

Scientific AI in Cosmology



Jason D. McEwen

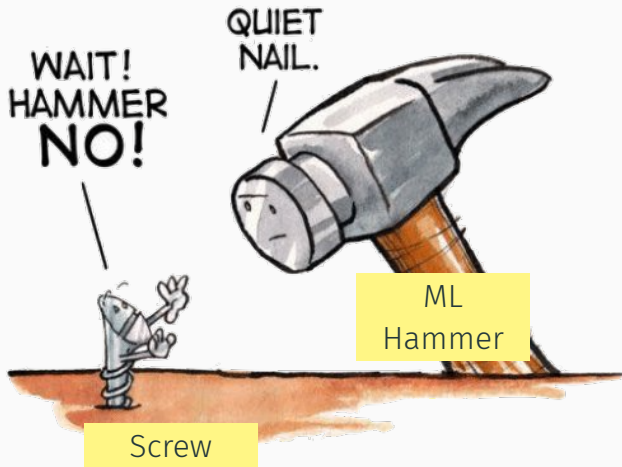
www.jasonmcewen.org

Scientific AI (SciAI) Group

Mullard Space Science Laboratory (MSSL), University College London (UCL)

Cosmo21, Chania, May 2024

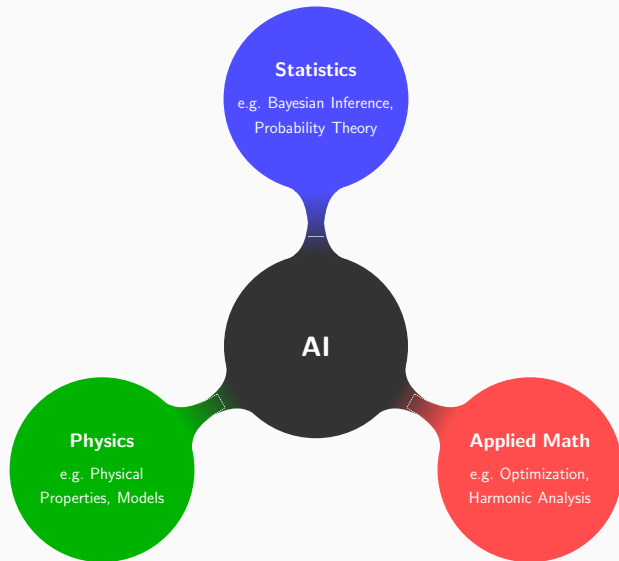
The machine learning hammer

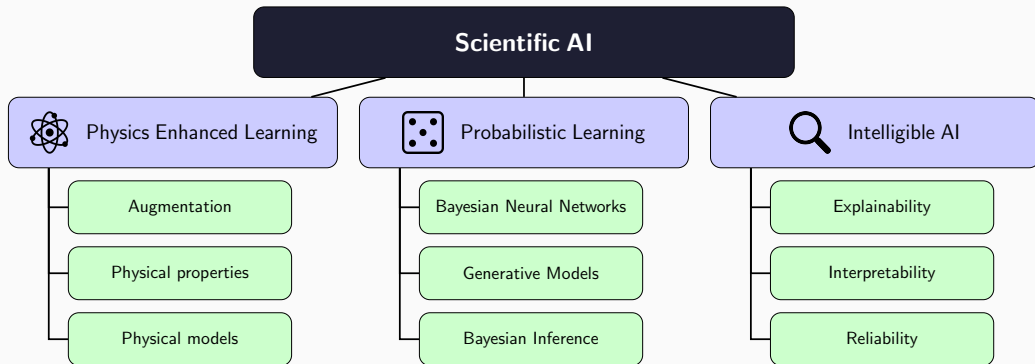


The machine learning cog



Merging paradigms



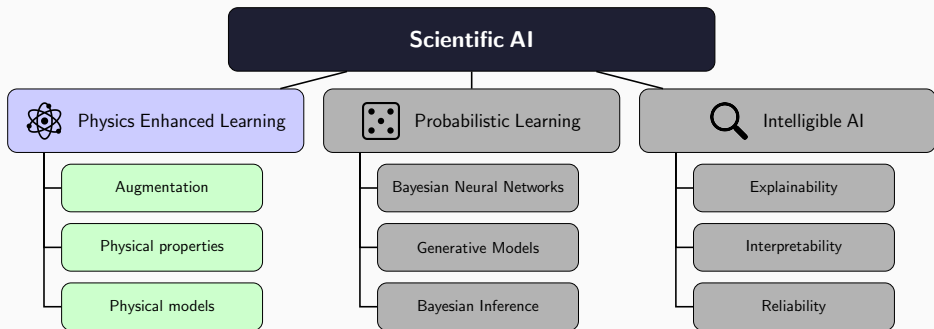


Physics Enhanced Learning

Physics Enhanced Learning

Embed physical understanding of the world into machine learning models.

(See review by Karniadakis *et al.* 2021.)



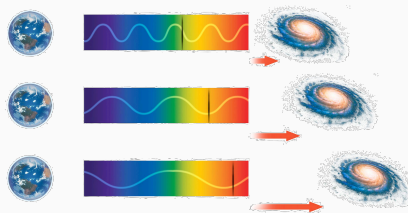


Apply **physical transformations** that data known to satisfy to augment training data \rightsquigarrow ML model **learns physics through training**.



Apply **physical transformations** that data known to satisfy to augment training data \rightsquigarrow ML model **learns physics through training**.

- ▷ Redshift augmentation of supernovae observations (Boone 2019, Alves *et al.* 2022, 2023)



Redshift augmentation



Apply **physical transformations** that data known to satisfy to augment training data \rightsquigarrow ML model **learns physics through training**.



▷ Data efficiency suffers: data “used” to learn physics, rather than problem.

Physical properties: geometries, symmetries, conservation laws



