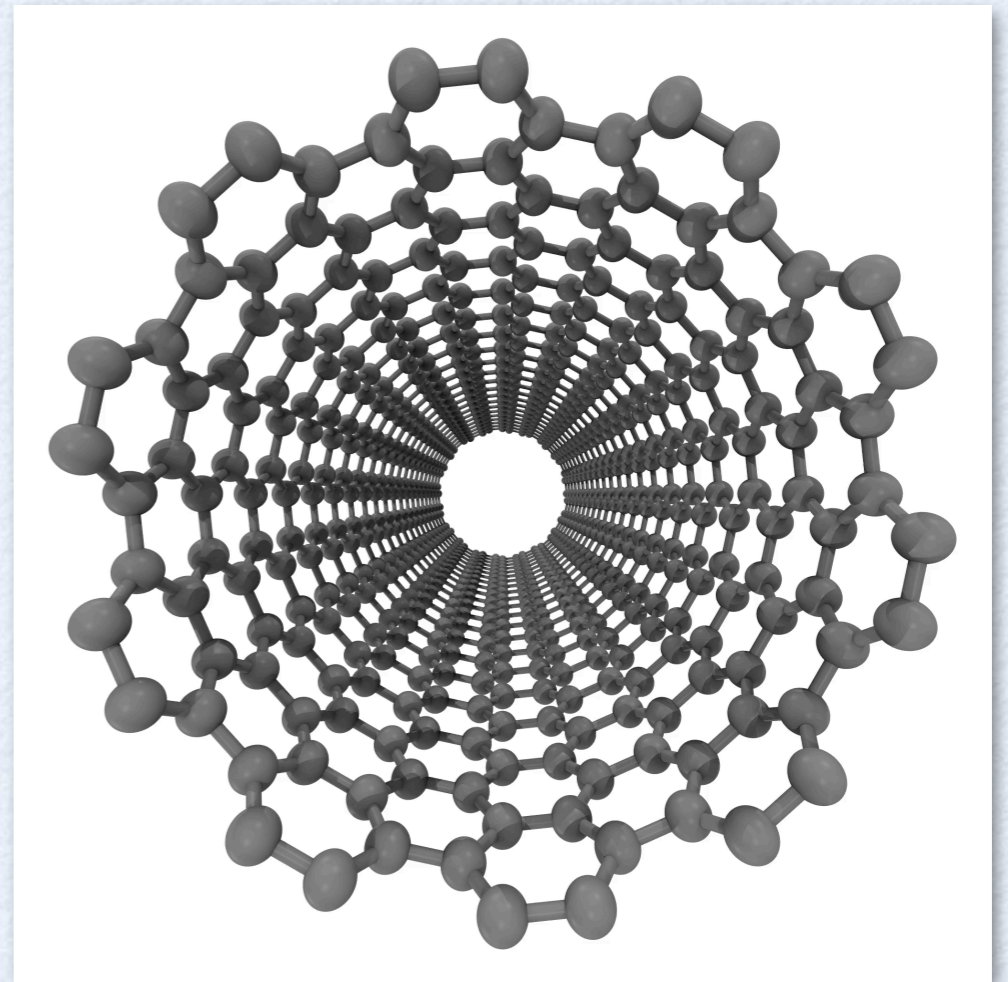
- Algorithms
- Graphical probability models for Bayesian updating
 - *Pooled data* - Water/Graphite interactions
 - *Structured data* - Argon
 - *Heterogeneous data* - Water
- HPC software - Pi4U

- **Summary and Outlook**

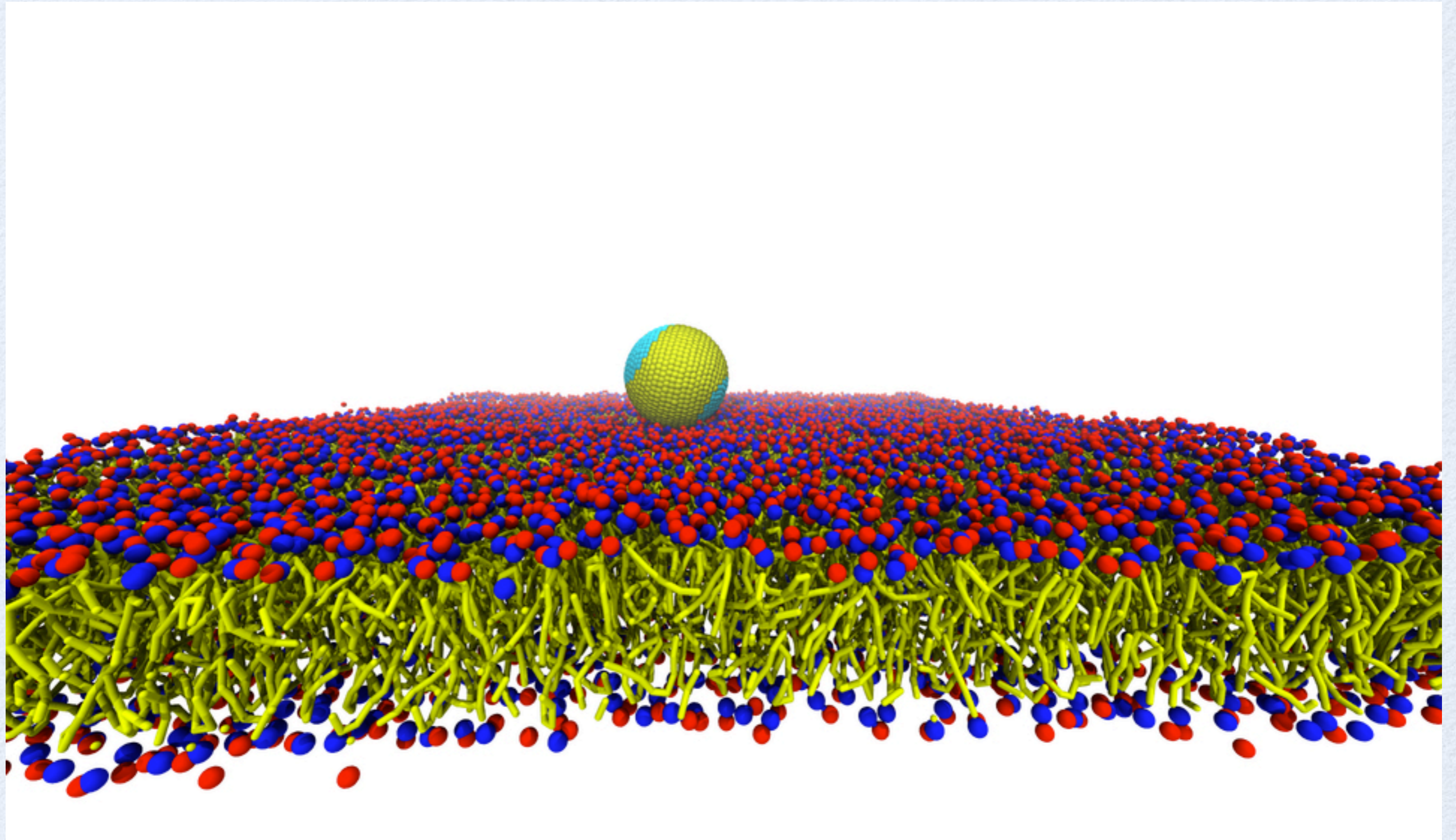
Molecular Dynamics Simulations



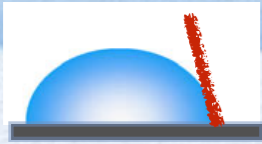
$$\frac{d\mathbf{x}_p}{dt} = \mathbf{u}_p$$
$$m \frac{d\mathbf{u}_p}{dt} = \sum_i^N \mathbf{F}_p \equiv - \frac{\partial U}{\partial \mathbf{x}_p}$$



Molecular Dynamics Simulations



Wetting Dependence on MD Potentials



NATURE | VOL 414 | 8 NOVEMBER 2001 | www.nature

Water conduction through the hydrophobic channel of a carbon nanotube

G. Hummer^{*}, J. C. Rasaiah^{*†} & J. P. Noworyta[†]

Molecular Dynamics Simulation of Contact Angles of Water Droplets in Carbon Nanotubes

Thomas Werder,^{*,†} Jens H. Walther,[†] Richard L. Jaffe,[‡] Timur Halicioglu,[§] Flavio Noca,^{||} and Petros Koumoutsakos^{†,-}

NANO LETTERS

2001
Vol. 1, No. 12
697–702



8 April 2002

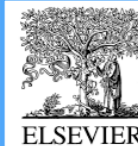
Chemical Physics Letters 355 (2002) 445–448

CHEMICAL PHYSICS LETTERS

www.elsevier.com/locate/cplett

Helical ice-sheets inside carbon nanotubes in the physiological condition

William H. Noon^a, Kevin D. Ausman^b, Richard E. Smalley^b, Jianpeng Ma^{a,c,d,*}



Chemical Physics 247 (1999) 413–430

Chemical Physics

www.elsevier.nl/locate/chemphys

Scattering of water from graphite: simulations and experiments

Nikola Marković^{*}, Patrik U. Andersson, Mats B. Någård, Jan B.C. Pettersson[†]



27 October 2000

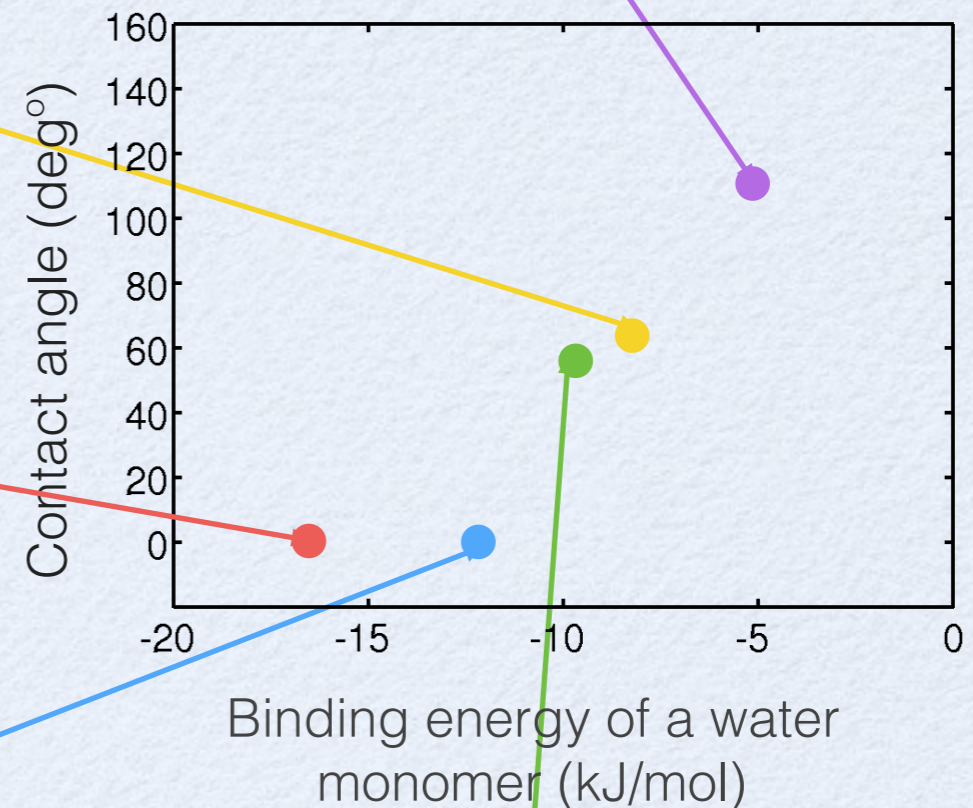
Chemical Physics Letters 329 (2000) 341–345

