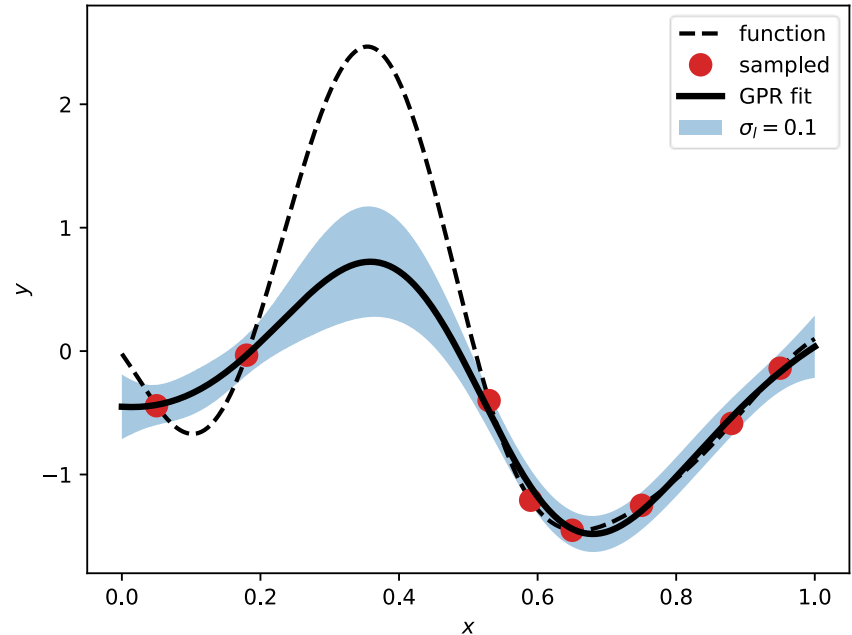


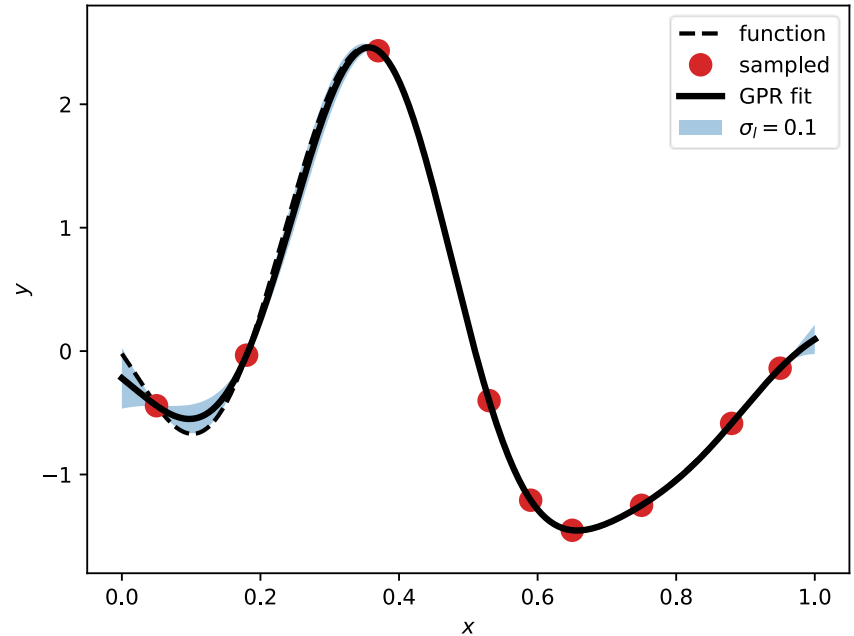
Trustworthy machine learning for materials

Federico Grasselli – COSMO Lab
EPFL

- Pillar of Scientific Method
- ML statistical nature
- Sources of uncertainty
- So far scarcely employed:
 - lack of standards
 - large cost (training and evaluation)
 - ad hoc training:
(MC dropout, deep and shallow ensembles, Gaussian mixture models, committees)

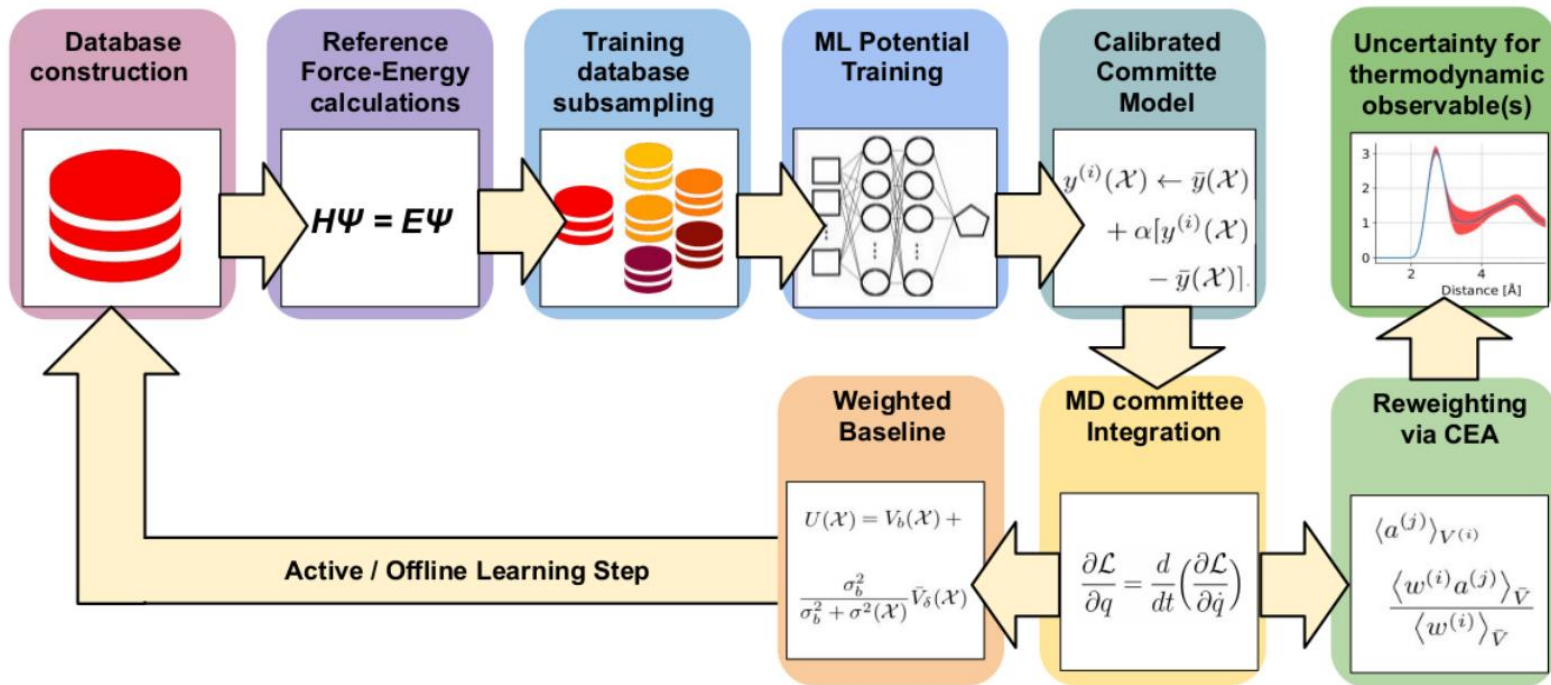


- Pillar of Scientific Method
- ML statistical nature
- Sources of uncertainty
- So far scarcely employed:
 - lack of standards
 - large cost (training and evaluation)
 - ad hoc training:
(MC dropout, deep and shallow ensembles, Gaussian mixture models, committees)



