

# Uncertainty Quantification of Machine Learning to Establish AI Trustworthiness in Nuclear Engineering Applications

Xu Wu

[xwu27@ncsu.edu](mailto:xwu27@ncsu.edu)

Assistant Professor  
Department of Nuclear Engineering  
North Carolina State University

Data Science and Artificial Intelligence Regulatory Applications Workshops  
Workshop #4: AI Characteristics for Regulatory Consideration  
Panel Session on “AI Safety, Security and Explainability”  
The U.S. Nuclear Regulatory Commission (NRC)

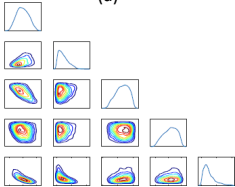
September 19, 2023

# Sources of uncertainties in physical modeling & simulation

## Parameter Uncertainty

unknown exact values of the model input parameters, and/or randomness

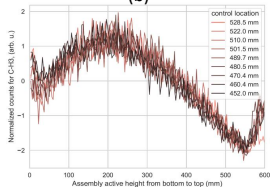
(a)



## Experimental/Data Uncertainty

noise or error in the measurement and/or data processing process

(b)



## Numerical Uncertainty

numerical approximation errors due to e.g., insufficient convergence, mesh

(c)



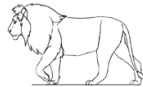
Observation



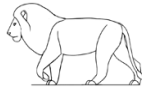
Good model



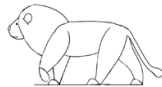
Model 1



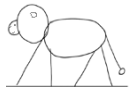
Model 2



Model 3



Model 4

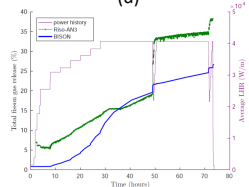


Model 5

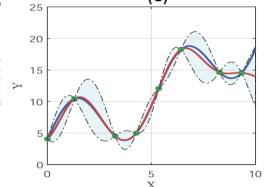
## Model Uncertainty (Bias, Discrepancy)

missing, inaccurate and/or incomplete underlying physics in the computer model

(d)



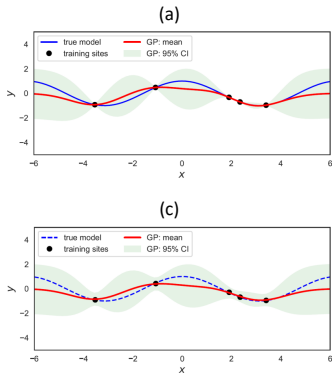
(e)



# Sources of uncertainties in data-driven Machine Learning models

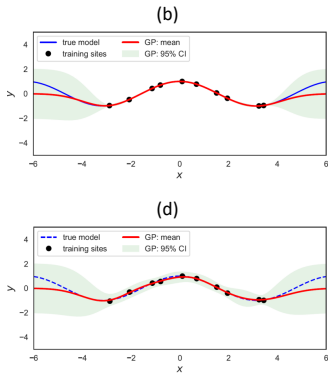
## Data Noise

noises in training data from either physical simulation models or experiments



## Data Coverage

few and/or gappy data that has incomplete coverage of training domain

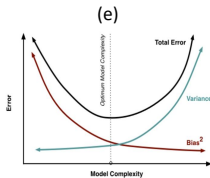


## Extrapolation

generalization to the extrapolated domains outside of the training domain

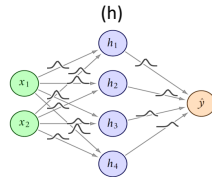
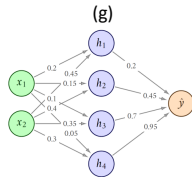
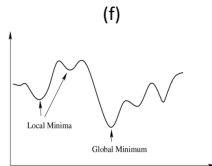
## Imperfect Model

ML model architecture is not properly defined, e.g., model is too simple or too complex