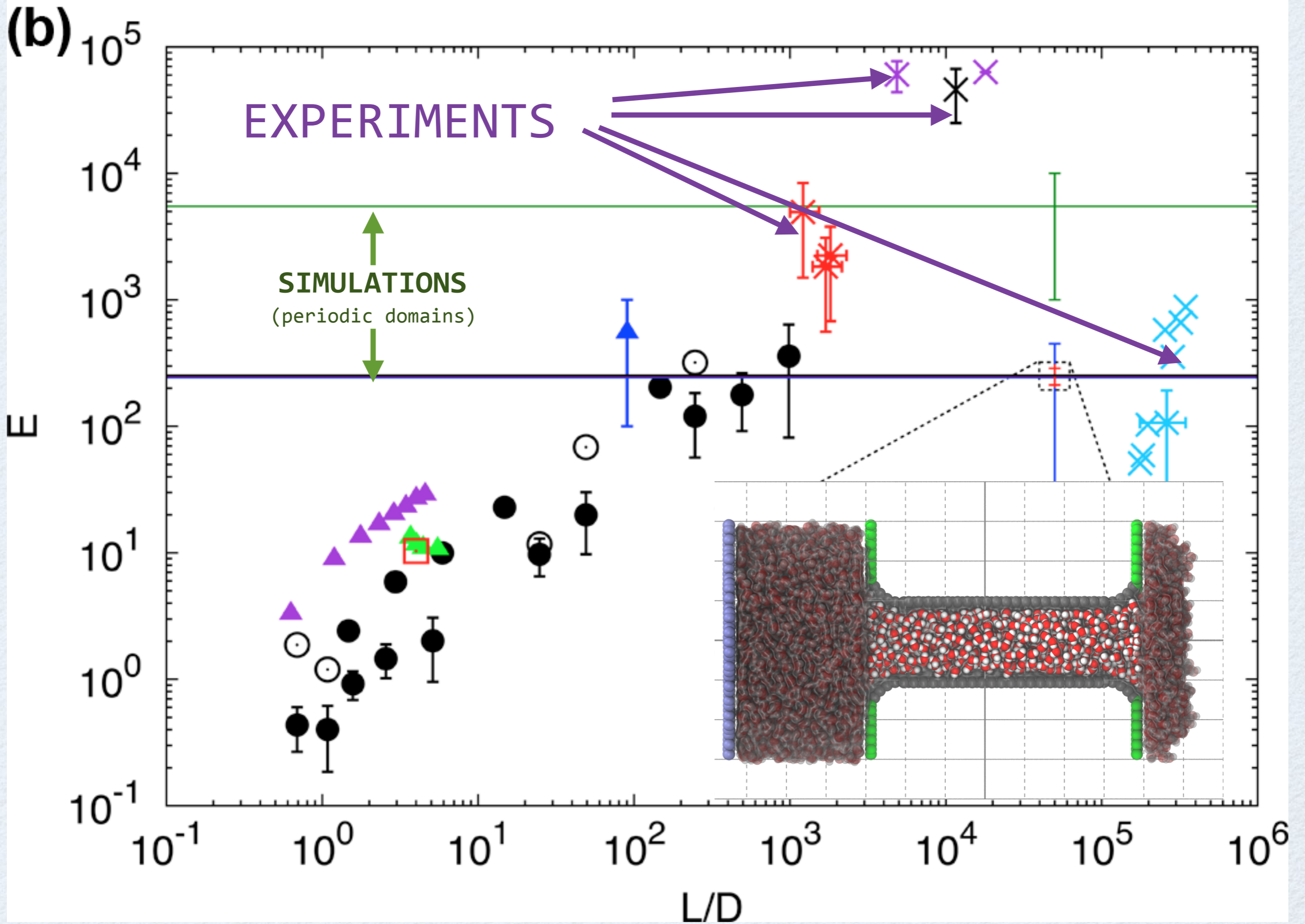
CHEMICAL PHYSICS LETTERS

www.elsevier.nl/locate/cplett

Hydrogen bond structure of liquid water confined in nanotubes

M.C. Gordillo, J. Martí^{*}





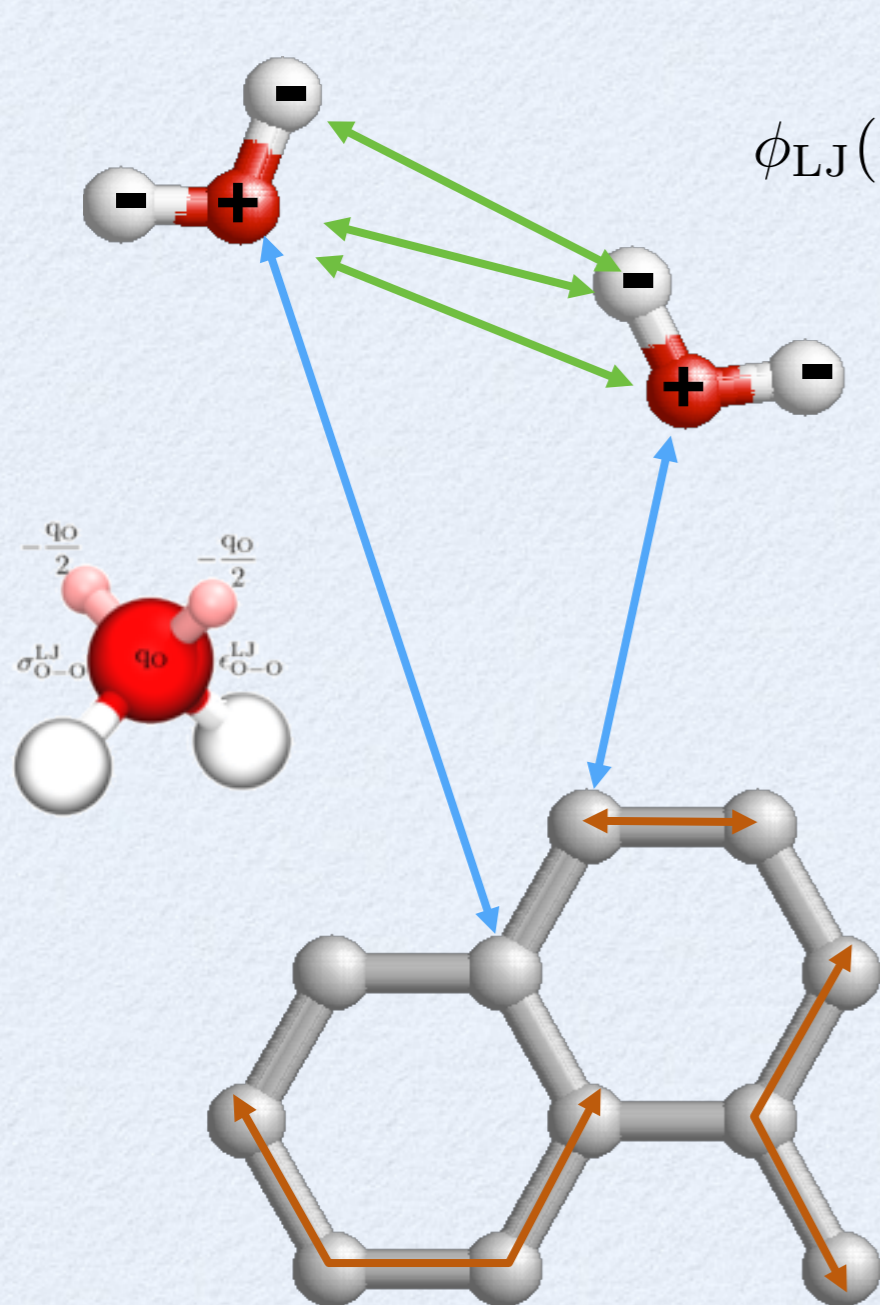
Sources of Uncertainty in Water-Graphite Systems

MODELLING

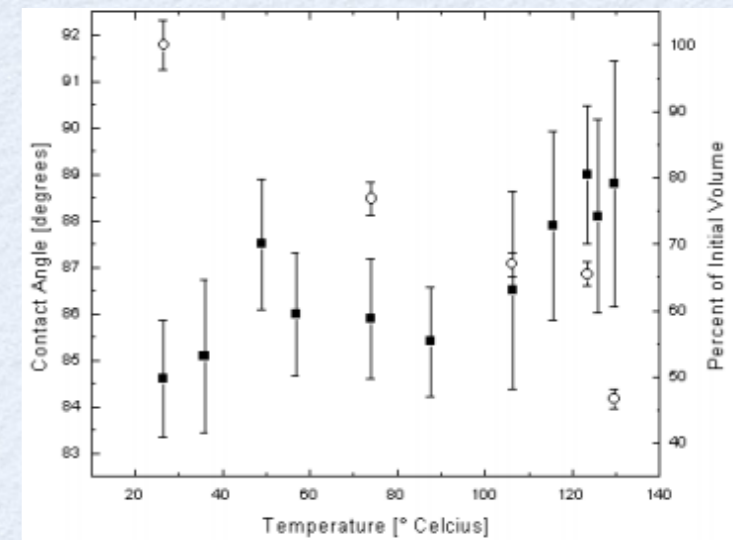
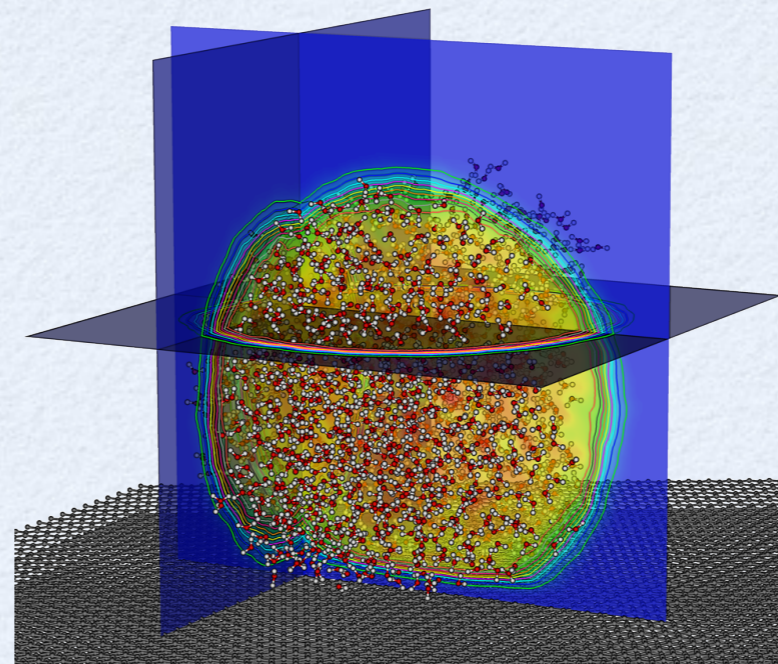
PARAMETRIC

COMPUTATIONAL

MEASUREMENT



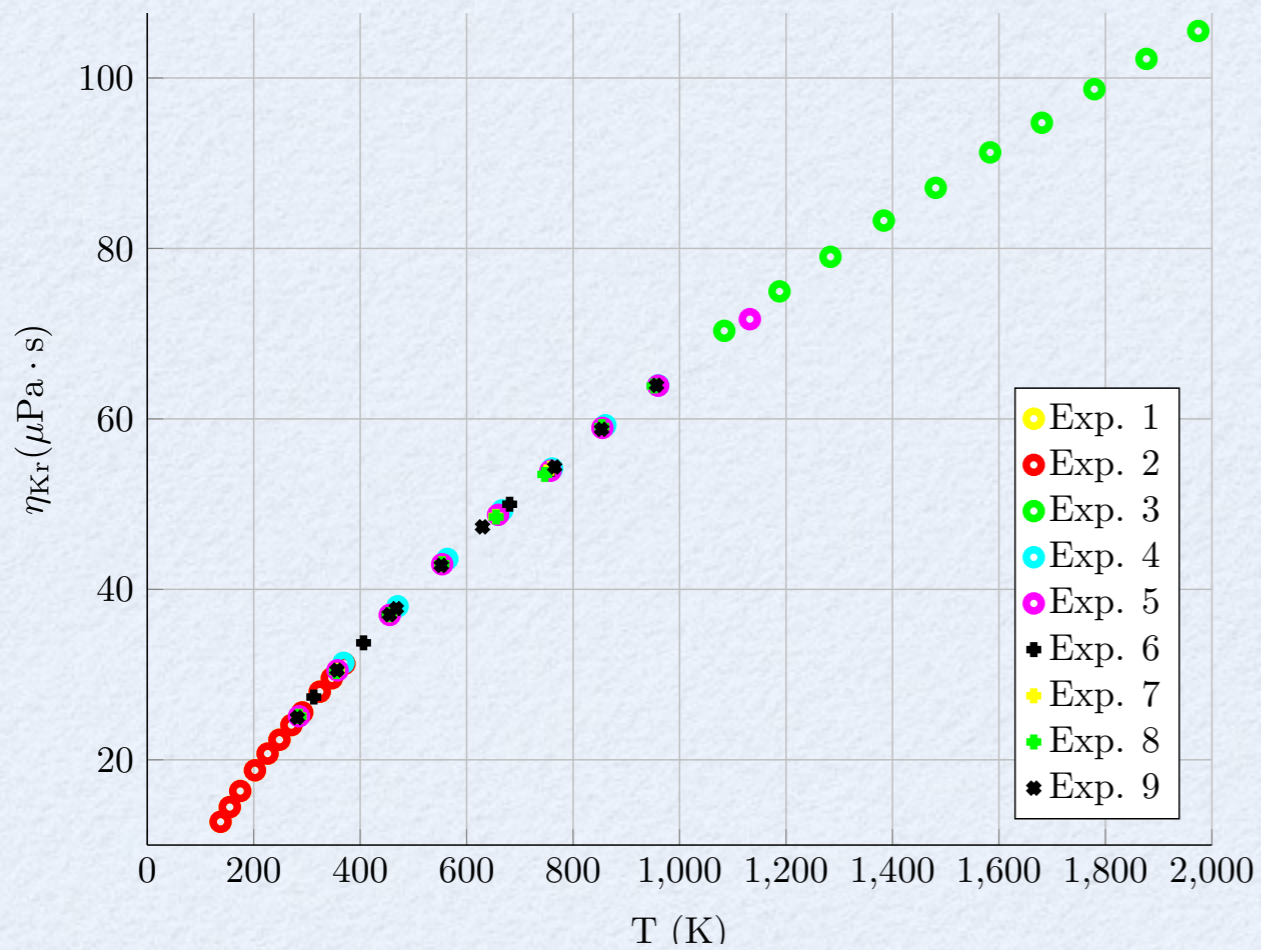
$$\phi_{LJ}(r_{ij}) = 4\epsilon_{LJ} \left[\left(\frac{\sigma_{LJ}}{r_{ij}} \right)^{12} - \left(\frac{\sigma_{LJ}}{r_{ij}} \right)^6 \right]$$