Encode physical properties of the world into ML models (e.g. geometry, symmetries, conservation laws) \rightsquigarrow **Physics embedded in architecture** of ML model.

Physical properties: geometries, symmetries, conservation laws



Encode physical properties of the world into ML models (e.g. geometry, symmetries, conservation laws) \rightsquigarrow **Physics embedded in architecture** of ML model.

- ▷ Geometric deep learning on the sphere (Cobb et al. 2021; McEwen et al. 2022; Ocampo, Price & McEwen 2023)

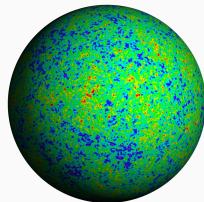
See Cosmo21 poster



Kevin Mulder



Matt Price



CMB observed on the celestial sphere

Physical models: PINNs and differentiable physics

Encode physical models of world into ML models:



1. Encode dynamics (differential equations) via loss functions (PINNs).
2. Embed full (differentiable) physical models inside ML model.

⇒ **Physics learned in training and embedded in model.**

Physical models: PINNs and differentiable physics

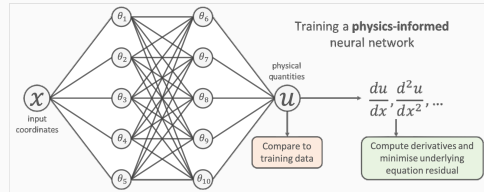
Encode physical models of world into ML models:



1. Encode dynamics (differential equations) via loss functions (PINNs).
2. Embed full (differentiable) physical models inside ML model.

↪ **Physics learned in training and embedded in model.**

- ▷ Physics informed neural networks (PINNs) encode differentiable equations (e.g. boundary conditions) in loss.