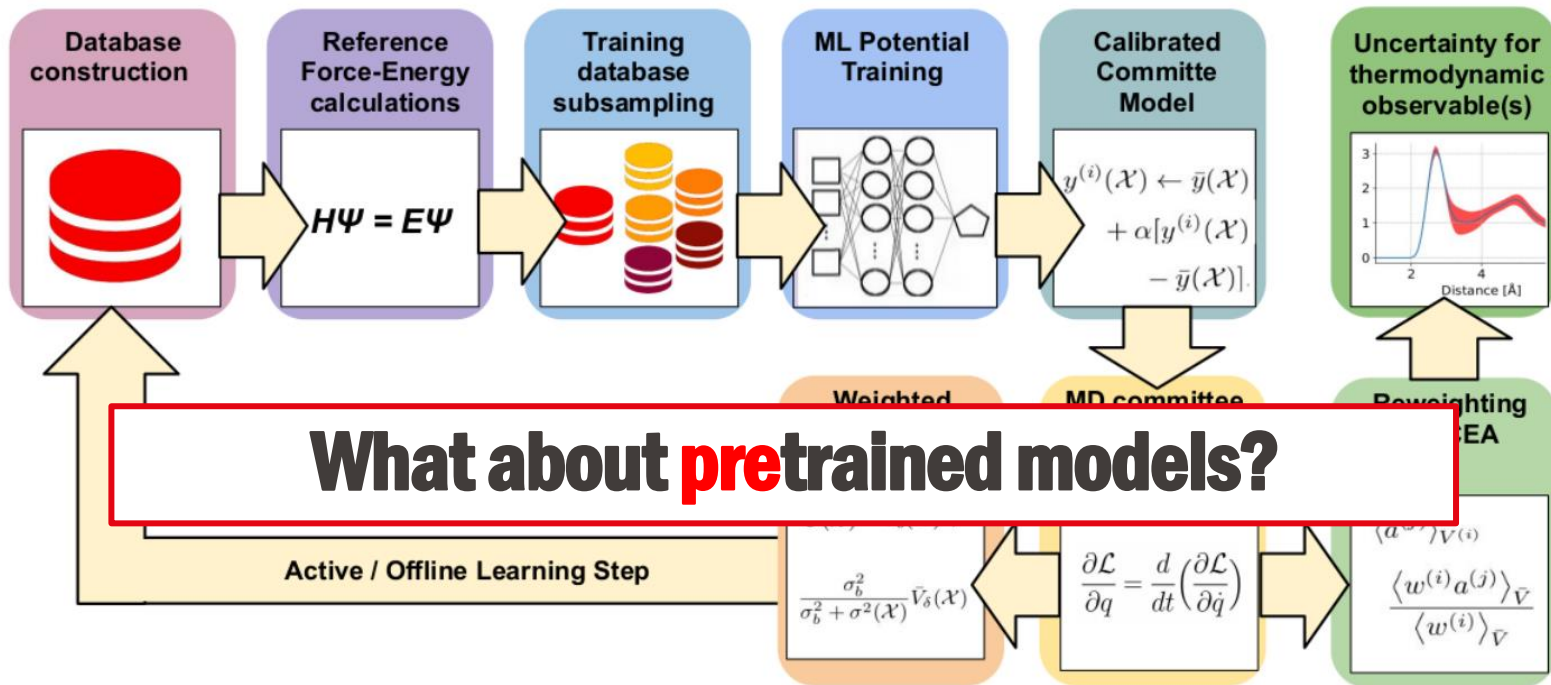
Use of UQ in atomistic simulations

- Uncertainty aware simulations
- Reweight via committee of models



Use of UQ in atomistic simulations

- Uncertainty aware simulations
- Reweight via committee of models



Prediction Rigidity (PR)

- Loss:

$$\mathcal{L}(\mathbf{w}|D) = \sum_{i=1}^{N_{\text{train}}} \ell[\tilde{y}_i(\mathbf{x}_i, \mathbf{w}), y_i]$$

- How “rigid” is the prediction for a given input \mathbf{x}_* ?

Prediction Rigidity (PR)

- Loss:

$$\mathcal{L}(\mathbf{w}|D) = \sum_{i=1}^{N_{\text{train}}} \ell[\tilde{y}_i(\mathbf{x}_i, \mathbf{w}), y_i]$$

- How “rigid” is the prediction for a given input \mathbf{x}_* ?
- Constrained minimization of $\mathcal{L}_c = \mathcal{L} + \lambda(\epsilon_* - \tilde{y}(\mathbf{x}_*, \mathbf{w}))$

Prediction Rigidity (PR)

- Loss:

$$\mathcal{L}(\mathbf{w}|D) = \sum_{i=1}^{N_{\text{train}}} \ell[\tilde{y}_i(\mathbf{x}_i, \mathbf{w}), y_i]$$

- How “rigid” is the prediction for a given input \mathbf{x}_* ?
- Constrained minimization of $\mathcal{L}_c = \mathcal{L} + \lambda(\epsilon_* - \tilde{y}(\mathbf{x}_*, \mathbf{w}))$
- For $\epsilon_* \approx \tilde{y}(\mathbf{x}_*, \mathbf{w}_o)$ we have

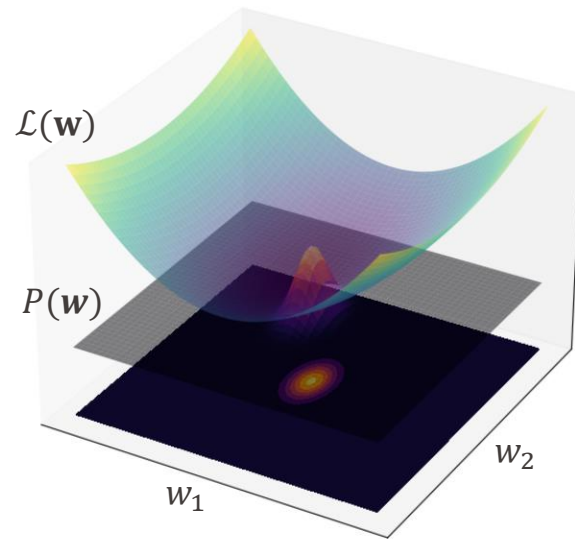
$$\mathcal{L}_c(\epsilon_*|D) \approx \mathcal{L}(\mathbf{w}_o|D) + \frac{1}{2} R_* (\epsilon_* - \tilde{y}(\mathbf{x}_*, \mathbf{w}_o))^2$$

- with

$$R_* \equiv \left. \frac{\partial^2 \mathcal{L}}{\partial \epsilon_*^2} \right|_{\epsilon_* = \tilde{y}(\mathbf{x}_*, \mathbf{w}_o)} = \left[\left. \frac{\partial \tilde{y}_*}{\partial \mathbf{w}} \right|_{\mathbf{w}_o} \left[\left. \frac{\partial^2 \mathcal{L}}{\partial \mathbf{w}^2} \right|_{\mathbf{w}_o} \right]^{-1} \left. \frac{\partial \tilde{y}_*}{\partial \mathbf{w}} \right|_{\mathbf{w}_o} \right]^{-1}$$

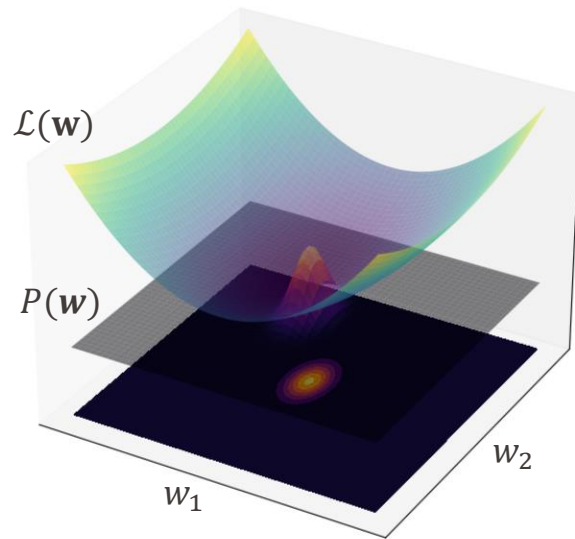
Prediction Rigidity (PR)

- Fitted models: $P(\mathbf{w}) \propto \exp(-\mathcal{L}(\mathbf{w}))$
- Laplace approximation:
 - 2nd-order approximation of the loss
 - Gaussian approx. of the probability density



Prediction Rigidity (PR)

- Fitted models: $P(\mathbf{w}) \propto \exp(-\mathcal{L}(\mathbf{w}))$
- Laplace approximation:
 - 2nd-order approximation of the loss
 - Gaussian approx. of the probability density