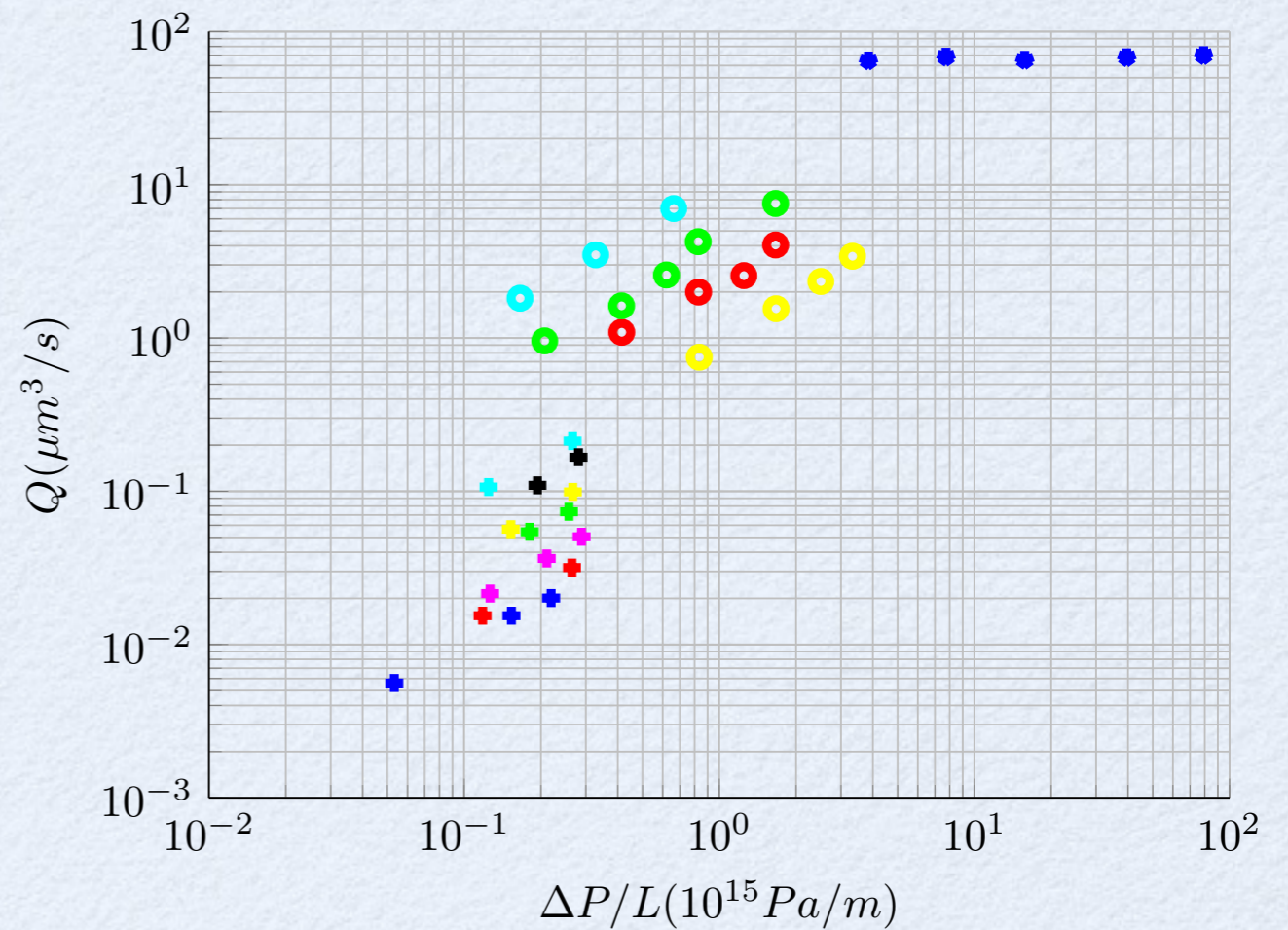
K. Osborne III (2009)

MD Data is Structured

Structure in the data and predictions



Argon viscosity



water flow rate inside carbon nanotubes

MD predictions are Heterogeneous

As stated above, the goal was to develop the simplest potential function, which reproduces well the density anomaly of liquid water, while simultaneously yielding good thermodynamic and structural properties near 25 °C and 1 atm. Additional studies of three- and four-site models, in-

Mahoney & Jorgensen, J. Phys. Chem. 2000

**Multiple
Objective
Functions**

- ❖ Pareto front
- ❖ Weighted sum of objective functions



How to
choose
weights?

“Everybody trusts an experiment, but the person that did it.

No-one believes *a simulation* but the person that did it

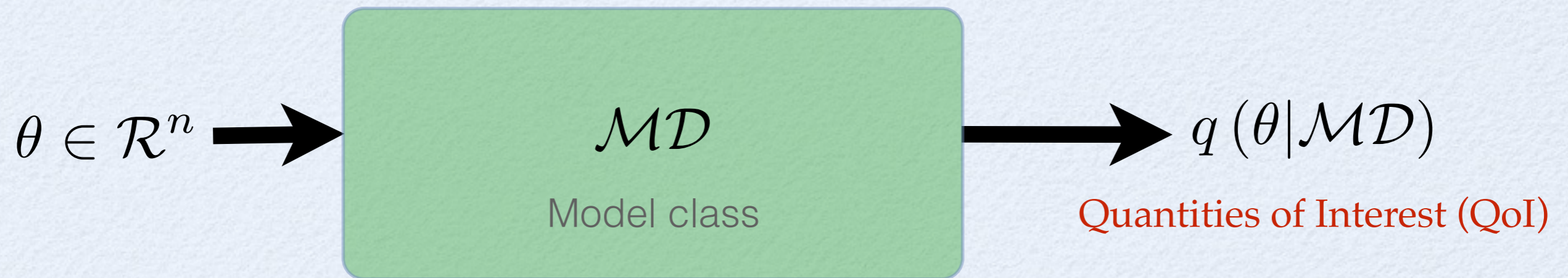
**Why not combine the two and get results everybody can
mistrust a little?”**

Tony Kordyban

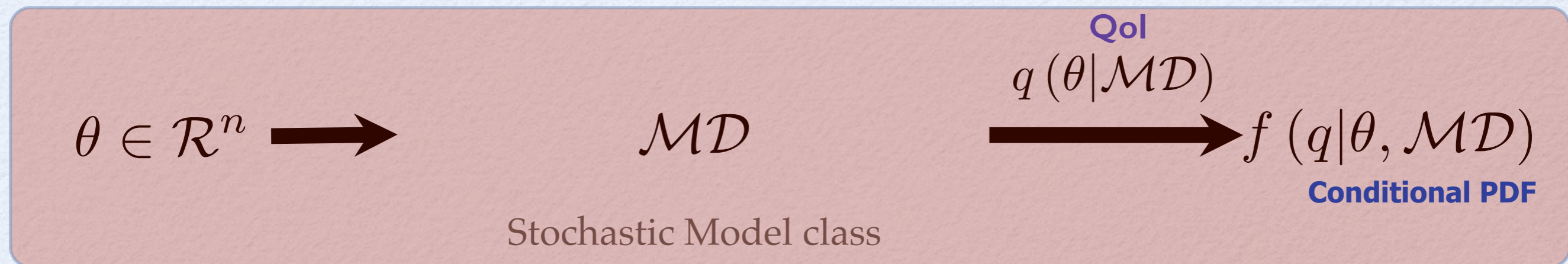
Bayesian Inference and Uncertainty Quantification

Embedding the model in a stochastic model class

From deterministic



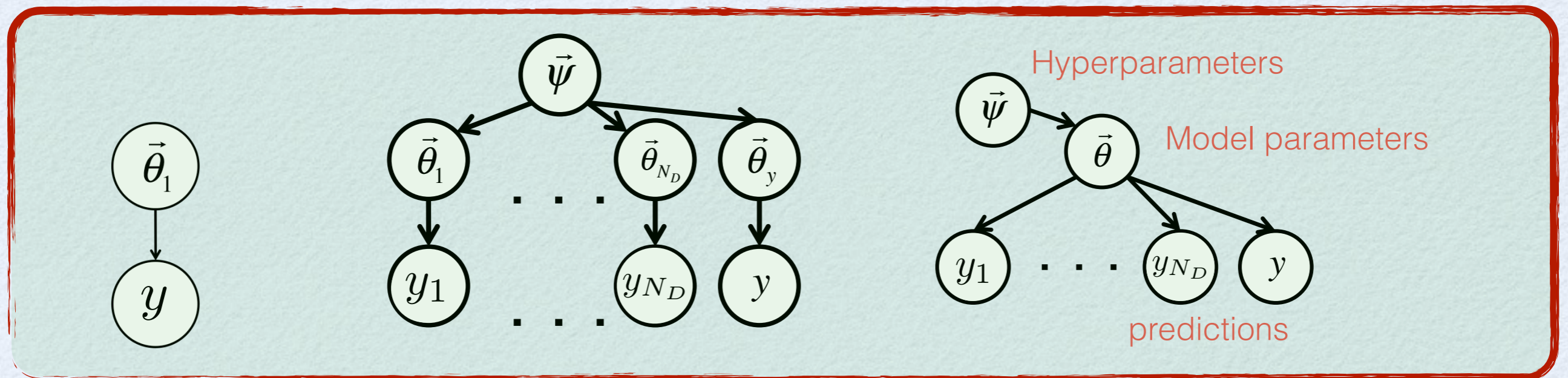
To Stochastic



Describing the stochastic model class : Graphical probabilistic models

Graph based models of a joint probability distribution

A way to describe **conditional independencies**- valid for any type of distributions (not just Gaussians)



e.g
$$p\left(\vec{\theta}_1, \vec{\theta}_2 | \vec{\psi}\right) = p\left(\vec{\theta}_1 | \vec{\psi}\right) \cdot p\left(\vec{\theta}_2 | \vec{\psi}\right)$$

- Types :
- 1) directed acyclic graphs (Bayesian networks)
 - 2) undirected graphs (Markov Random Fields)

Bayesian UQ : Calibration and Model Selection

Experimental Data: D

Use observations to select the model classes and estimate their parameter values such that the model predictions best fit the data

PARAMETER ESTIMATION

$$f(\theta_i | D, \mathcal{MD}_i) = \frac{f(D | \theta_i, \mathcal{MD}_i) \pi(\theta_i | \mathcal{MD}_i)}{f(D | \mathcal{MD}_i)}$$

Experiments Physical limitations
Past studies
Expert elicitation

MODEL CLASS SELECTION

$$Pr(\mathcal{MD}_i | D) = \frac{f(D | \mathcal{MD}_i) Pr(\mathcal{MD}_i)}{f(D)}$$

Evidence of Model

$$f(D | \mathcal{MD}_i) = \int f(D | \theta_i, \mathcal{MD}_i) \pi(\theta_i | \mathcal{MD}_i) d\theta_i$$

Bayesian Uncertainty Propagation



QUANTITIES OF INTEREST: Posterior Robust Predictions: PDF

$$f(q|D, \mathcal{MD}) = \int \underbrace{f(q|\theta, \mathcal{MD})}_{\text{Conditional PDF}} \underbrace{f(\theta|D, \mathcal{MD})}_{\text{Posterior PDF}} d\theta$$

$$f(q|D, \mathcal{MD}) = \frac{1}{N} \sum_{i=1}^N \underbrace{f(q|\theta^{(i)}, \mathcal{MD})}_{\text{Conditional PDF}} \quad \theta^{(i)} \sim f(\theta|D, \mathcal{MD})$$

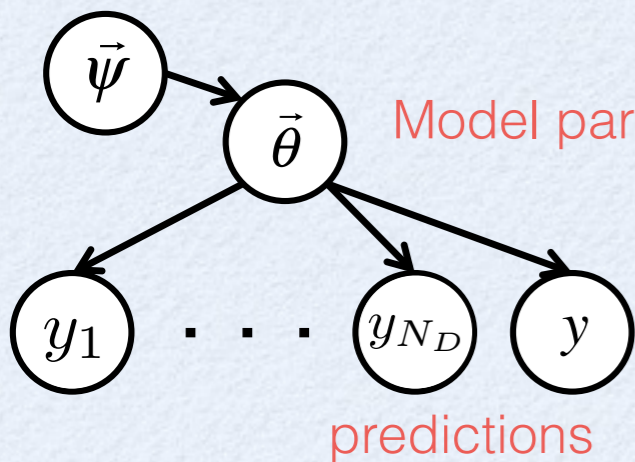
Samples drawn from Posterior PDF

Basic data structure : pooled data

DISCREPANCY BETWEEN
MODEL PREDICTIONS AND
DATA

$$\hat{y} = g(\theta | \mathcal{MD}) + e$$

$$e = e^d + e^f + e^m$$



- **Measurement Error**

$$e^d \sim N(\mu^d, \Sigma^d)$$

$$\Sigma^d = \text{diag}(s_r^2 \hat{y}_r^2)$$

- **Computational Error**

$$e^f \sim N(\mu^f, \Sigma^f)$$

$$\Sigma^f = \text{diag}(\sigma_r^2 \hat{y}_r^2)$$

- **Model Error**

$$e^m \sim N(\mu^m, \Sigma^m)$$

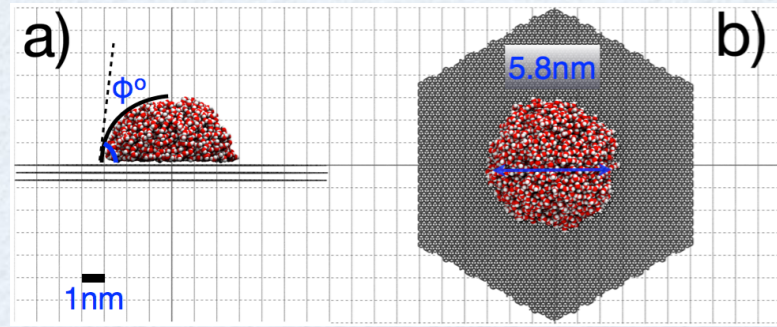
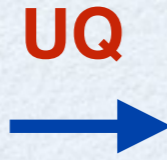
$$\Sigma^m = \text{diag}(\lambda_r^2 \hat{y}_r^2)$$