PINNs

Physical models: PINNs and differentiable physics

Encode physical models of world into ML models:

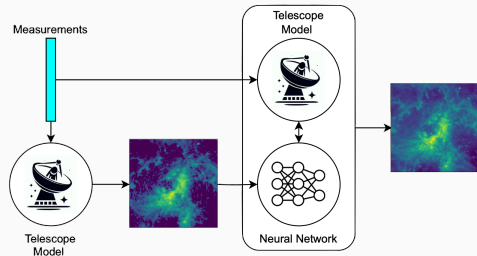


1. Encode dynamics (differential equations) via loss functions (PINNs).
2. Embed full (differentiable) physical models inside ML model.

⇒ **Physics learned in training and embedded in model.**

▷ Differentiable physical models

- ▶ Radio interferometric telescope
(Mars *et al.* 2023, 2024)
- ▶ Optical PSF
(Liaudat *et al.* 2023)
- ▶ JAX-Cosmo
(Campagne *et al.* 2023)



Physics inside AI models for imaging data from radio interferometric telescopes (Mars *et al.* 2024) 8

Physical models: PINNs and differentiable physics

Encode physical models of world into ML models:



1. Encode dynamics (differential equations) via loss functions (PINNs).
2. Embed full (differentiable) physical models inside ML model.

↪ Physics learned in training and embedded in model.

▷ Differentiable mathematical methods

- ▶ Spherical harmonic transforms
(`s2fft`; Price & McEwen 2024)
- ▶ Spherical wavelet transforms
(`s2wav`; Price *et al.* 2024)

See Cosmo21 poster



Matt Price



Alicja Polanska



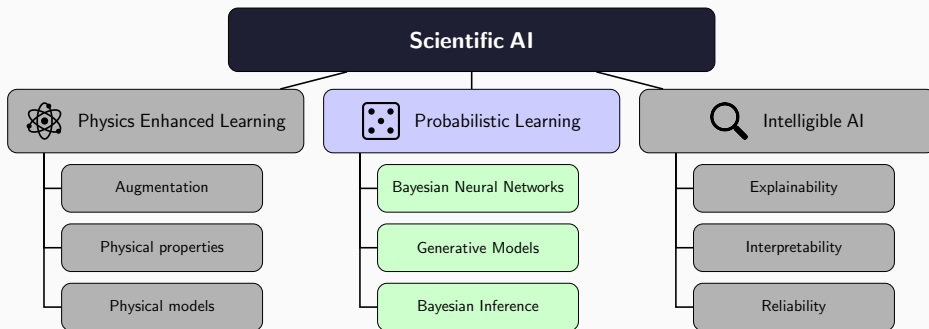
Jess Whitney

Probabilistic Learning

Probabilistic Learning

Embed a probabilistic representation of data, models and/or outputs.

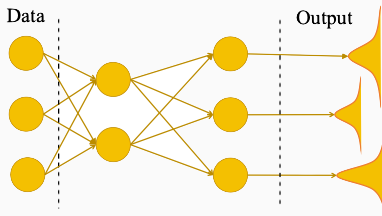
(See Murray 2022.)



Bayesian neural networks for uncertainty quantification



Bayesian neural networks incorporate **probabilistic representation** to quantify **uncertainty of outputs** (idea pioneered by MacKay 1992).

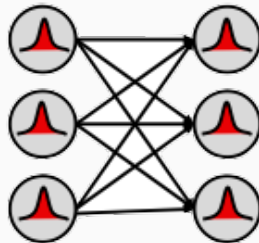


Bayesian neural networks for uncertainty quantification



Bayesian neural networks incorporate **probabilistic representation** to quantify **uncertainty of outputs** (idea pioneered by MacKay 1992).

- ▷ MC Dropout (Gal & Ghahramani 2016): drop nodes probabilistically to sample an ensemble of networks.

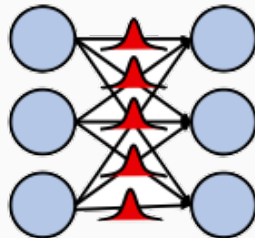


Bayesian neural networks for uncertainty quantification



Bayesian neural networks incorporate **probabilistic representation** to quantify **uncertainty of outputs** (idea pioneered by MacKay 1992).