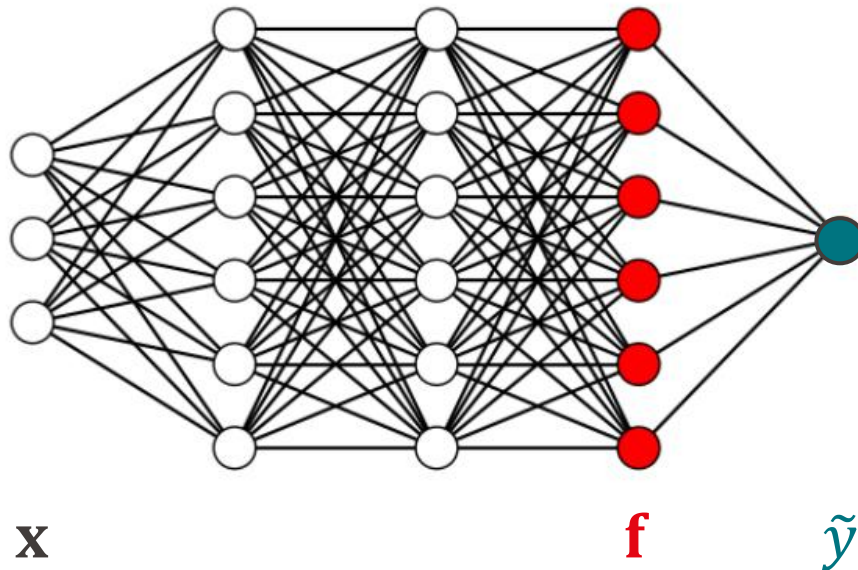
- Linear regression: $\tilde{y} = \mathbf{w} \cdot \mathbf{x} \rightarrow \frac{\partial \tilde{y}}{\partial \mathbf{w}} \equiv \mathbf{x}$
- PR is simple:

$$R_{\star} = [\mathbf{x}_{\star} \cdot [\mathbf{X}^T \mathbf{X}]^{-1} \mathbf{x}_{\star}]^{-1}$$

- $\frac{1}{R_{\star}}$ has got the same shape of variance in Gaussian process regression

What about Neural Networks?

- Statistical theory of NNs. Training and over-parametrization
- Central Limit Theorem → infinitely wide NNs as Gaussian processes
- Last-layer features \mathbf{f}



[1] Lee et al. arXiv:1711.00165 (2017)

[2] Jacot et al, NIPS (2018)

[3] Lee et al, NIPS (2019)

Last-Layer Prediction Rigidity

- Uncertainty simplifies to (the inverse of) **Last-Layer PR**:

The diagram shows the equation for predictive uncertainty, $\sigma_{\star}^2 = \alpha^2 \mathbf{f}_{\star}^{\top} (\mathbf{F}^{\top} \mathbf{F} + \zeta^2 \mathbf{I})^{-1} \mathbf{f}_{\star}$, with red arrows pointing to various components and their labels:

- predictive uncertainty** points to the left side of the equation, σ_{\star}^2 .
- scale** points to the coefficient α^2 .
- last-layer features for the training set** points to the matrix $\mathbf{F}^{\top} \mathbf{F}$.
- regularization** points to the matrix $\zeta^2 \mathbf{I}$.
- last-layer features for the prediction** points to the vector \mathbf{f}_{\star} on the right side of the equation.

Last-Layer Prediction Rigidity

- Uncertainty simplifies to (the inverse of) **Last-Layer PR**:

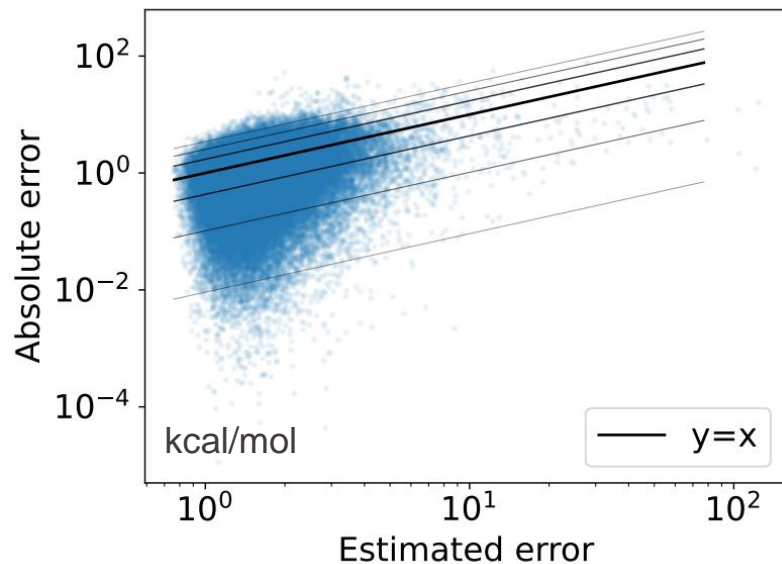
The diagram shows the equation for predictive uncertainty: $\sigma_{\star}^2 = \alpha^2 \mathbf{f}_{\star}^{\top} (\mathbf{F}^{\top} \mathbf{F} + \zeta^2 \mathbf{I})^{-1} \mathbf{f}_{\star}$. Red arrows point from text labels to parts of the equation: 'predictive uncertainty' points to the left side; 'scale' points to α^2 ; 'last-layer features for the training set' points to $\mathbf{F}^{\top} \mathbf{F}$; 'regularization' points to $\zeta^2 \mathbf{I}$; and 'last-layer features for the prediction' points to \mathbf{f}_{\star} .

$$\sigma_{\star}^2 = \alpha^2 \mathbf{f}_{\star}^{\top} (\mathbf{F}^{\top} \mathbf{F} + \zeta^2 \mathbf{I})^{-1} \mathbf{f}_{\star}$$

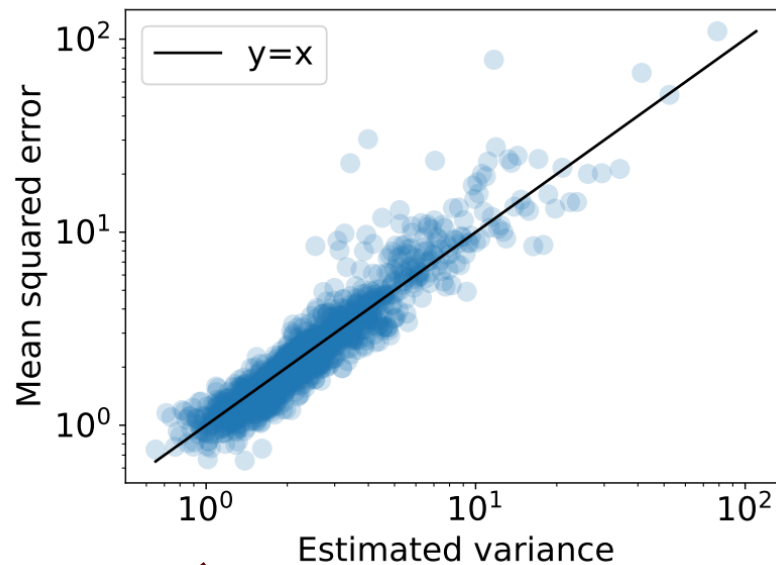
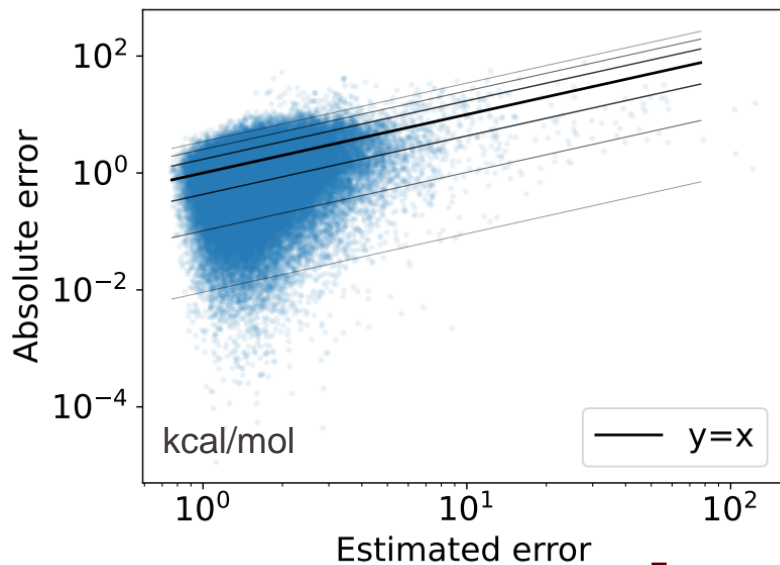
- Effectively the uncertainty of a linear model
- Very easy to calculate, two hyperparameters must be tuned
- Doesn't depend on target values of the training set