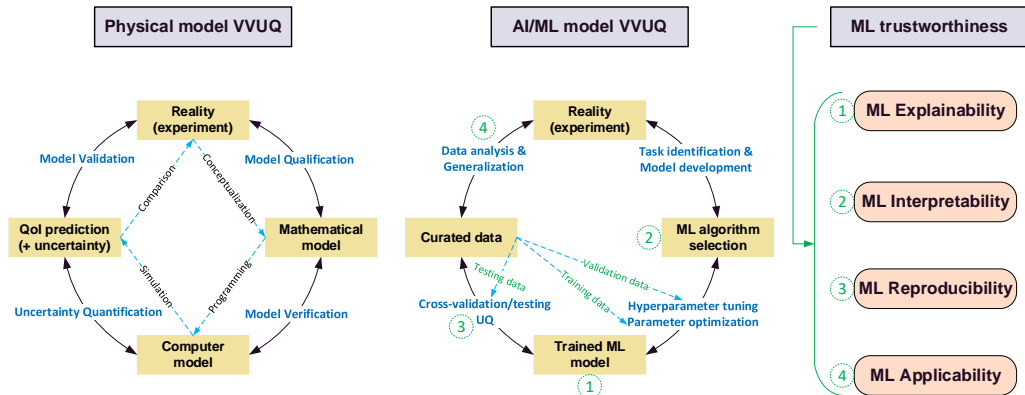


## Training

random initialization, convergence to local minima, hyper-parameter tuning, etc



- **Application-agnostic algorithms**, or those designed for more traditional ML applications such as computer vision and natural language processing, cannot typically be directly applied to scientific data in nuclear applications and require non-trivial, task-specific modifications.
- **Low-consequence error-tolerant settings** → **high-consequence nuclear systems**, need to establish **ML trustworthiness**, including accuracy, robustness (reproducibility, applicability), algorithmic fairness, algorithmic transparency (explainability, interpretability), and privacy.



- **Explainability**<sup>1,2</sup>: the ability to ensure that algorithmic decisions as well as any data driving those decisions can be explained to end-users and stakeholders in non-technical terms. It helps stakeholders and decision makers to understand ML solutions by “opening the black-box”.
- **Interpretability**<sup>3</sup>: the degree that an ML model obeys structural knowledge of the domain, such as monotonicity, causality, structural constraints, additivity, or physical constraints that come from domain knowledge.
- **Reproducibility**<sup>4,5</sup>: the ability of being able to replicate the ML model from data processing to model design, reporting, model analysis, or evaluation to successful deployment.
- **Applicability**<sup>6</sup>: the usability of ML for new scenarios such as unseen domains.
- **Other definitions: NIST framework on AI trustworthiness**<sup>7,8</sup> consists of nine factors: **accuracy, reliability, resiliency, objectivity, security, explainability, safety, accountability and privacy.**

---

<sup>1</sup>Barocas, S., Friedler, S., Hardt, M., Kroll, J., et al. (2018). The FAT-ML Workshop Series on Fairness, Accountability, and Transparency in Machine Learning.

<sup>2</sup>Gunning, D., Vorm, E., Wang, J. Y., & Turek, M. (2021). DARPA's explainable AI (XAI) program: A retrospective. *Applied AI Letters*, 2(4), e61.

<sup>3</sup>Rudin, C. (2019). Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature Machine Intelligence*, 1(5), 206-215.

<sup>4</sup>Heil, B. J., Hoffman, M. M., Markowitz, F., Lee, S. I., et al. (2021). Reproducibility standards for machine learning in the life sciences. *Nature Methods*, 18(10), 1132-1135.

<sup>5</sup>Beam, Andrew L., Arjun K. Manrai, and Marzyeh Ghassemi. "Challenges to the reproducibility of machine learning models in health care." *Jama* 323.4 (2020): 305-306.

<sup>6</sup>Wang, J., Lan, C., Liu, C., Ouyang, Y., Zeng, W., & Qin, T. (2021). Generalizing to unseen domains: A survey on domain generalization. *arXiv preprint arXiv:2103.03097*.

<sup>7</sup>Stanton, B., & Jensen, T. (2021). Trust and artificial intelligence. Draft NIST Interagency/Internal Report (NISTIR) 8332, National Institute of Standards and Technology, Gaithersburg, MD, URL: <https://doi.org/10.6028/NIST.IR.8332-draft>.

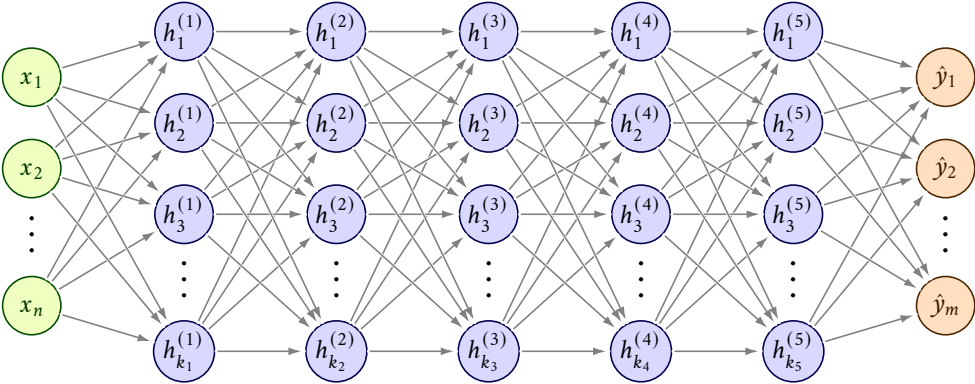
<sup>8</sup>Phillips, P. J., Hahn, C. A., Fontana, P. C., Broniatowski, D. A., & Przybocki, M. A. (2020). Four Principles of Explainable Artificial Intelligence, NIST Interagency/Internal Report (NISTIR) 8312, National Institute of Standards and Technology, Gaithersburg, MD, <https://doi.org/10.6028/NIST.IR.8312>