- ▷ Bayes by Backprop (Blundel *et al.* 2015): model distribution of weights (by variational inference).



Bayesian neural networks for uncertainty quantification



Bayesian neural networks incorporate **probabilistic representation** to quantify **uncertainty of outputs** (idea pioneered by MacKay 1992).



- ▷ Encode epistemic uncertainty of model.
- ▷ But what does the output distribution represent?
- ▷ Requires careful consideration of training data.



Generative models **learn a prior distribution** from data for sampling and/or evaluating probabilities.

Generative models



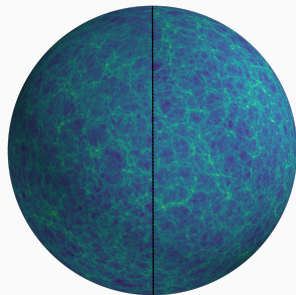
Generative models **learn a prior distribution** from data for sampling and/or evaluating probabilities.

- ▷ Emulation: sample from learned prior
(Perraudin *et al.* 2020, Allys *et al.* 2020, Price *et al.* 2023, Price *et al.* in prep., Mousset, Price, Allys, McEwen in prep.)

See Cosmo21 poster



Matt Price



Emulated LSS
(Mousset, Price, Allys, McEwen in prep.)



ML techniques can be integrated into Bayesian frameworks to **enhance accuracy and computational efficiency**, making some approaches accessible that were previously intractable.



ML techniques can be integrated into Bayesian frameworks to **enhance accuracy and computational efficiency**, making some approaches accessible that were previously intractable.

▷ Enhanced Bayesian model selection
(**harmonic**; McEwen *et al.* 2021, Polanska *et al.* 2023, 2024, Piras *et al.* in prep.)

- ▶ Only requires posterior samples.
- ▶ Agnostic to sampling technique:
 - ↔ Leverage efficient samplers.
 - ↔ Variational inference.
- ▶ Scale to high dimensions.

See Cosmo21 talk



Alicja Polanska



Davide Piras



Matt Price

Bayesian inference



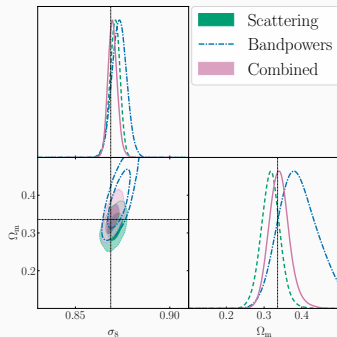
ML techniques can be integrated into Bayesian frameworks to **enhance accuracy and computational efficiency**, making some approaches accessible that were previously intractable.

- ▷ Simulation-based inference (SBI)
(Alsing *et al.* 2018, Cranmer *et al.* 2021, Lin *et al.* 2022, in prep., von Wietersheim-Kramsta *et al.* 2024)
- ▷ Model selection for SBI (Spurio Mancini *et al.* 2022)

See Cosmo21 poster



Kiyam Lin



SBI with scattering transform
(Lin *et al.* in prep.)



ML techniques can be integrated into Bayesian frameworks to **enhance accuracy and computational efficiency**, making some approaches accessible that were previously intractable.

▷ Variational inference (Whitney *et al.* in prep.)