0

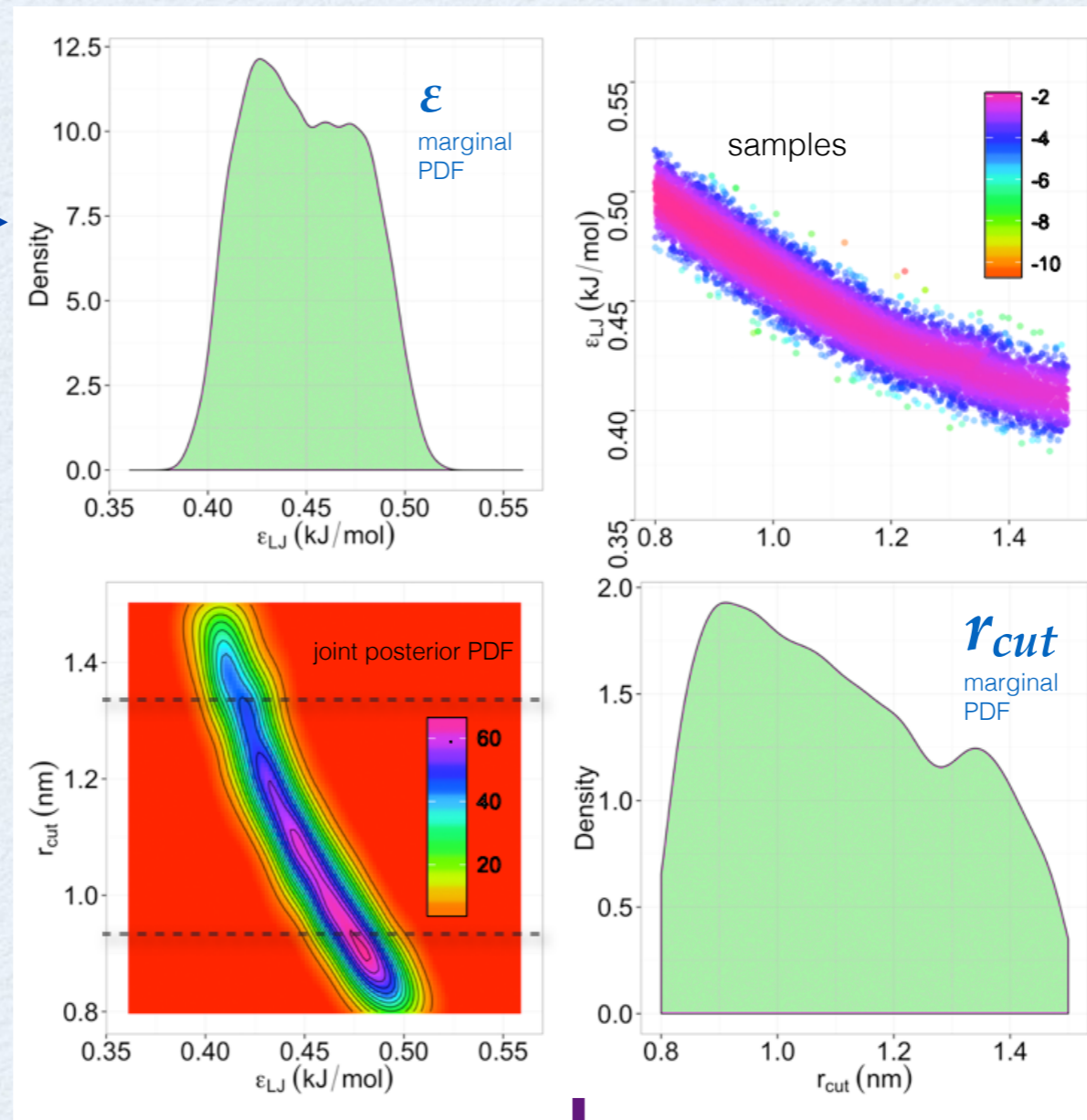
LIKELIHOOD: Assume independent model prediction errors $\Sigma(\theta_e) = \Sigma^d + \Sigma^f + \Sigma^m$

$$p(D|\theta, \mathcal{MD}) = \frac{|\Sigma|^{-1/2}}{(2\pi)^{m/2}} \exp \left[-\frac{1}{2} [\hat{y} - f(\theta_m | \mathcal{MD}) - \mu(\theta_e)]^T \Sigma^{-1}(\theta_e) [\hat{y} - f(\theta_m | \mathcal{MD}) - \mu(\theta_e)] \right]$$

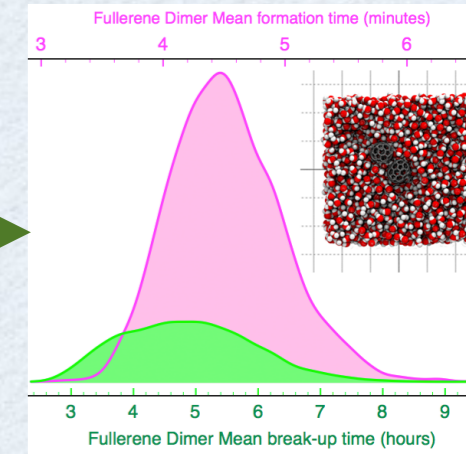
water contact angle



parameters: ϵ , r_{cut}



buckyball aggregation

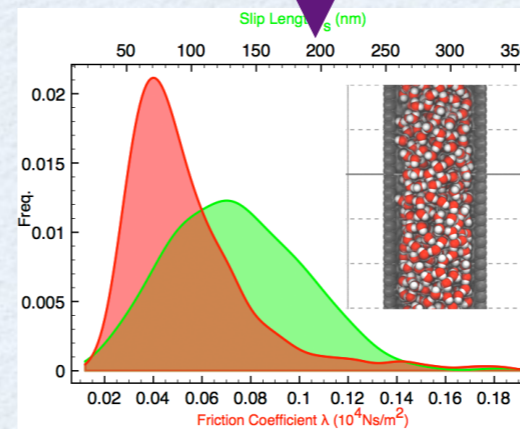


200 samples - 7 days

HPC challenges for UQ in nanoscale flows:

large, variable, computational cost per calibration or propagation sample

Resources: 1200 cores, 32GPUs



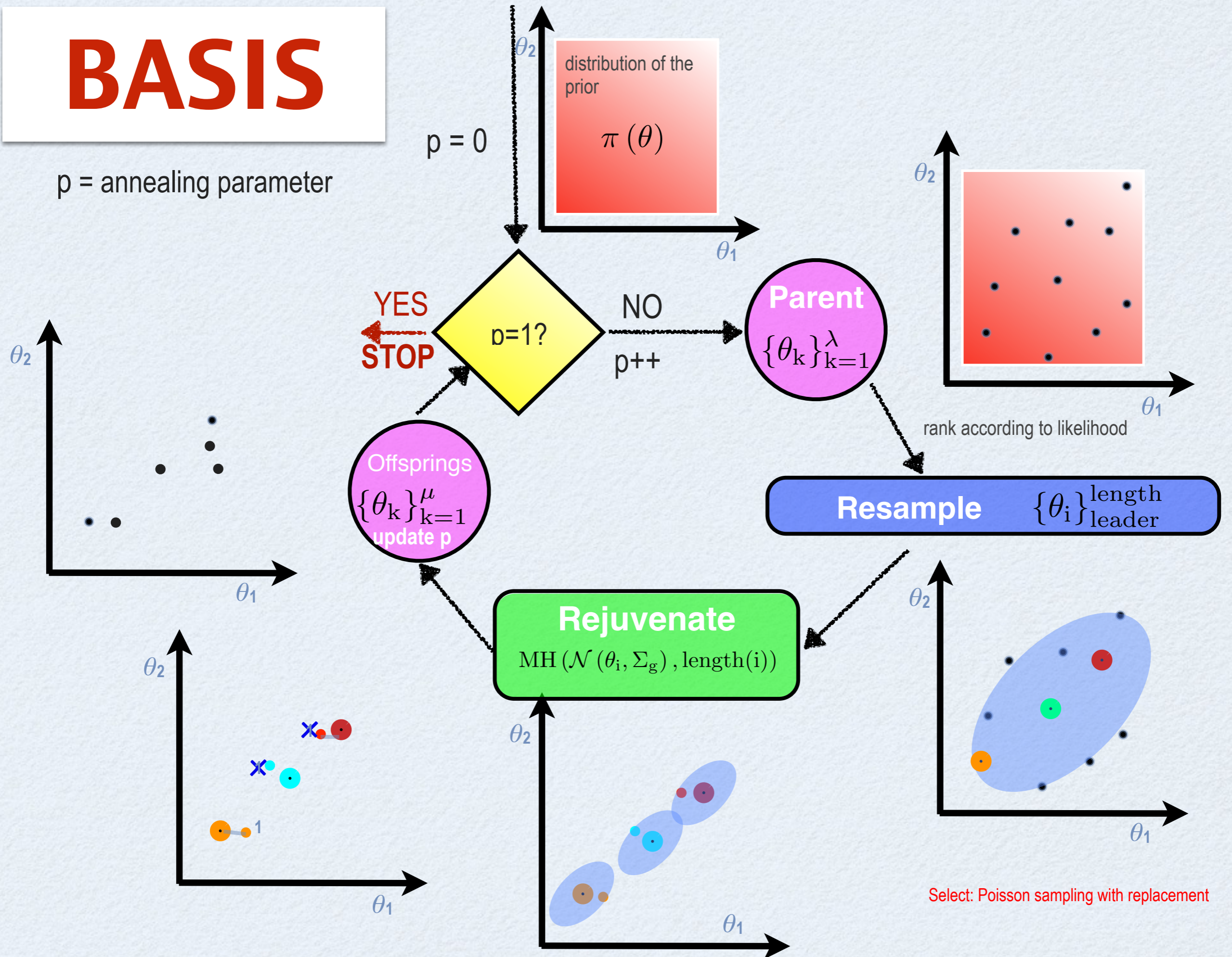
Friction coefficient and slip length of water inside CNTs

500 samples - 2 days

water transport in CNTs

BASIS

p = annealing parameter



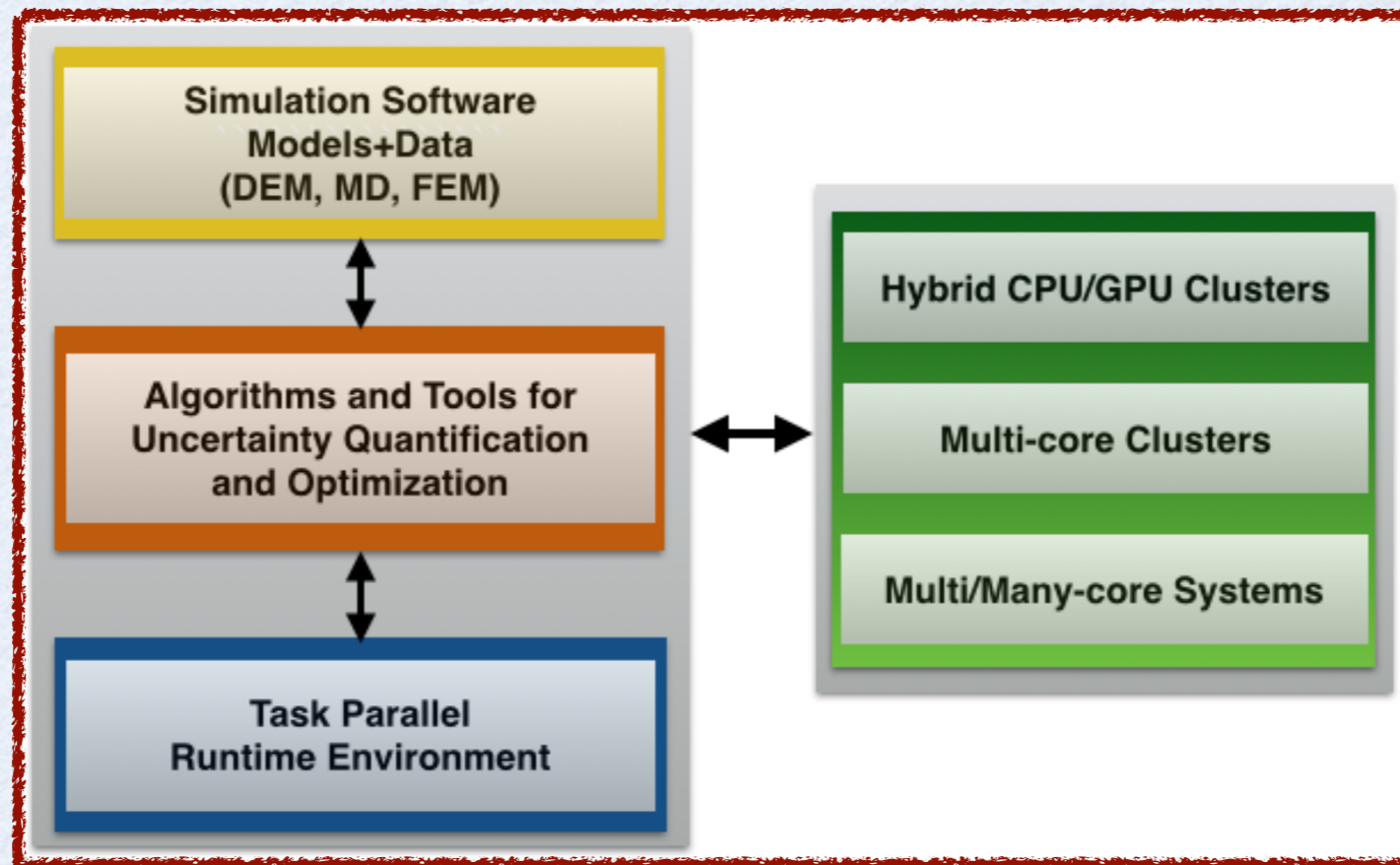
Π4U: HPC Framework for Bayesian UQ

Cluster Setup:

Piz Daint Cray
XC30, 512 nodes
x 8-core Intel
Xeon E5-2670 + 1
Tesla GPU

=> **4,096 cores +
512 GPUs in total**

96% efficiency



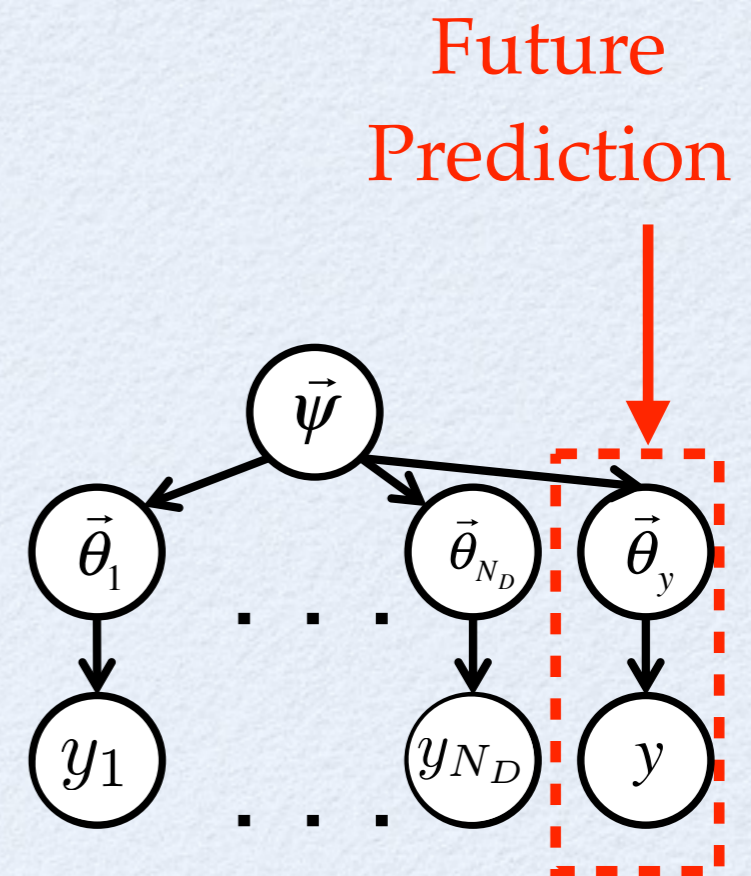
- ✓ Automatic Load balancing
- ✓ Multi-level Nested parallelism
- ✓ Extensible/ built upon tasking Library TORC
- ✓ Slurm/ LSF queuing systems compatibility
- ✓ Intel Phi support

Algorithms implemented:

- ✓ BASIS
- ✓ SubSet simulation
- ✓ ABC-Subsim
- ✓ CMA-ES/AMalgam

“Hierarchical” Bayesian Framework

- Given N sets of heterogeneous data:
 1. Calibrate θ_i for each data set D_i
 2. Link all θ_i with hyperparameters ψ
 3. Combine heterogeneous data sets based on Evidence of D_i in Bayesian inference
- Evidence combines data-fitting & Ockham's razor



Pitfalls of more complicated probability graphs

$$p(\theta|D) = \frac{p(D|\theta)p(\theta)}{p(D)}$$

$$p(\theta|D) = \int p(\theta|\psi)p(\psi|D) d\psi$$

Double
Integral

Nested
MCS

Very High

Computational cost

$$p(\psi|D) = \frac{p(D|\psi)p(\psi)}{p(D)}$$

$$p(D|\psi) = \prod_{i=1}^N p(D_i|\psi) \quad \text{Evidence of } D_i$$

$$= \prod_{i=1}^N \int p(D_i|\theta_i, \psi)p(\theta_i|\psi) d\theta_i$$

Importance Sampling for marginal likelihood

$$p(D_i|\psi, \sigma_y) = \int \frac{p(D_i|\theta_i, \sigma_y)p(\theta_i|\psi)}{\pi(\theta_i|D_i)} \pi(\theta_i|D_i) d\theta_i$$

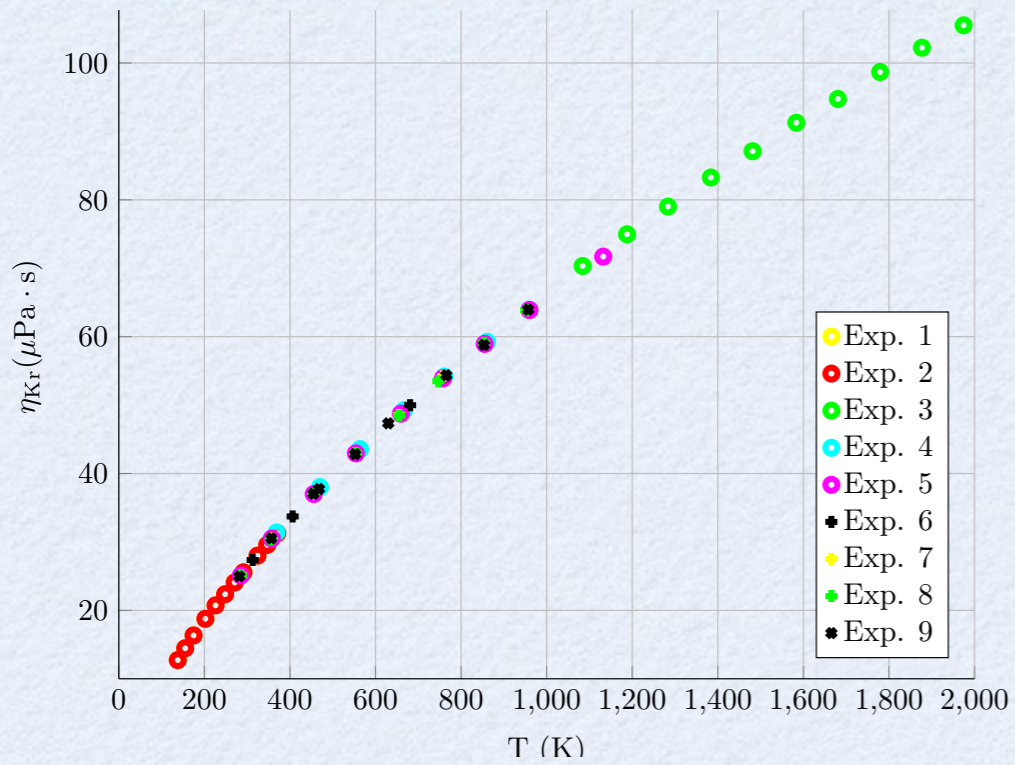
$$\left[\theta_i^{(j)} \sim \pi(\theta_i|D_i) \right] \approx \frac{1}{N_{s,i}} \sum_{j=1}^{N_{s,i}} \frac{p(D_i|\theta_i^{(j)}, \sigma_y)p(\theta_i^{(j)}|\psi)}{\pi(\theta_i^{(j)}|D_i)}$$

Basic idea: use posterior samples for each data D_i to estimate the integral with varying (ψ, σ_y)

Pros: Usefull to do UQ on individual data set anyway, HB becomes a post-processing step with no extra likelihood evaluations

Cons: potentially high variance for the estimation

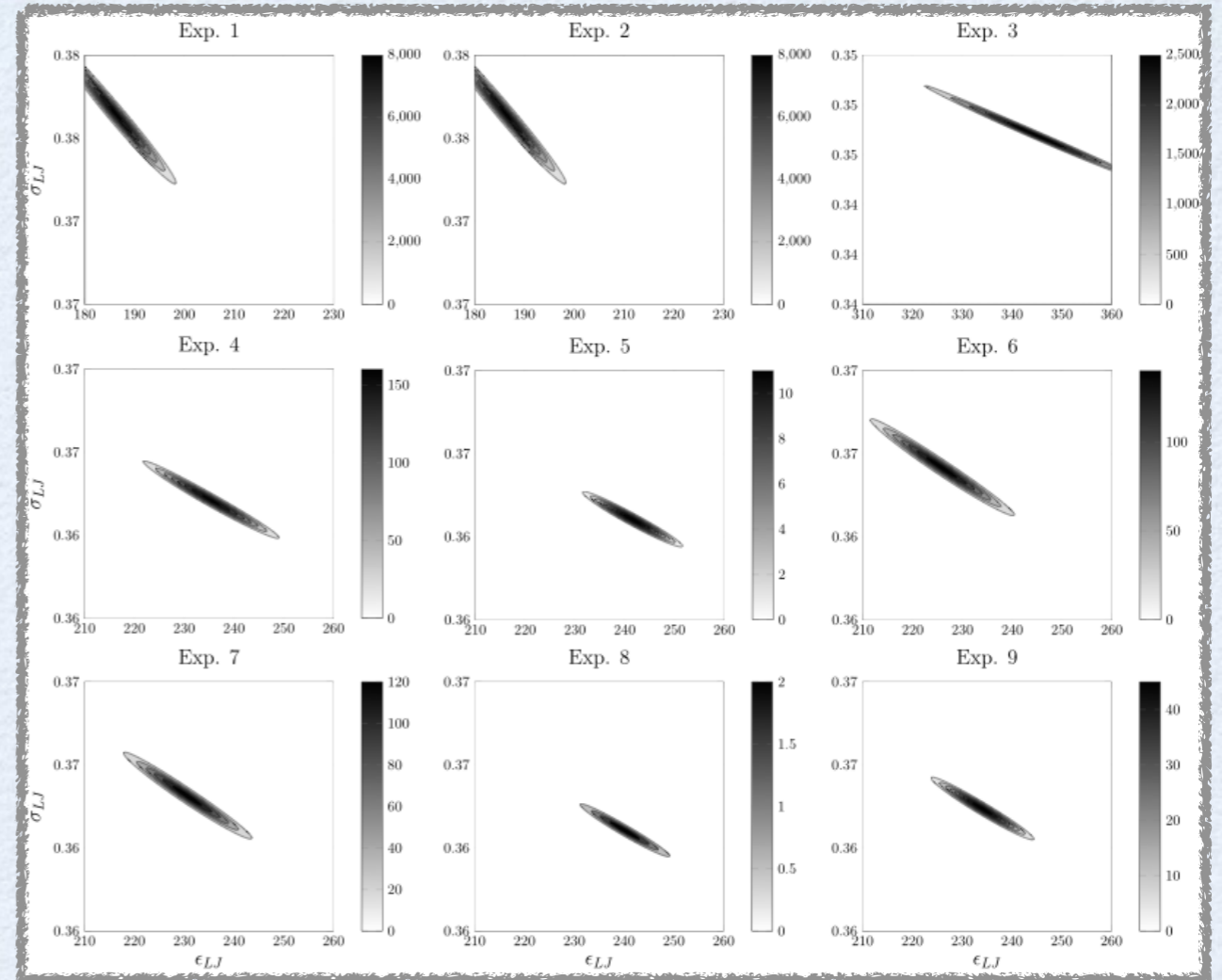
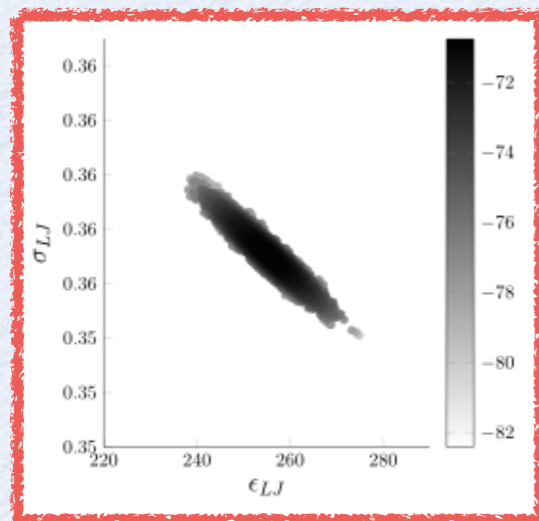
Data structure matters – Argon