Last-Layer PR: results

QM9 dataset



Last-Layer PR: results QM9 dataset

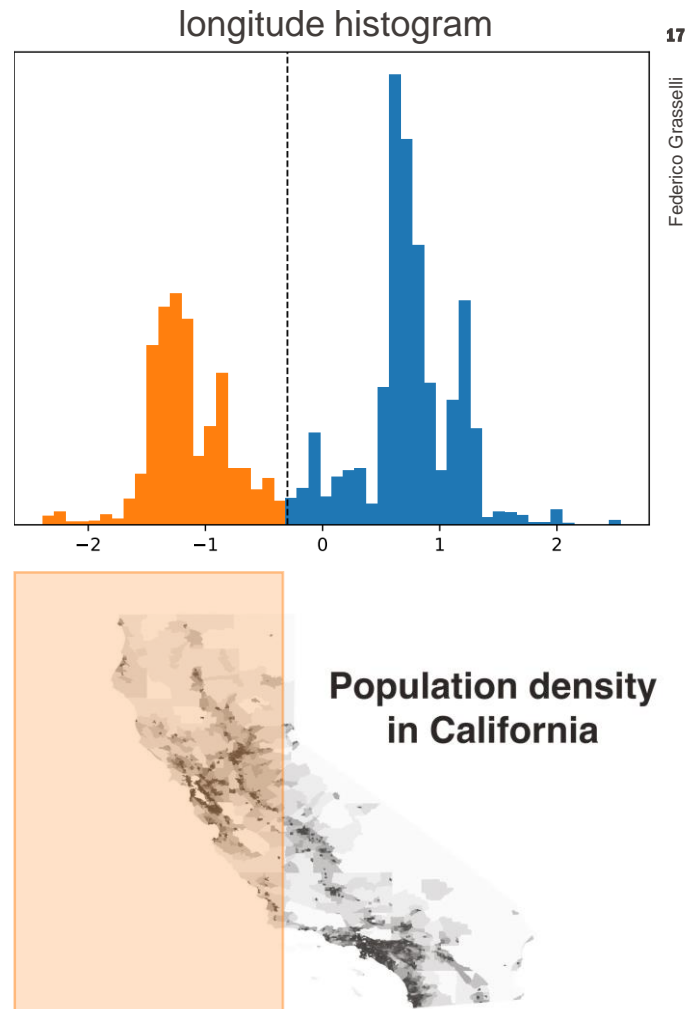


binning

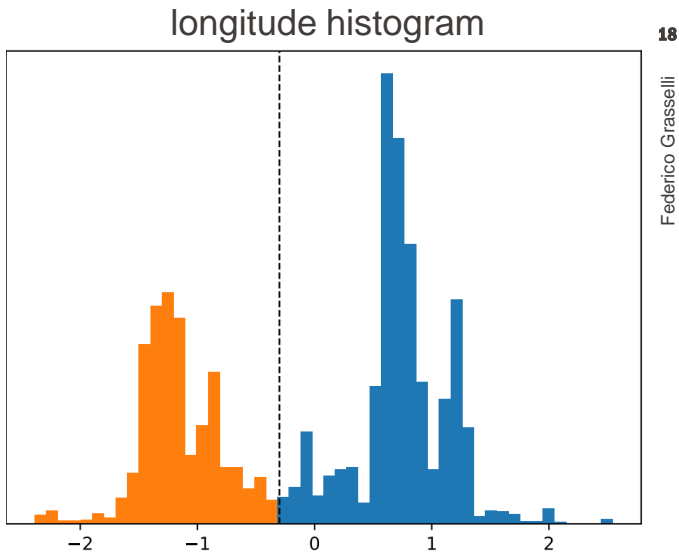
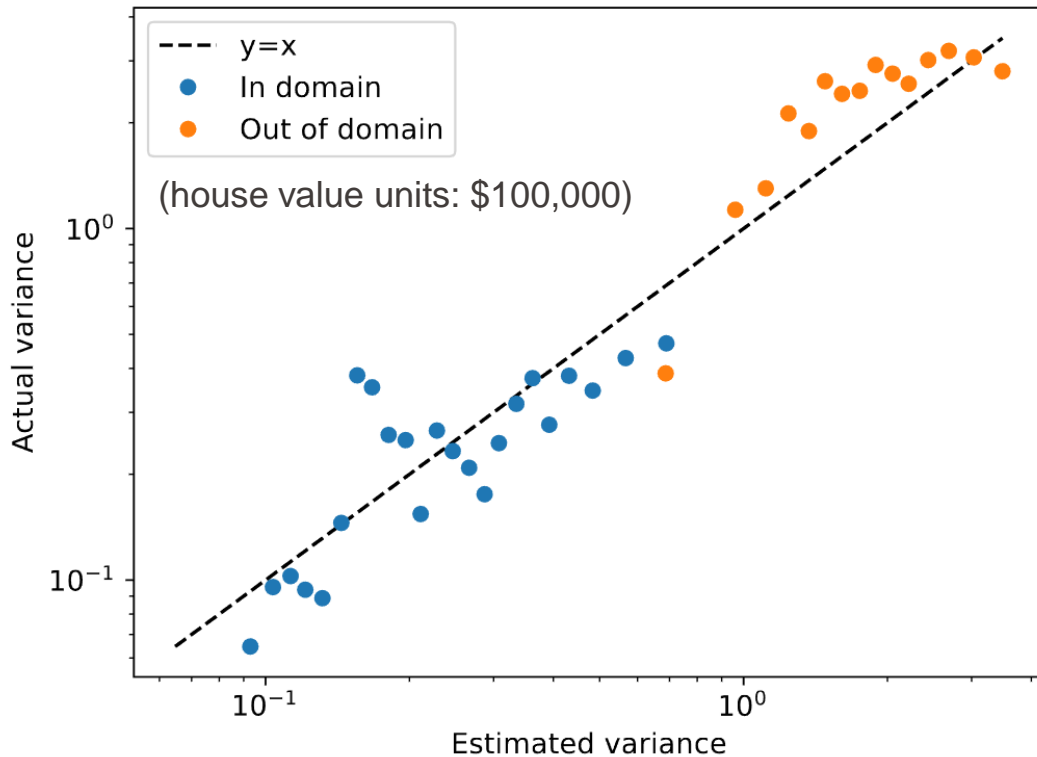
Last-Layer PR results: California housing \$



Last-Layer PR results: California housing \$

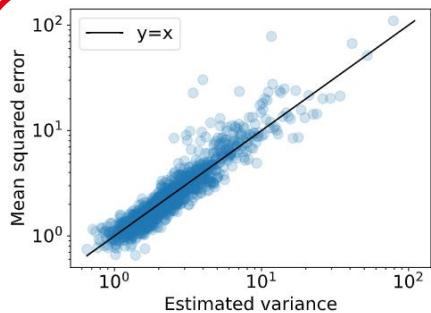


Last-Layer PR results: California housing \$

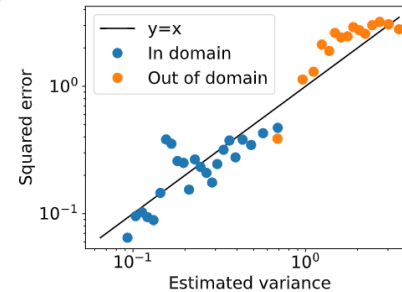


Last-Layer PR results

- Very good uncertainty estimates across a wide variety of problems



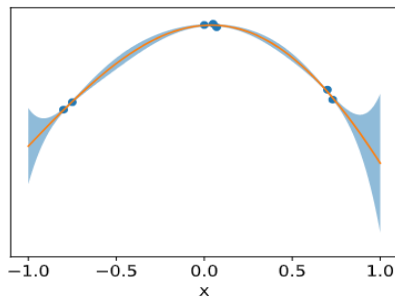
Chemistry



Housing prices

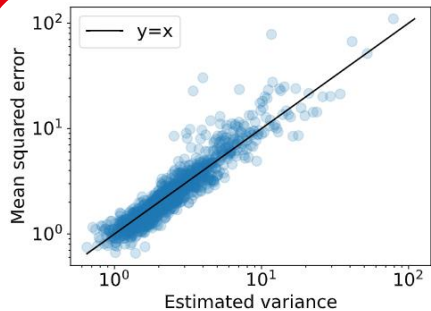
- Very good uncertainty estimates across a wide variety of problems