See Cosmo21 talk



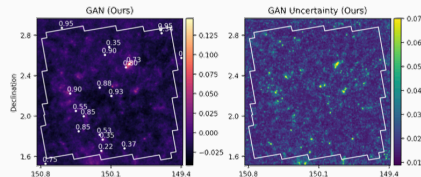
Jess Whitney



Tobias Liaudat



Matt Price



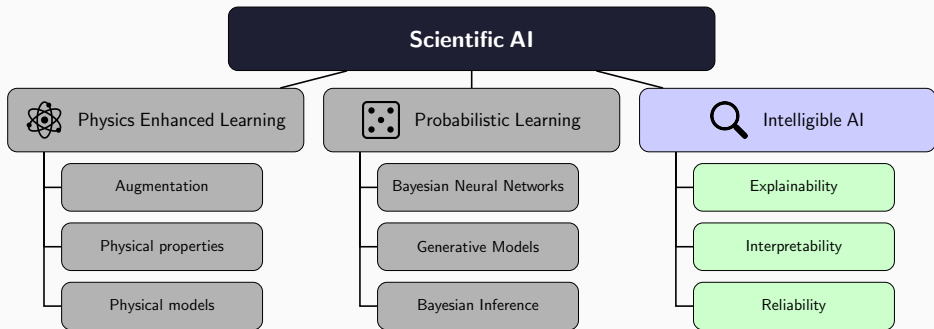
Mass mapping with uncertainties by variational inference (Whitney *et al.* in prep.)

Intelligible AI

Intelligible AI

Machine learning methods that are able to be understood by humans.

(See Weld & Bansal 2018, Ras *et al.* 2020.)



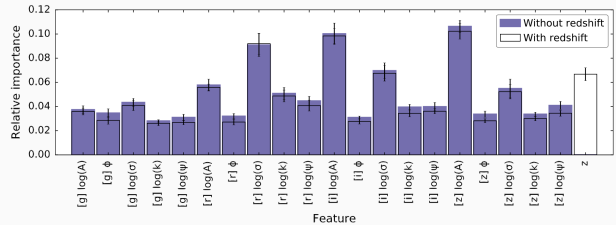


Explainable ML techniques may or may not be interpretable themselves but their outputs can be explained to humans.



Explainable ML techniques may or may not be interpretable themselves but their outputs can be explained to humans.

- ▷ Feature importances (Lochner *et al.* 2016)

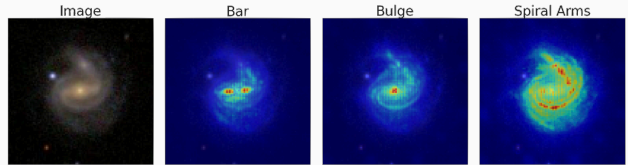


Supernova feature importances



Explainable ML techniques may or may not be interpretable themselves but their outputs can be explained to humans.

- ▷ Saliency maps
(Bhambra *et al.* 2022)



Galaxy saliency mapping



Explainable ML techniques may or may not be interpretable themselves but their outputs can be explained to humans.



Poking the black box: may provide some explanation of outputs but humans still not able to comprehend underlying process.



Interpretable ML models are **white boxes** that can be understood by humans.



Interpretable ML models are **white boxes** that can be understood by humans.

- ▷ Deep priors learned from training data (hybrid model-based and data-driven) (Remy *et al.* 2022, McEwen *et al.* 2023)

See Cosmo21 poster



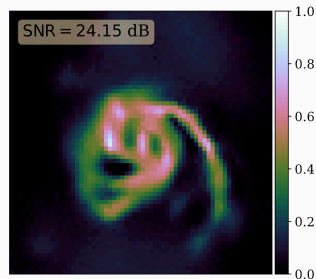
Tobias Liaudat



Henry Aldridge



Matt Price



Compute Bayesian evidence for model selection (proxnest, McEwen *et al.* 2023)



Interpretable ML models are **white boxes** that can be understood by humans.

- Interpretable constraints on ML models, *e.g.* convexity (Liaudat *et al.* 2023)

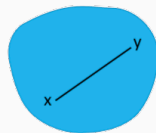
See Cosmo21 talk



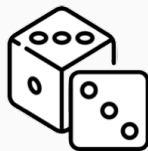
Tobias Liaudat



Matt Price



Convexity



Uncertainty Quantification

Impose convexity on learned model

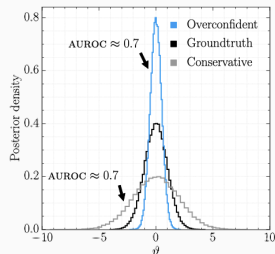


Reliability **critical for science** in order for humans to have confidence in results of ML models. Closely coupled with a **meaningful statistical distribution** of outputs.