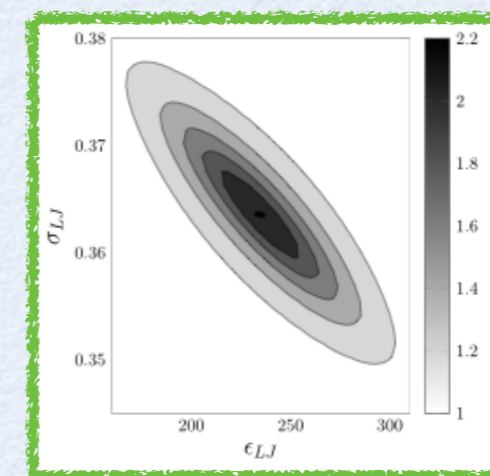
Argon viscosity

Kestin et al. 1984

Pooled Data



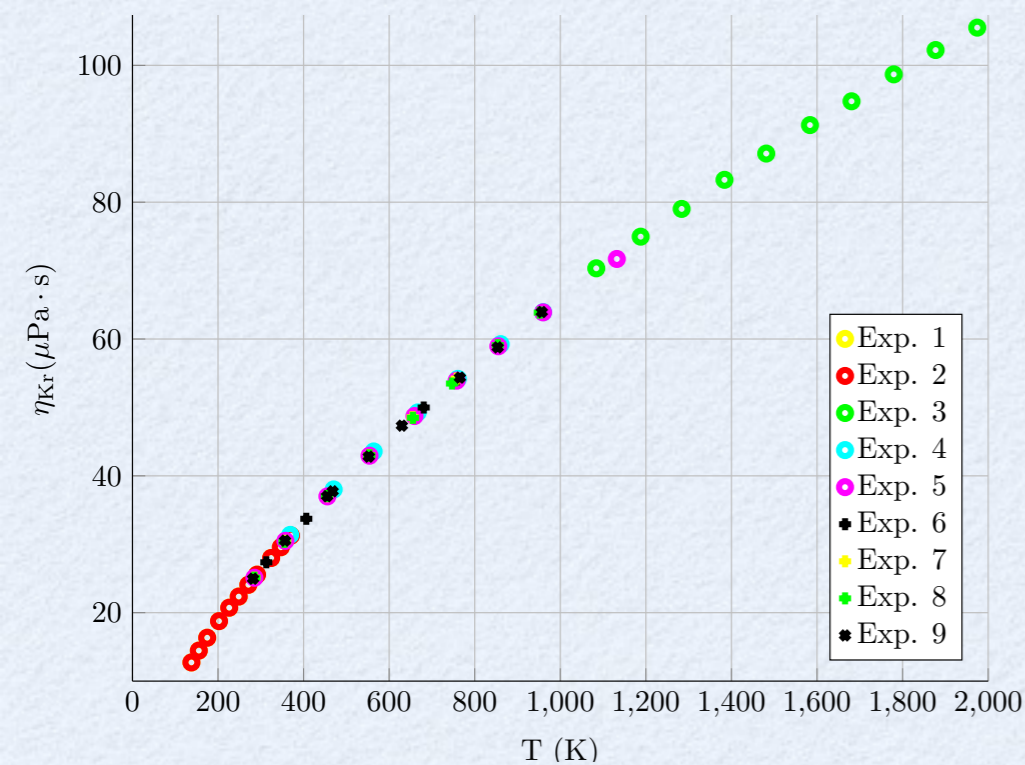
with structured Data



Identifying the plausible structure of the Data

Parameters: Lennard-Jones

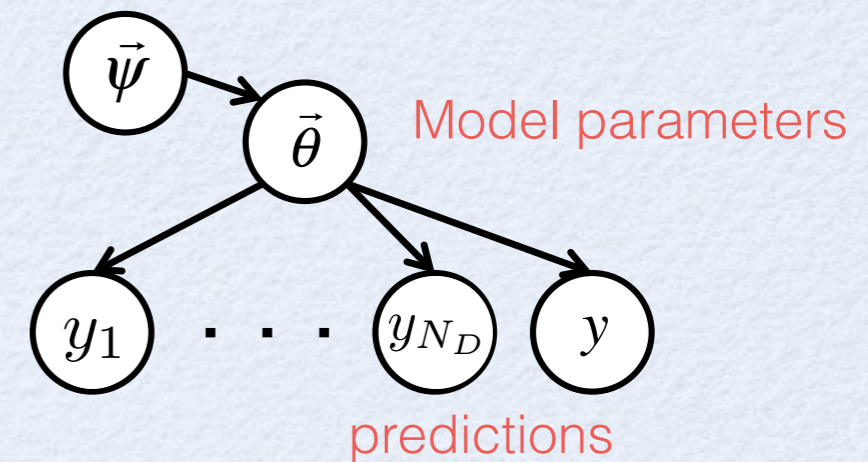
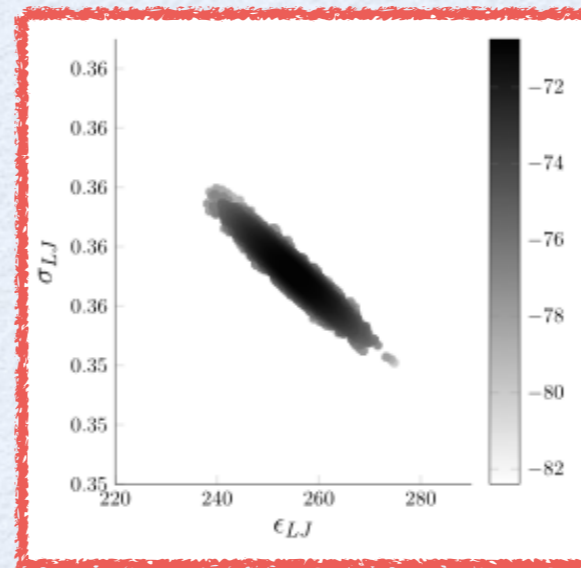
$\epsilon_{LJ}, \sigma_{LJ}$



Argon viscosity
Kestin et al. 1984

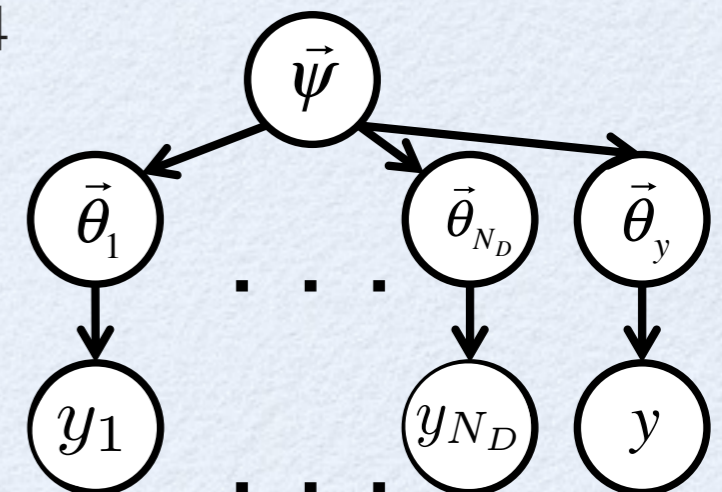
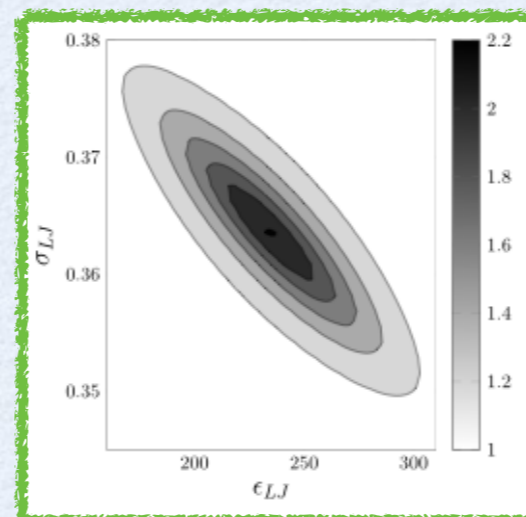
Pooled Data