A toy model

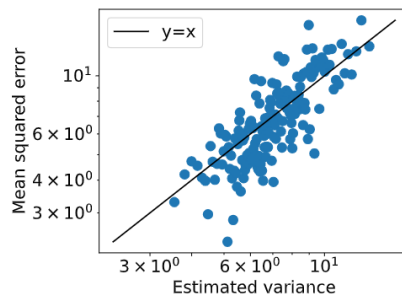


Dataset	RMSE				NLL			
	PBP	MCD	DE	LLPR	PBP	MCD	DE	LLPR
Concrete	5.67 _{0.00}	5.23 _{0.02}	6.03 _{0.13}	5.26 _{0.25}	3.16 _{0.02}	3.04 _{0.02}	3.06 _{0.00}	3.09 _{0.07}
Energy	1.80 _{0.00}	1.66 _{0.00}	2.09 _{0.06}	0.49 _{0.03}	2.04 _{0.02}	1.99 _{0.02}	1.38 _{0.00}	0.69 _{0.07}
Kin8nm	0.10 _{0.00}	0.10 _{0.00}	0.09 _{0.00}	0.08 _{0.00}	-0.90 _{0.00}	-0.95 _{0.00}	-1.20 _{0.00}	-1.12 _{0.00}
Naval	0.01 _{0.00}	0.01 _{0.00}	0.00 _{0.00}	0.00 _{0.00}	-3.73 _{0.00}	-3.80 _{0.00}	-5.63 _{0.00}	-7.07 _{0.00}
Power	4.12 _{0.00}	4.02 _{0.00}	4.11 _{0.00}	3.94 _{0.07}	2.84 _{0.00}	2.80 _{0.00}	2.79 _{0.00}	2.83 _{0.00}
Protein	4.73 _{0.00}	4.36 _{0.02}	4.71 _{0.00}	4.18 _{0.02}	2.97 _{0.00}	2.89 _{0.00}	2.83 _{0.00}	2.91 _{0.00}
Wine	0.64 _{0.00}	0.62 _{0.00}	0.64 _{0.00}	0.63 _{0.02}	0.97 _{0.00}	0.93 _{0.00}	0.94 _{0.00}	1.02 _{0.00}
Yacht	1.02 _{0.00}	1.11 _{0.00}	1.58 _{0.11}	1.19 _{0.04}	1.63 _{0.02}	1.55 _{0.00}	1.18 _{0.00}	1.58 _{0.20}
Year	8.88 _{0.00}	8.86 _{0.00}	8.89 _{0.00}	8.91 _{0.00}	3.60 _{0.00}	3.59 _{0.00}	3.35 _{0.00}	3.61 _{0.00}

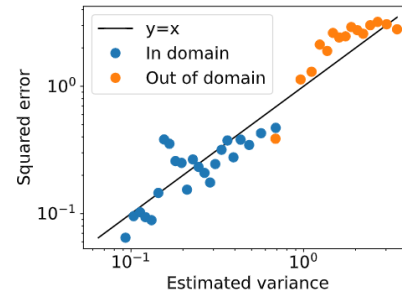
A standard benchmark



Chemistry



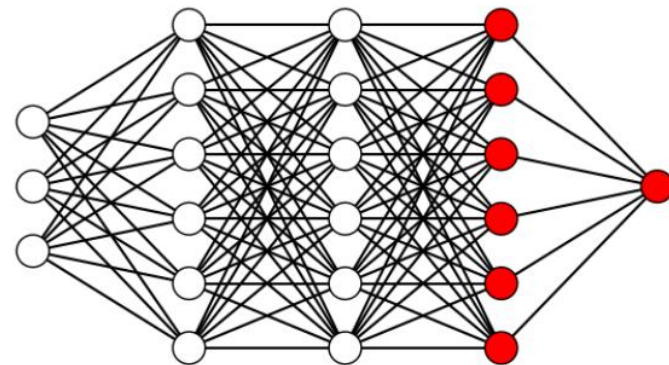
Weather prediction



Housing prices

Last-Layer PR: Summary

- Cheap, practical, scalable, a-posteriori
- Explains the success of last-layer approximations
- Pre-print on arxiv¹, code available @ <https://github.com/frostedoyster/llpr>



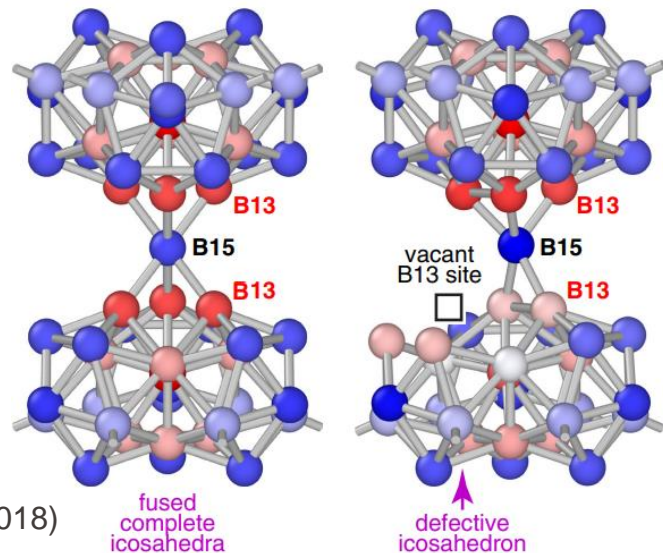
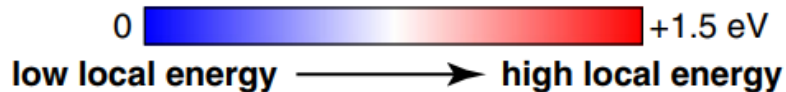
LLPR

¹ Bigi, Chong, Ceriotti & Grasselli, arXiv:2403.02251 (2024)

PR for local predictions

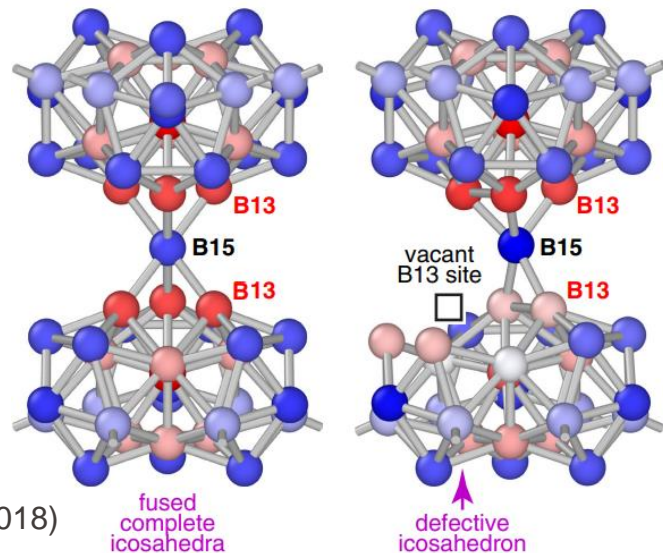
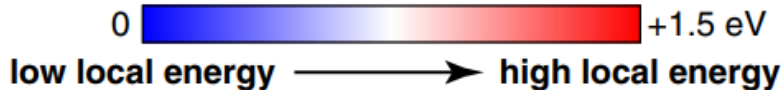
- Atomistic models: local energies are **not** observables
- Yet used in
 - constructing ML models $E(A) = \sum_{i \in A} E_i$
 - heuristic analyses

- Atomistic models: local energies are **not** observables
- Yet used in
 - constructing ML models $E(A) = \sum_{i \in A} E_i$
 - heuristic analyses