Reliability **critical for science** in order for humans to have confidence in results of ML models. Closely coupled with a **meaningful statistical distribution** of outputs.

- ▷ Validity of statistical distributions
(Hermans *et al.* 2022, Lemos *et al.* 2023)



Validity of distribution
(Hermans *et al.* 2022)



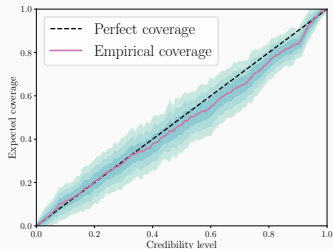
Reliability **critical for science** in order for humans to have confidence in results of ML models. Closely coupled with a **meaningful statistical distribution** of outputs.

- ▷ Validity of statistical distributions
(Hermans *et al.* 2022, Lemos *et al.* 2023)

See Cosmo21 poster



Kiyam Lin



Coverage analysis for SBI with scattering (Lin *et al.* in prep.)



Reliability **critical for science** in order for humans to have confidence in results of ML models. Closely coupled with a **meaningful statistical distribution** of outputs.

- ▷ Diversity (avoiding mode-collapse)
(Price *et al.* 2023, Whitney *et al.* in prep.)



Jess Whitney

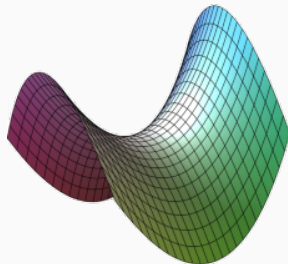


Tobias Liaudat



Matt Price

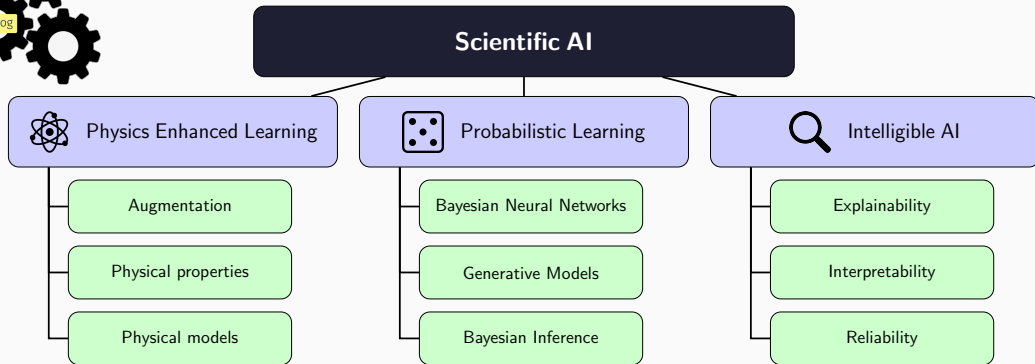
See Cosmo21 talk



Recover probability distribution over full underlying manifold

Summary

Summary



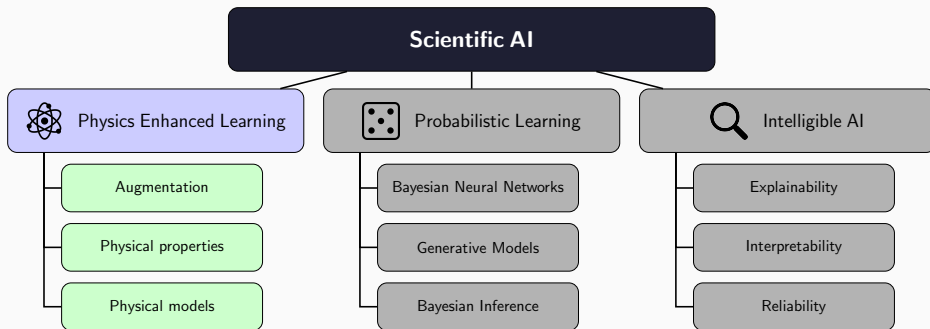
With great power comes great responsibility!