M_1 Evidence : -77.02



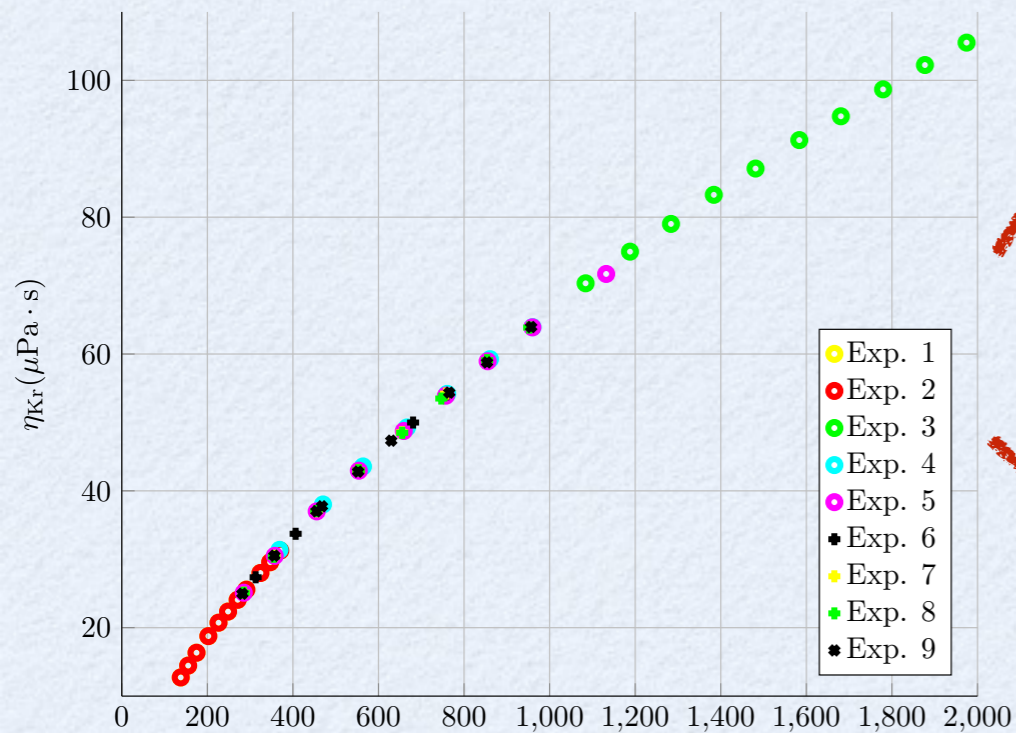
Tree Data Structure

M_2 Evidence : -20.64

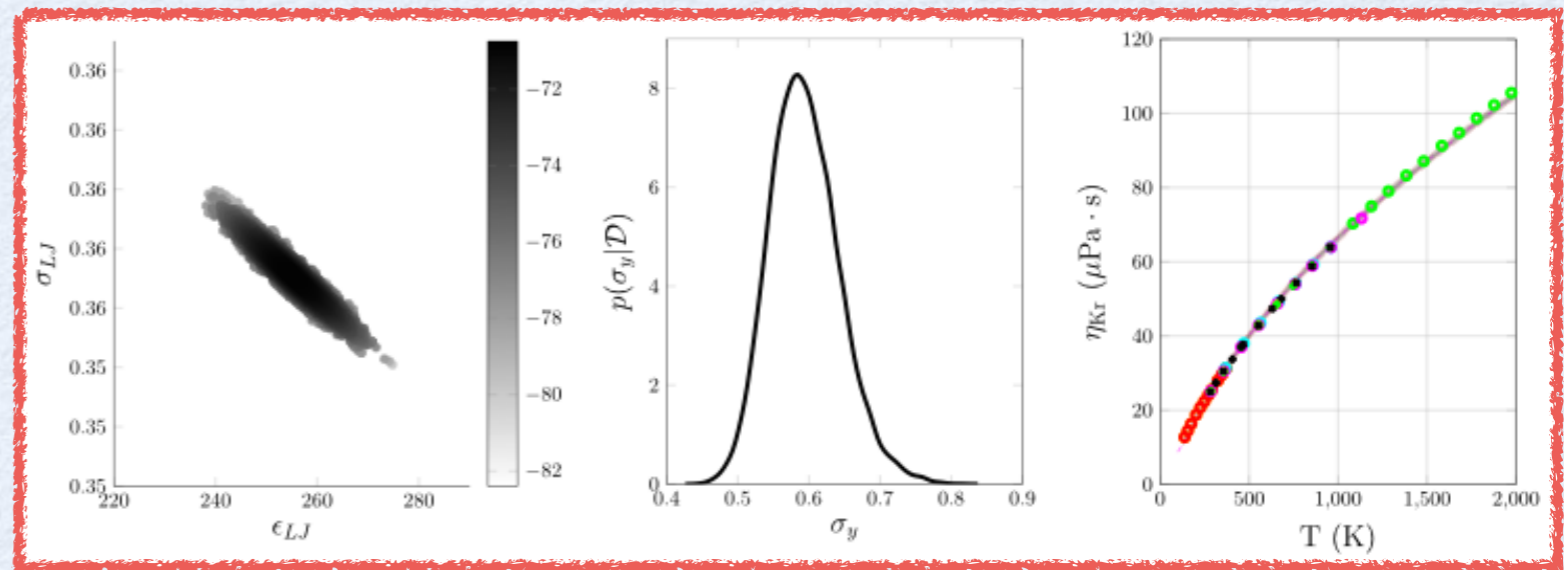


Effect of graphical model on prediction

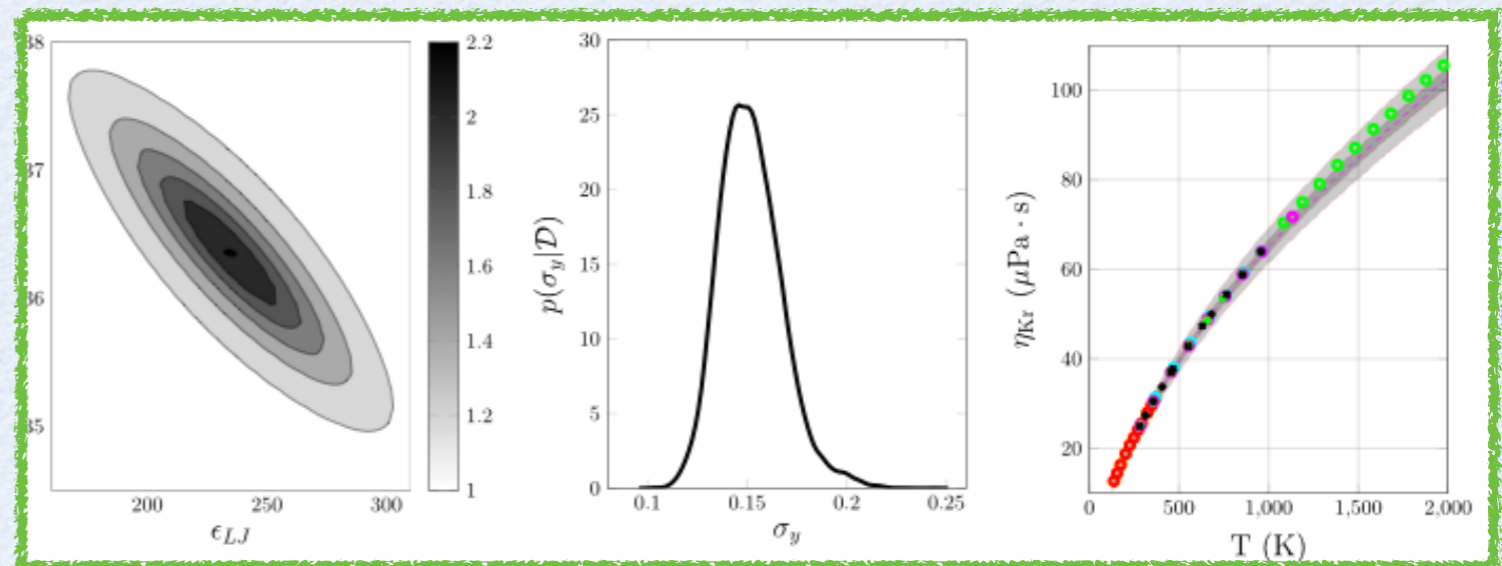
Effect on robust prediction



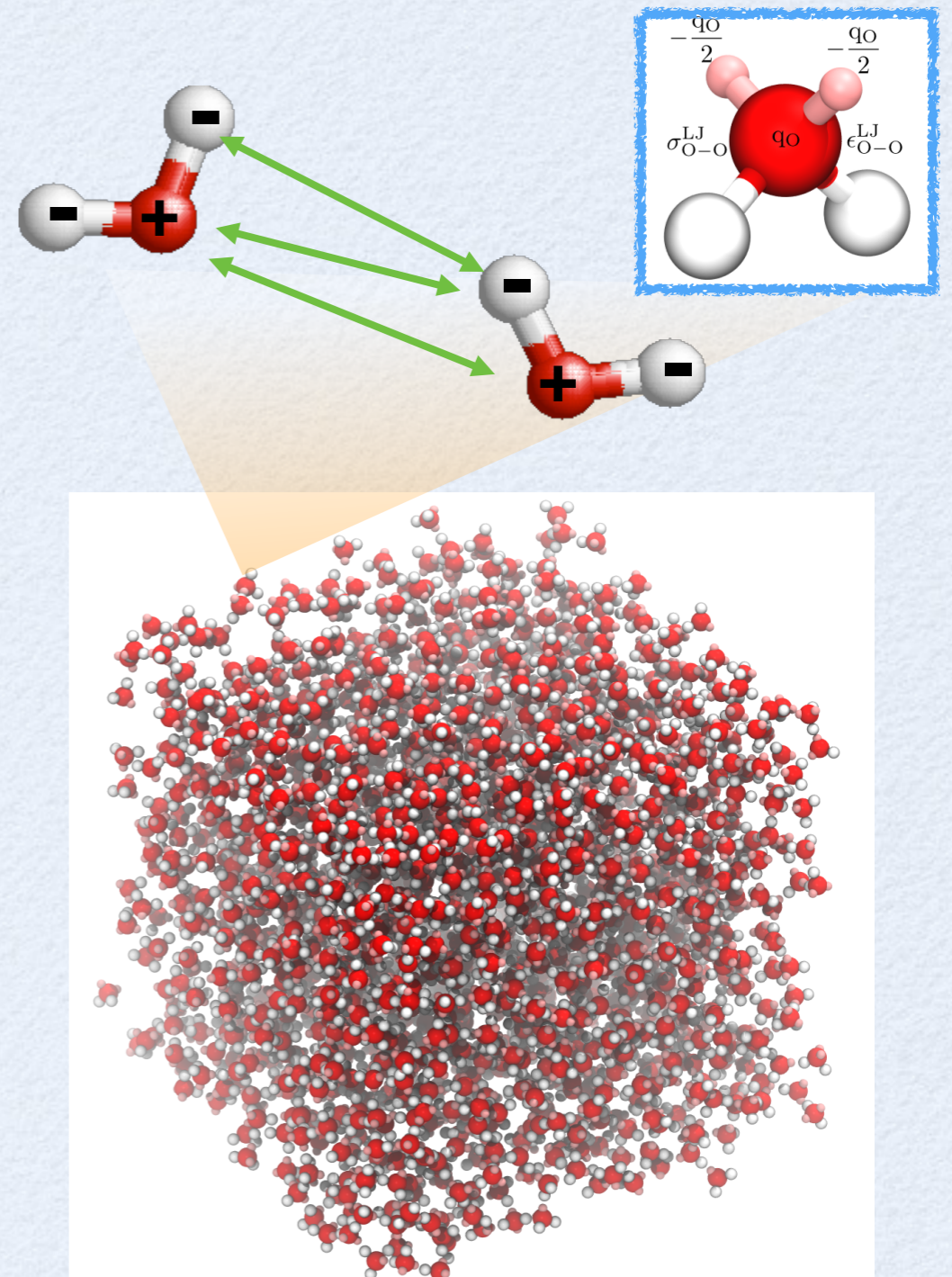
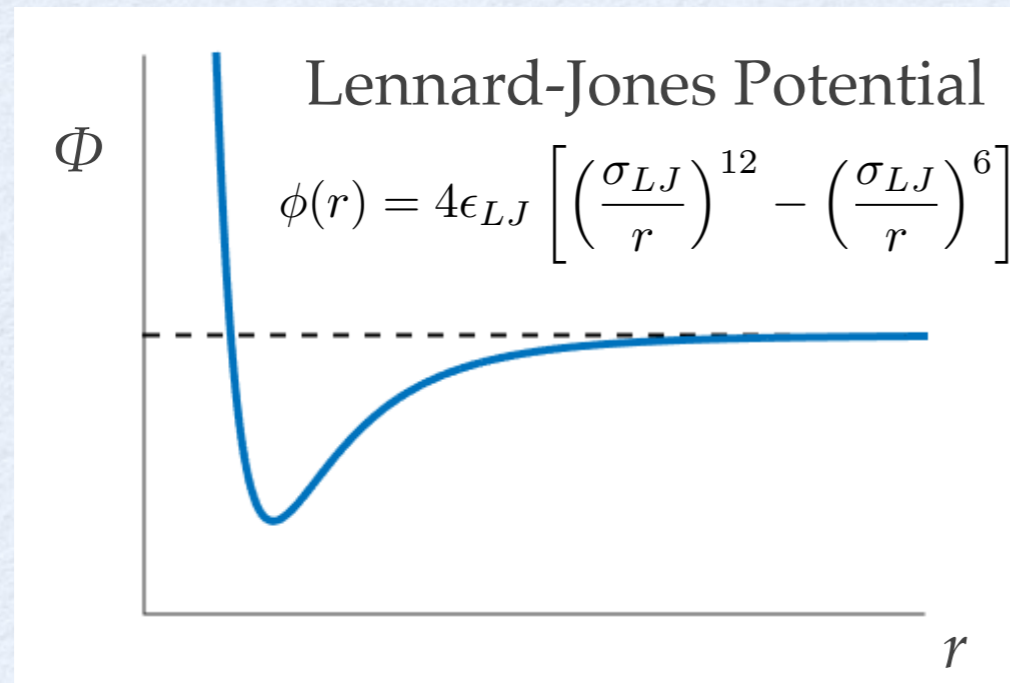
without Hierarchical Data Structure



with Hierarchical Data Structure



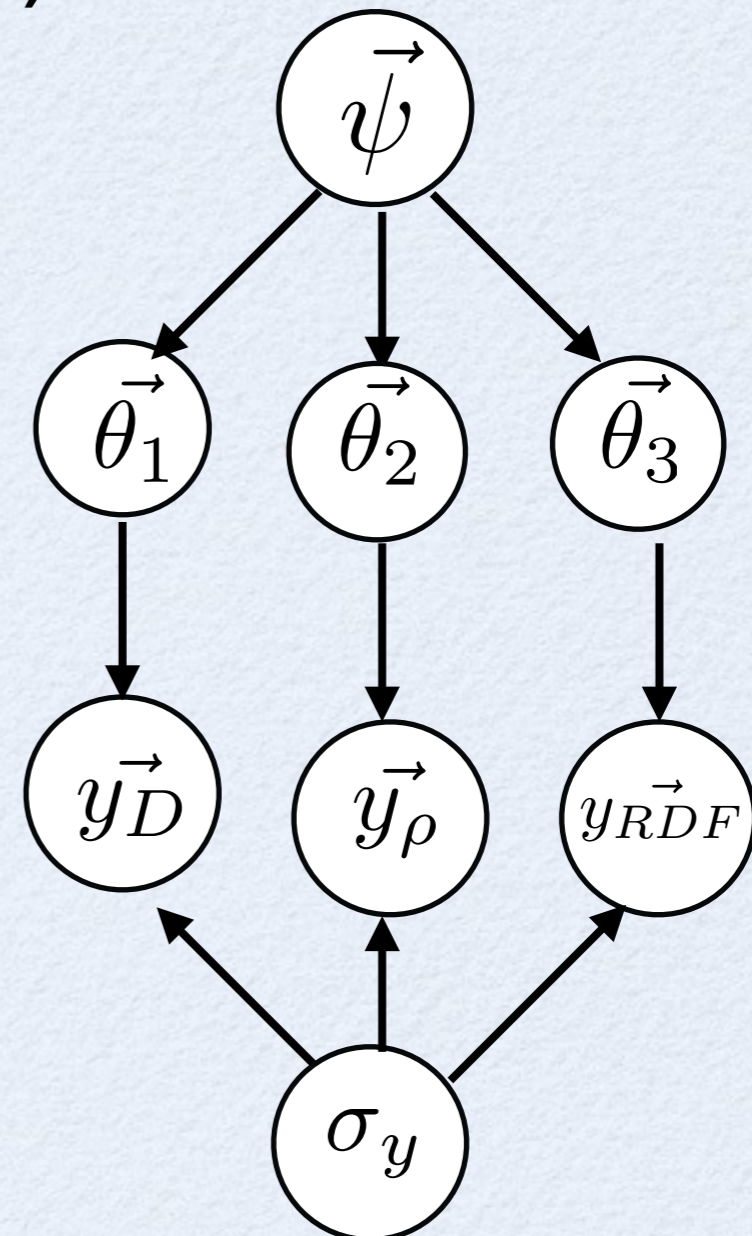
EXAMPLE: MD Simulations for Water



- MD parameter calibration - θ
 - O-O Lennard-Jones (ϵ_{LJ} , σ_{LJ})
 - O-O, O-H Coulomb charges (q)

DATA for Water

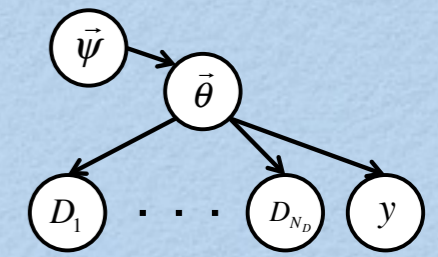
- Heterogeneous Data Sets
(as a function of temperature)
- Diffusion coefficient
- Density
- Radial Distribution Function (RDF)