# How to quantify the approximation/prediction uncertainties in deep neural networks<sup>9</sup>?

input layer

hidden layers

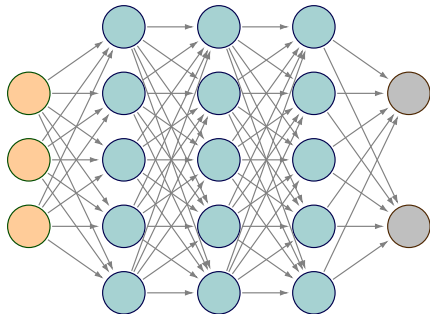
output layer



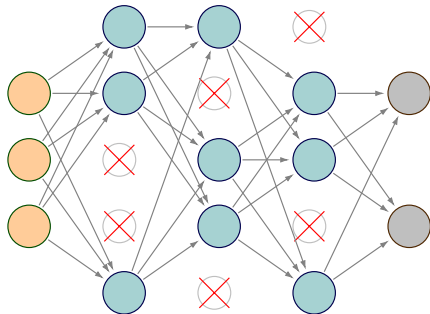
<sup>9</sup>Yaseen, M., & Wu, X. (2023). Quantification of Deep Neural Network Prediction Uncertainties for VVUQ of Machine Learning Models. Nuclear Science and Engineering, 197(5), 947-966.

## Monte Carlo Dropout (MCD)<sup>10</sup> introduces randomness to prediction in addition to training

- The **training** step is performed in the regular way, using stochastic gradient descent methods and re-evaluating the dropout matrices before each learning step.
- At the **prediction** step, we again evaluate the dropout matrices before every forward pass resulting in random network outputs.



Regular DNN



DNN after dropout

<sup>10</sup>Gal, Y., & Ghahramani, Z. (2016). Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In international conference on machine learning (pp. 1050-1059).

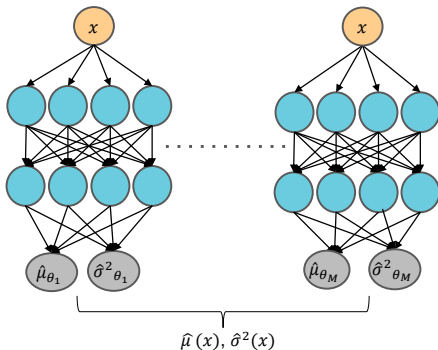
## Deep ensembles (DE)<sup>11</sup> changes network predictions as distributional parameters

- DE assumes the data to have a given **parameterized distribution** (e.g., Gaussian) where the distribution parameters depend on the input.
- Use the **negative log-likelihood function** of the Gaussian distribution as the cost function:

$$\mathcal{L}_\theta(\mathbf{x}, y) = -\log \phi_\theta(y|\mathbf{x}) = \frac{\log \hat{\sigma}_\theta^2(\mathbf{x})}{2} + \frac{(y - \hat{\mu}_\theta(\mathbf{x}))^2}{2\hat{\sigma}_\theta^2(\mathbf{x})} + c$$

- With an ensemble of  $M$  neural networks, the joint Gaussian has mean and variance given by:

$$\hat{\mu}(\mathbf{x}) = \frac{1}{M} \sum_{i=1}^M \hat{\mu}_{\theta_i}(\mathbf{x})$$
$$\hat{\sigma}^2(\mathbf{x}) = \frac{1}{M} \sum_{i=1}^M (\hat{\sigma}_{\theta_i}^2(\mathbf{x}) + \hat{\mu}_{\theta_i}^2(\mathbf{x})) - \hat{\mu}^2(\mathbf{x})$$



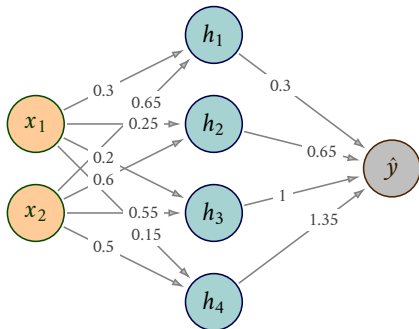
<sup>11</sup>Lakshminarayanan, B., Pritzel, A., & Blundell, C. (2017). Simple and scalable predictive uncertainty estimation using deep ensembles. In Advances in neural information processing systems (pp. 6402-6413).