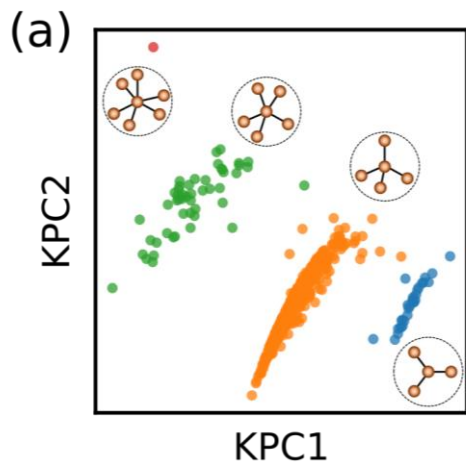
Adapted from Deringer, Pickard, Csányi. PRL 120 156001 (2018)

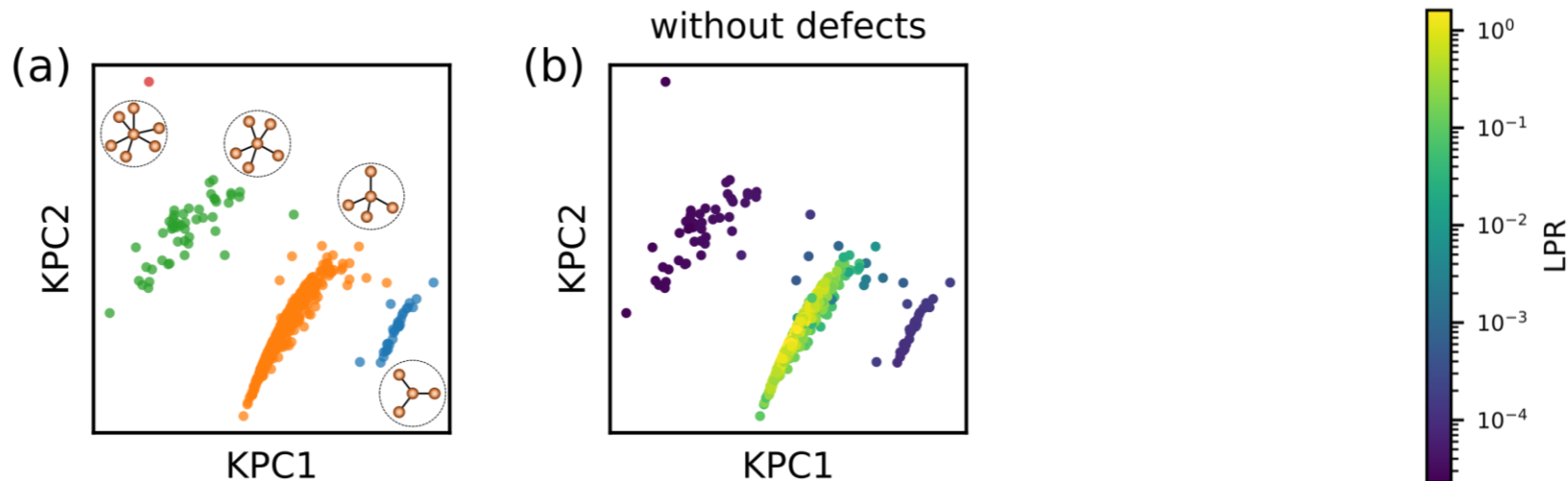
- Atomistic models: local energies are **not** observables
- Yet used in
 - constructing ML models $E(A) = \sum_{i \in A} E_i$
 - heuristic analyses



Adapted from Deringer, Pickard, Csányi. PRL 120 156001 (2018)

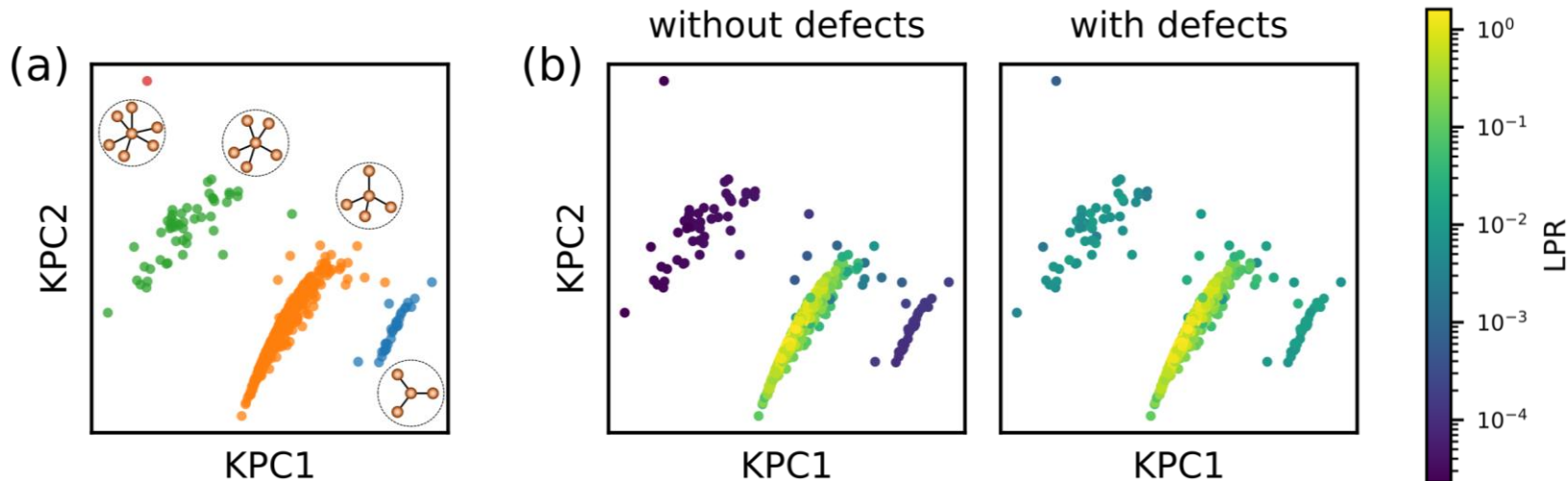
- How robust are these? Use the same formulation of constrained loss, but now for *local predictions*



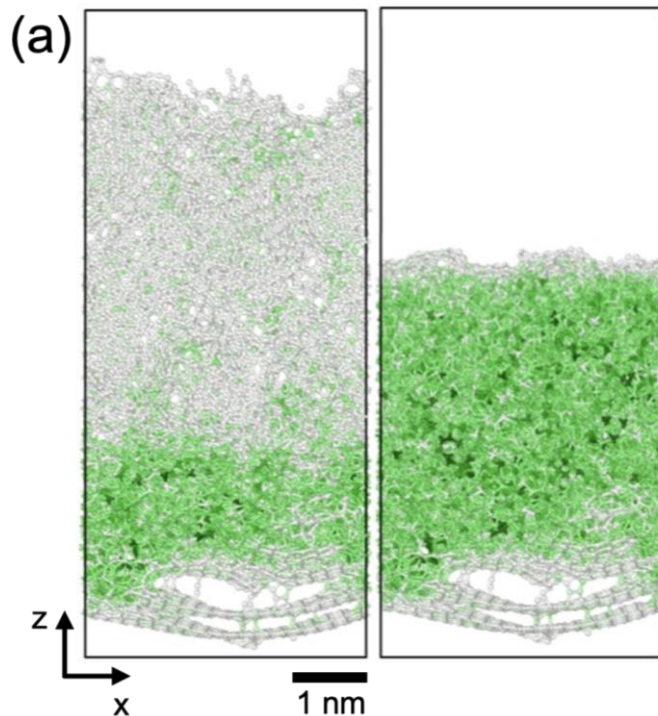


Local PR: amorphous silicon

- Dramatic increase of local PR by adding structures containing under/over coordinated environments