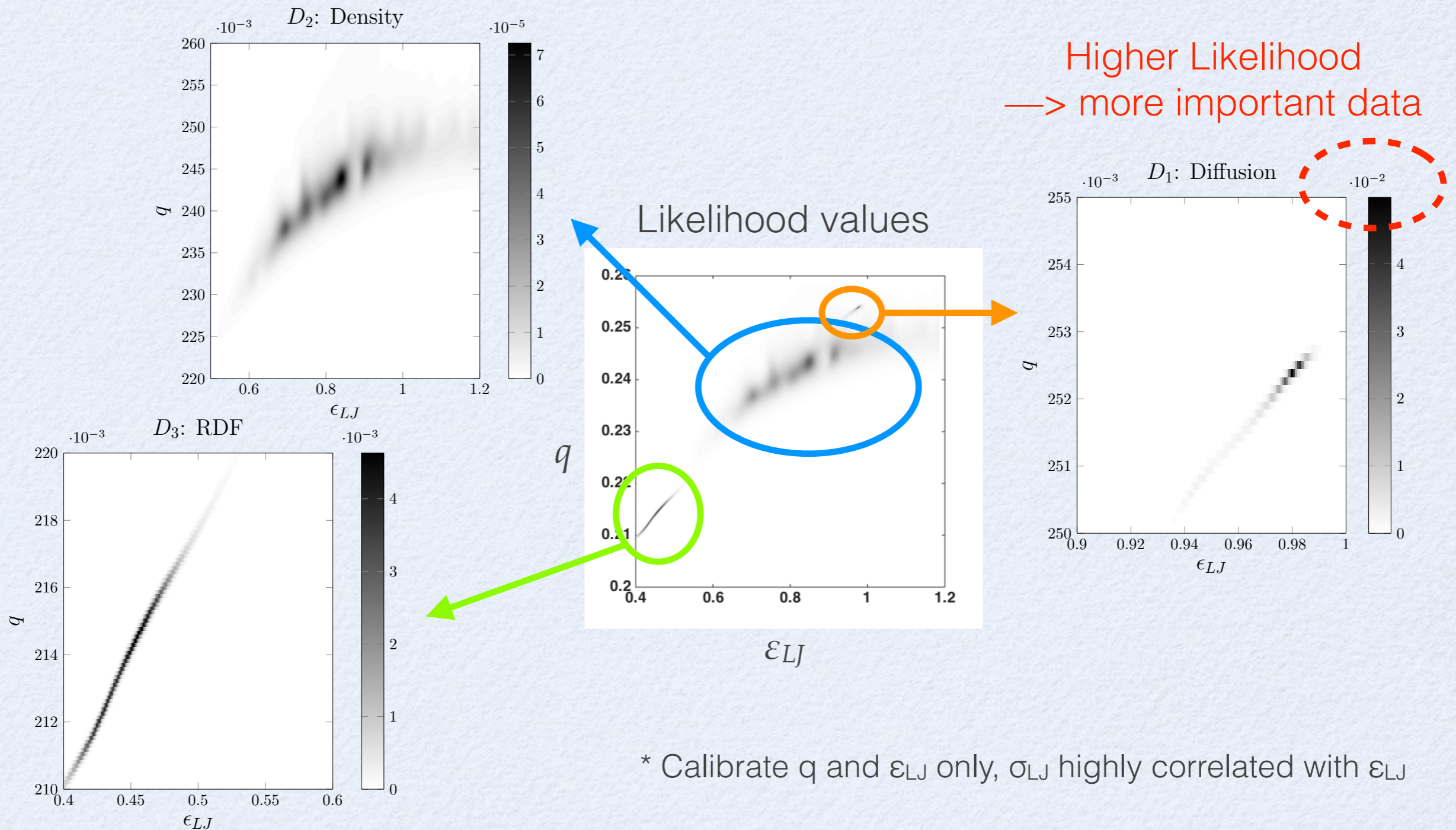
Data sources:

Holz et al. 2000; Jones & Harris 1992; Soper 2013

Bayesian Inference for pooled data



Calibrate for each data set individually...



* Calibrate q and ϵ_{LJ} only, σ_{LJ} highly correlated with ϵ_{LJ}

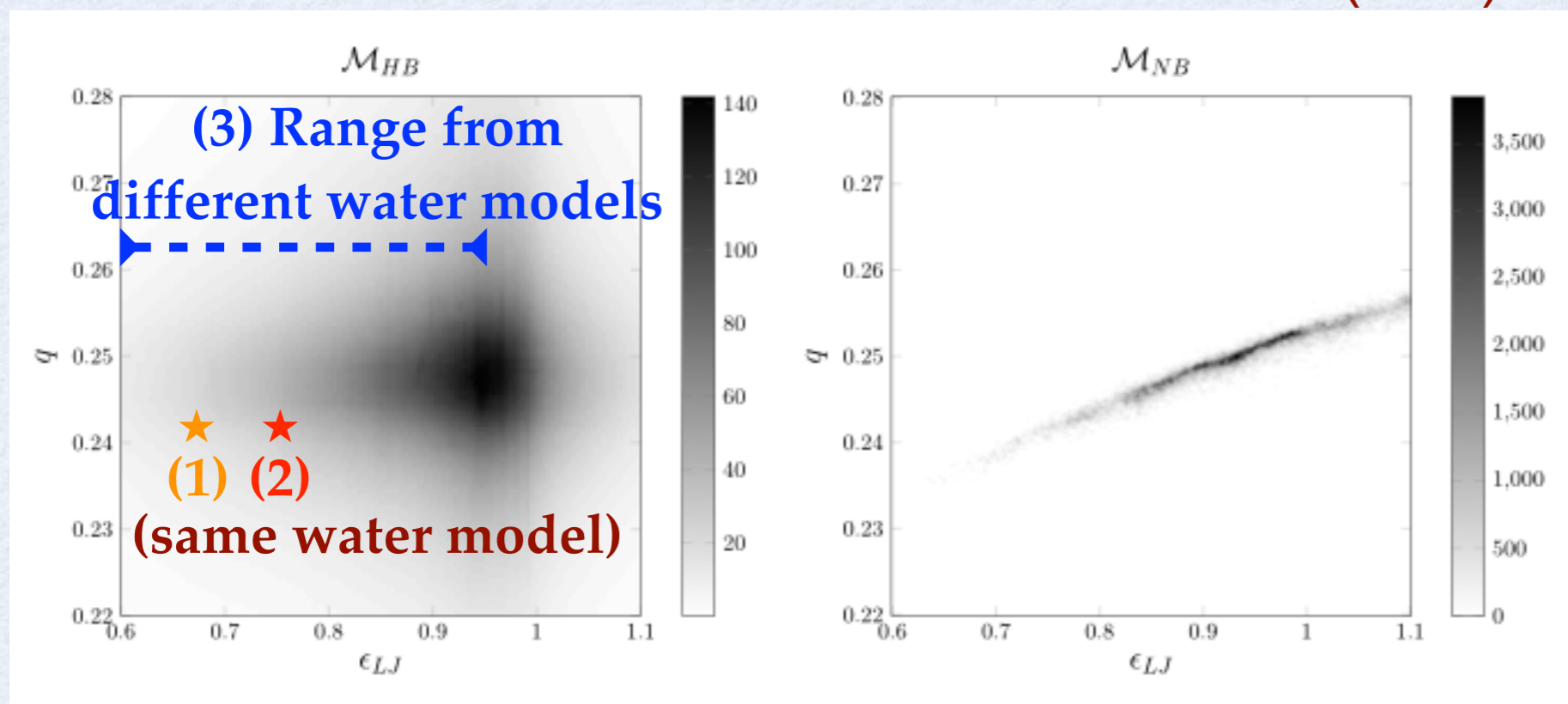
Posterior of Model Parameters

Compare Hierarchical Bayesian model with Independent Likelihoods model

Hierarchical (HB)

Independent Likelihoods (NB)

Model Selection



Model P(Model|D)

HB 0.98

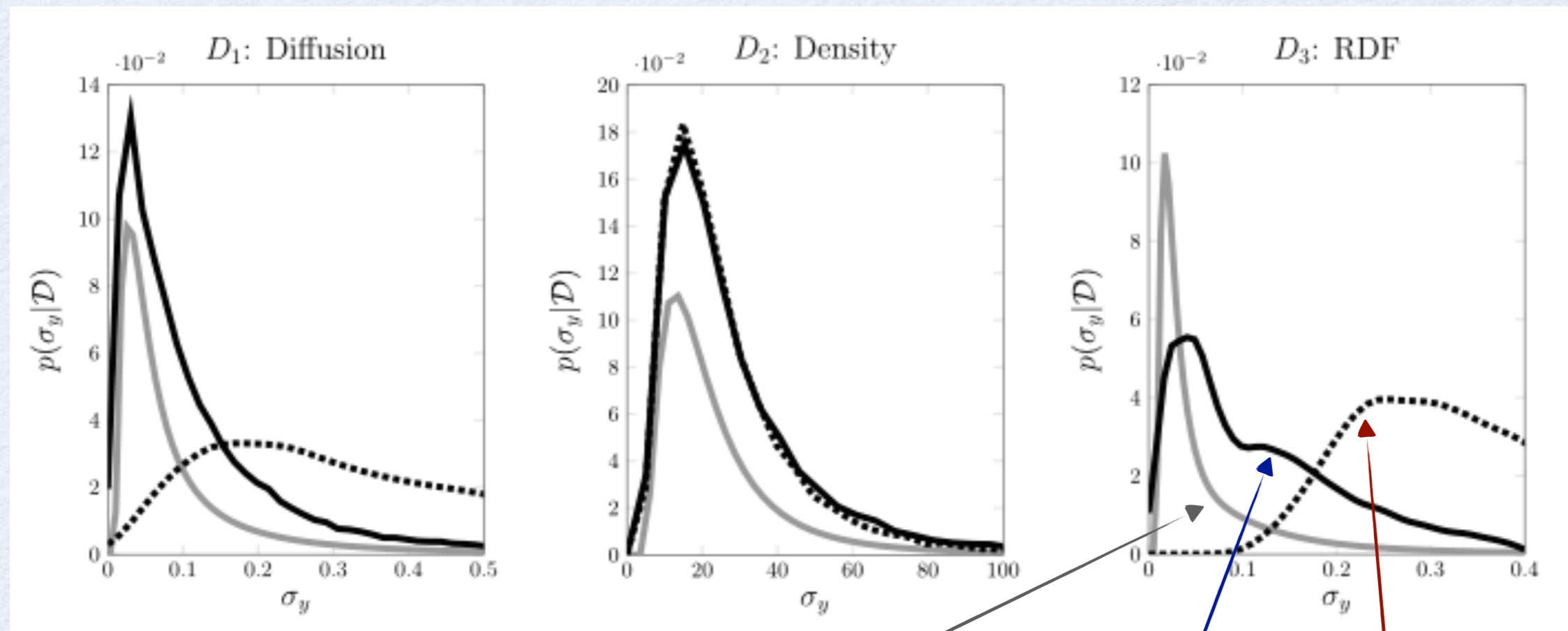
NB 0.02

(1) Mahoney & Jorgensen 2000 (2) Rick et al. 2004

(3) http://www1.lsbu.ac.uk/water/water_models.html

Posterior of Measurement Noise

Results from Hierarchical Bayesian model overlaps with results from individual data set



single data set

Hierarchical

Independent Likelihoods

SUMMARY

- Bayesian UQ in MD: **Data Structure** and **Computational Intensity**
 - **Data Structuring**
 - Calibration data for MD are heterogeneous →
 - Assuming independence for all data is NOT ENOUGH to reach good calibration and prediction!
 - Combine heterogeneous data based on Evidence
 - balance data-fitting and information gain
 - **Computational intensity**
 - “Evolutionary” algorithms (BASIS/SMC)
 - HPC framework for Bayesian UQ (Pi4U)
 - Surrogate models

Thank you!

Some References:

- Wu, S., Angelikopoulos, P., Moser, R., Papadimitriou, C. and Koumoutsakos, P. (2015) "A Hierarchical Bayesian Framework for Force Field Selection in Molecular Dynamics Simulations", *Philosophical Transactions of the Royal Society B*.
- Angelikopoulos, P., Papadimitriou, C. and Koumoutsakos, P. (2015) "X-TMCMC: Adaptive Kriging for Bayesian inverse modeling", *Computer Methods in Applied Mechanics and Engineering*
- Hadjidoukas, P.E., Angelikopoulos, P., Papadimitriou, C. and Koumoutsakos, P. (2015) "Pi4U: A high performance computing framework for Bayesian uncertainty quantification of complex models", *Journal of Computational Physics*
- Hadjidoukas, P.E., Angelikopoulos, P., Rossinelli, D., Alexeev, D., Papadimitriou, C. and Koumoutsakos, P. (2014) "Bayesian Uncertainty Quantification and propagation for Discrete Element Simulations of granular materials", *Computer Methods in Applied Mechanics and Engineering*
- Angelikopoulos, P., Papadimitriou, C. and Koumoutsakos, P. (2013) "Data Driven, Predictive Molecular Dynamics for Nanoscale Flow Simulations Under Uncertainty", *Journal of Physical Chemistry B*
- Angelikopoulos, P., Papadimitriou, C. and Koumoutsakos, P. (2012) "Bayesian uncertainty quantification and propagation in molecular dynamics simulations: A high performance computing framework", *Journal of Chemical Physics*