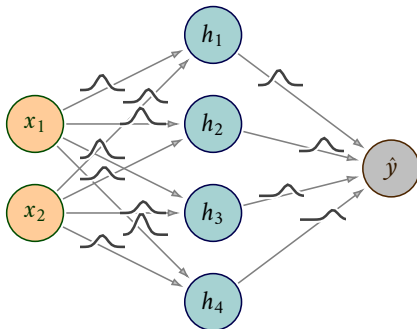
## Bayesian Neural Networks (BNNs)<sup>12</sup> - neural networks with distributions over parameters

- In BNNs, **prior distributions** are specified upon the parameters (weights, bias) of neural networks.
- Given the training data, the **posterior distributions** over the parameters are computed, which are used to quantify the predictive uncertainty.

Regular Neural Network



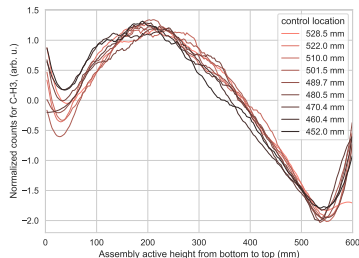
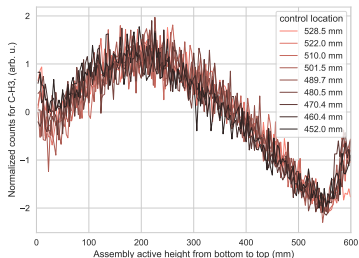
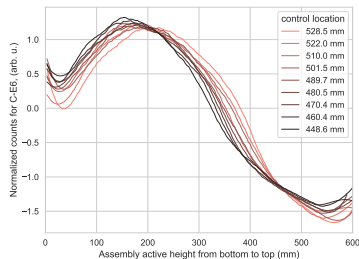
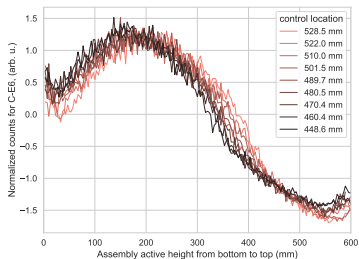
Bayesian Neural Network



<sup>12</sup>Blundell, C., Cornebise, J., Kavukcuoglu, K., & Wierstra, D. (2015). Weight uncertainty in neural network. In International conference on machine learning (pp. 1613-1622). PMLR.

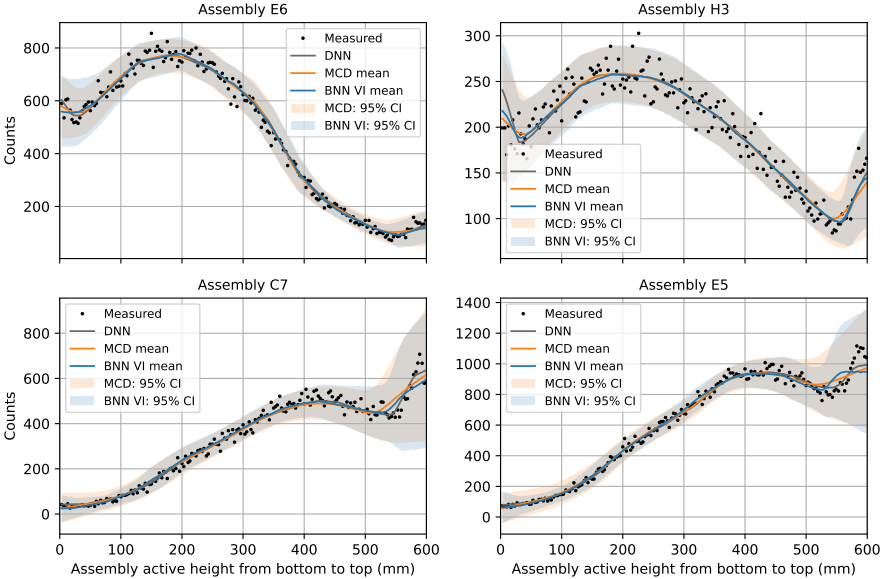
## Example: using DNNs to predict the axial neutron flux profiles given the control rod bank position

- **Training data**<sup>13</sup>: copper-wire activation measurements and measured control bank positions obtained from the SAFARI-1 research reactor (South Africa) historical cycles.



<sup>13</sup> Moloko, L. E., Bokov, P. M., Wu, X., & Ivanov, K. N. (2023). Prediction and uncertainty quantification of SAFARI-1 axial neutron flux profiles with neural networks. *Annals of Nuclear Energy*, 188, 109813.

# Example: the DNN predictions are made on assemblies and cycles that are unseen during training



Thank you for your attention!  
Questions and comments?