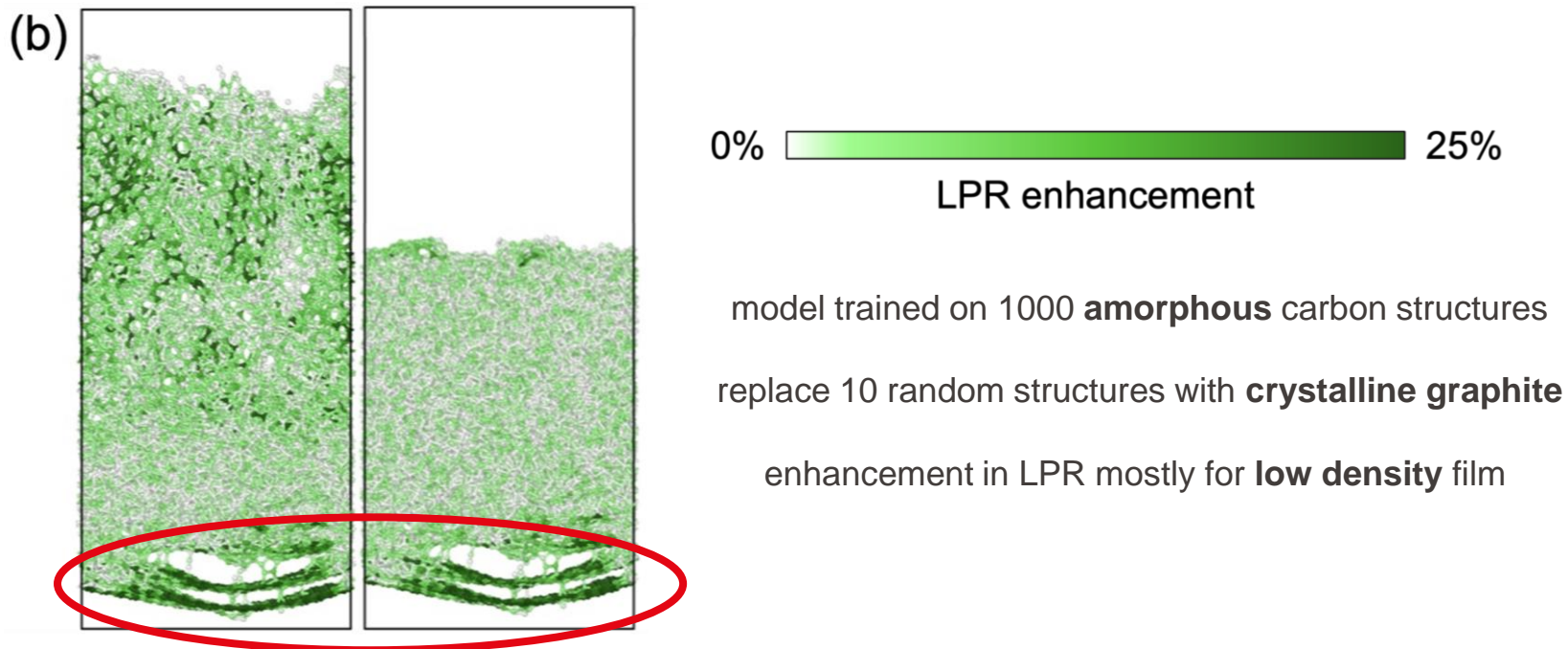
- Selective increase of local PR in low/high density carbon films



0%  25%
LPR enhancement

model trained on 1000 **amorphous** carbon structures
replace 10 random structures with **crystalline diamond**
enhancement in LPR mostly for **high density** film

- Selective increase of local PR in low/high density carbon films

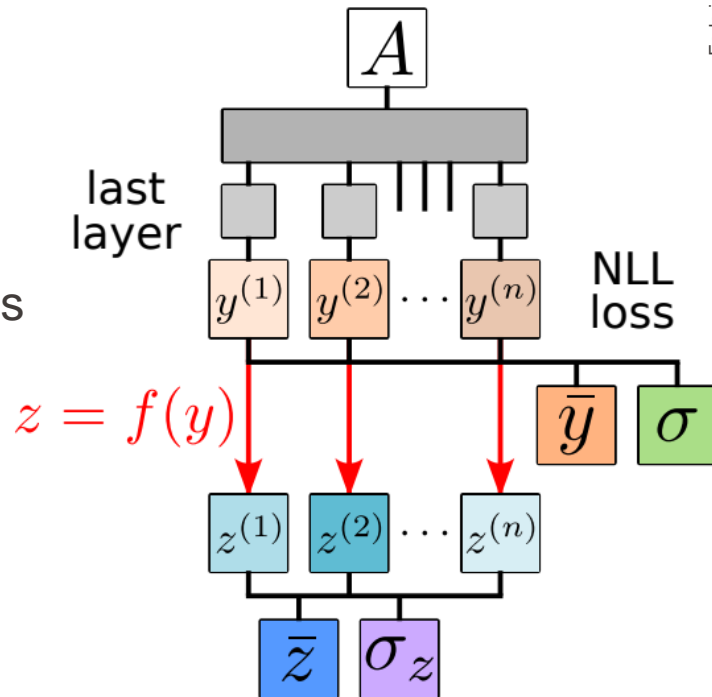


Remarks, Conclusions and Outlook

- Low-cost uncertainty on new predictions for pretrained models
- Constrained loss minimization is just a theoretical tool:
 - no need to train a model with constrained loss to get PR and uncertainties
 - no need for target values
- Rigidity of local predictions is readily obtained with same formalism

Remarks, Conclusions and Outlook

- Sample last-layer weights according to Laplace approximation \rightarrow ensemble
- Propagate uncertainty to derived quantities
- Use it on thermodynamic observables



Kellner & Ceriotti, <https://arxiv.org/abs/2402.16621> (2024)

Acknowledgements



Sanggyu Chong, EPFL



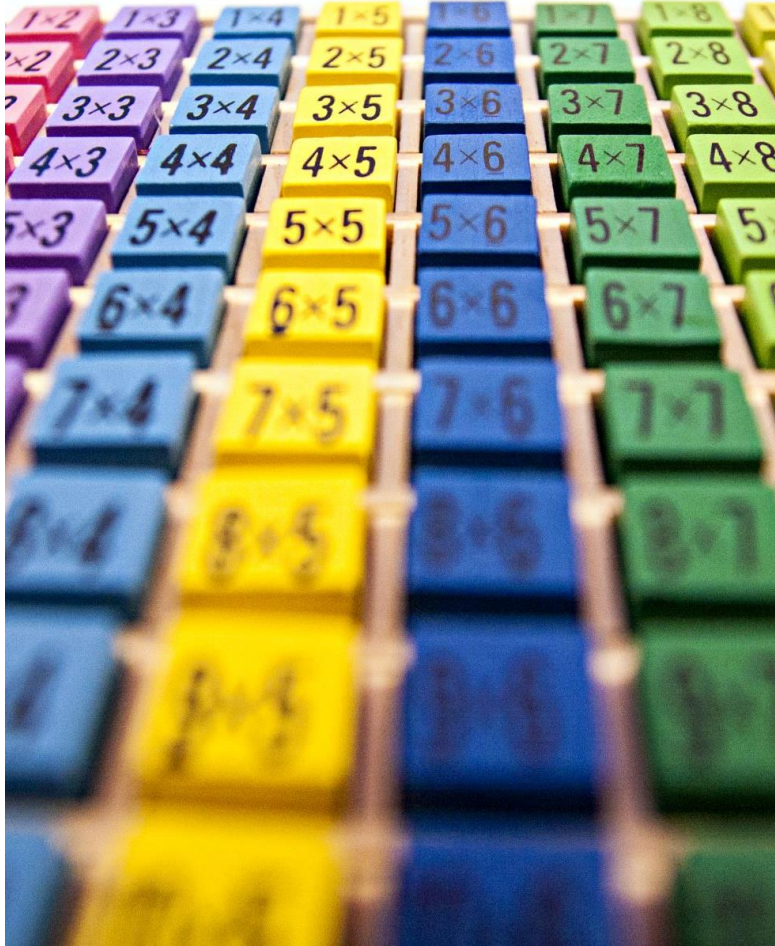
Filippo Bigi, EPFL



Chiheb Ben Mahmoud, EPFL,
now Oxford University



Michele Ceriotti, EPFL

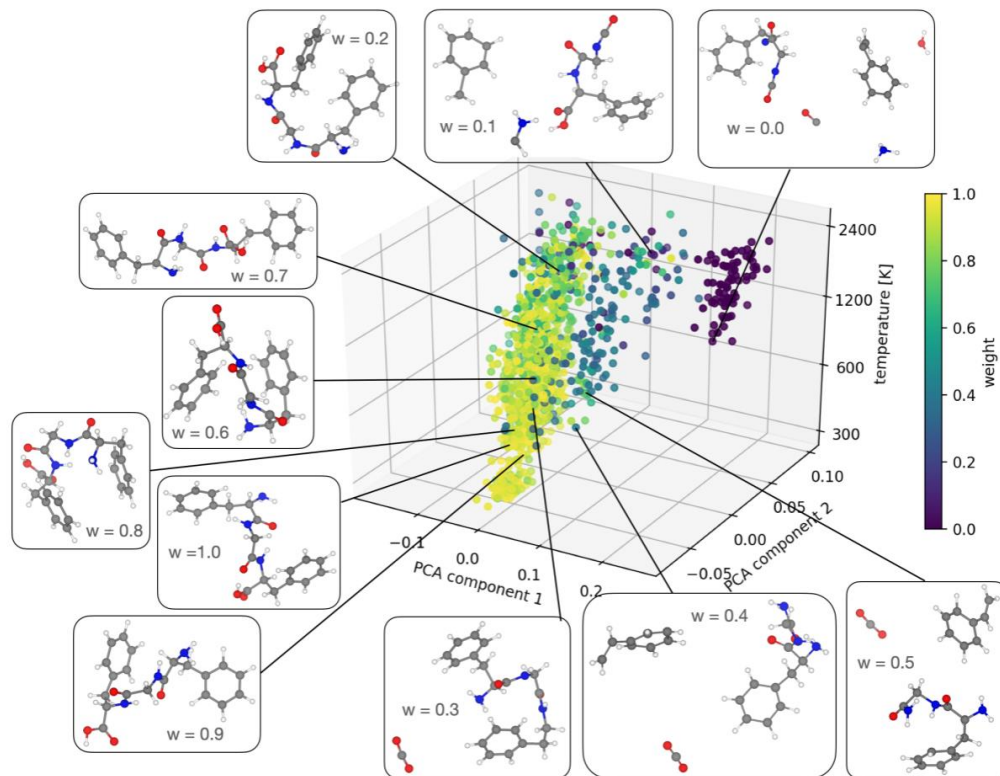


Backup slides

Federico Grasselli – COSMO Lab
EPFL

Use of UQ in atomistic simulations

- Uncertainty aware simulations



- Statistical theory of NNs. Training and over-parametrization
- Central Limit Theorem → infinitely wide NNs as Gaussian processes
- Two deterministic kernels
 - Neural-Network Gaussian Process (NNGP)¹: initialization
 - Neural Tangent Kernel (NTK)²: training
- Distribution of predictions³ during training is Gaussian
- Evolution of mean and variance is deterministic!

$$\mu_{\star} = \mathbf{k}_{\text{NTK}}(\star, \mathcal{D}) \mathbf{K}_{\text{NTK}}^{-1} (\mathbf{I} - e^{-\eta \mathbf{K}_{\text{NTK}} t}) \mathbf{y}$$

$$\sigma_{\star}^2 = k_{\text{NNGP}}(\star, \star)$$

$$+ \mathbf{k}_{\text{NTK}}(\star, \mathcal{D}) \mathbf{K}_{\text{NTK}}^{-1} (\mathbf{I} - e^{-\eta \mathbf{K}_{\text{NTK}} t}) \mathbf{K}_{\text{NNGP}} (\mathbf{I} - e^{-\eta \mathbf{K}_{\text{NTK}} t}) \mathbf{K}_{\text{NTK}}^{-1} \mathbf{k}_{\text{NTK}}(\mathcal{D}, \star)$$

$$- \mathbf{k}_{\text{NTK}}(\star, \mathcal{D}) \mathbf{K}_{\text{NTK}}^{-1} (\mathbf{I} - e^{-\eta \mathbf{K}_{\text{NTK}} t}) \mathbf{k}_{\text{NNGP}}(\mathcal{D}, \star)$$

$$- \mathbf{k}_{\text{NNGP}}(\star, \mathcal{D}) (\mathbf{I} - e^{-\eta \mathbf{K}_{\text{NTK}} t}) \mathbf{K}_{\text{NTK}}^{-1} \mathbf{k}_{\text{NTK}}(\mathcal{D}, \star)$$

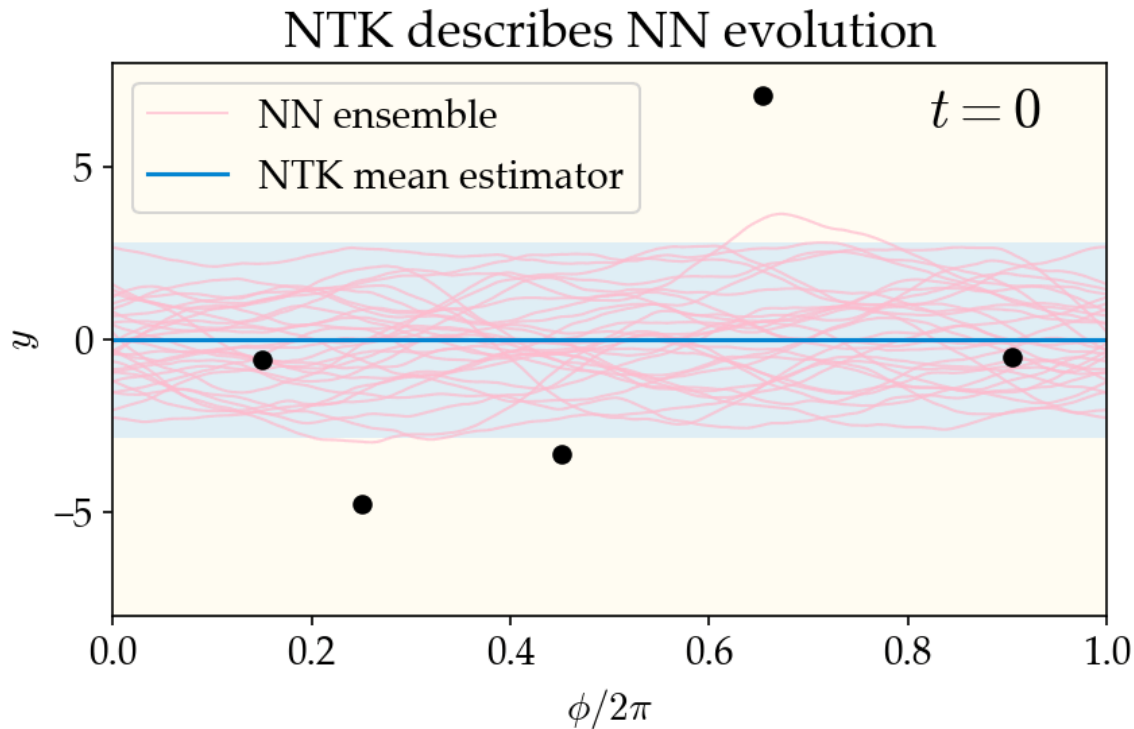
[1] Lee et al. arXiv:1711.00165 (2017)

[2] Jacot et al, NIPS (2018)

[3] Lee et al, NIPS (2019)

Extension to neural networks

- Toy example (from Wikipedia)



Last-layer Prediction Rigidity (LLPR)

- A last-layer approximation of the prediction rigidity recovers:
 - the NNGP exactly at initialization

$$K_{\text{NNGP}}(\mathbf{x}_i, \mathbf{x}_j) \approx \sigma_w^2 \mathbf{f}_i^\top \mathbf{f}_j \implies \mathbf{K}_{\text{NNGP}}(\mathcal{D}, \mathcal{D}) \approx \sigma_w^2 \mathbf{F}\mathbf{F}^\top$$

- the NTK to a good approximation

$$K_{\text{NTK}}(\mathbf{x}_i, \mathbf{x}_j) \approx c \left(\frac{\partial \tilde{y}(\mathbf{x}_i, \mathbf{w})}{\partial \mathbf{w}_L} \right)^\top \frac{\partial \tilde{y}(\mathbf{x}_j, \mathbf{w})}{\partial \mathbf{w}_L} = c \mathbf{f}_i^\top \mathbf{f}_j \implies \mathbf{K}_{\text{NTK}}(\mathcal{D}, \mathcal{D}) \approx c \mathbf{F}\mathbf{F}^\top$$

California housing \$: NN width test