Extra Slides

Physics Enhanced Learning

Embed physical understanding of the world into machine learning models.

(See review by Karniadakis *et al.* 2021.)





Apply **physical transformations** that data known to satisfy to augment training data \rightsquigarrow ML model **learns physics through training**.



▷ Data efficiency suffers: data “used” to learn physics, rather than problem.

Physical properties: geometries, symmetries, conservation laws



Encode physical properties of the world into ML models (e.g. geometry, symmetries, conservation laws) \rightsquigarrow **Physics embedded in architecture** of ML model.



- ▷ Highly computationally demanding.
- ▷ Always required?



- ▷ Develop efficient algorithms (e.g. Ocampo, Price & McEwen 2023).
- ▷ Inductive biases not enforced.

Physical models: PINNs and differentiable physics

Encode physical models of world into ML models:



1. Encode dynamics (differential equations) via loss functions (PINNs).
2. Embed full (differentiable) physical models inside ML model.

↪ **Physics learned in training and embedded in model.**



- ▷ PINNs only capture limited dynamics via loss.
- ▷ Full physical models requires differentiable programming frameworks.

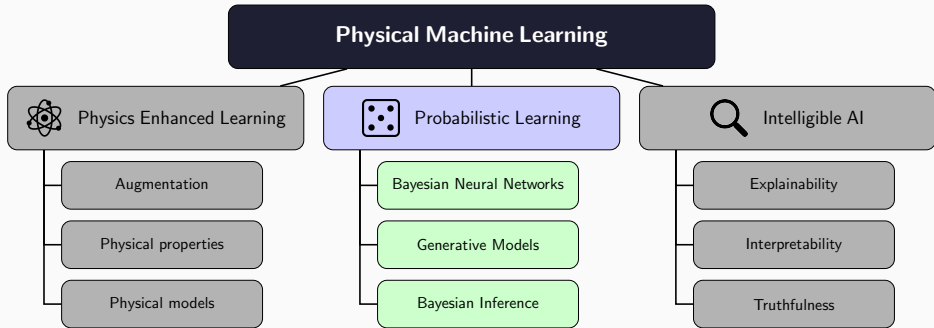


- ▷ Capture full physics with differentiable models!
- ▷ Emulators also provide differentiability (e.g. `CosmoPower`; Spurio Mancini et al. 2021).
- ▷ Write new differentiable codes (e.g. `s2fft`; Price & McEwen 2023).

Probabilistic Learning

Embed a probabilistic representation of data, models and/or outputs.

(See Murray 2022.)



Bayesian neural networks for uncertainty quantification



Bayesian neural networks incorporate **probabilistic representation** to quantify **uncertainty of outputs** (idea pioneered by MacKay 1992).



- ▶ Encode epistemic uncertainty of model.
- ▶ But what does the output distribution represent?
- ▶ Requires careful consideration of training data.



- ▶ Statistical validation (hold that thought... see upcoming Reliability section).

Generative models



Generative models **learn a prior distribution** from data for sampling and/or evaluating probabilities.



- ▷ Availability and representativeness of training data.
- ▷ Reliability, *e.g.* diversity of ML model often lacking.



- ▷ Public datasets/benchmarks (*e.g.* BASE, IllustrisTNG, CAMELS, Quijote, CosmoGrid).
- ▷ Meta sampling to recover distribution over manifold (*e.g.* Price *et al.* 2023).
- ▷ Reliability (hold that thought... see upcoming Reliability section).



ML techniques can be integrated into Bayesian frameworks to **enhance accuracy and computational efficiency**, making some approaches accessible that were previously intractable.



- ▷ Availability and representativeness of training data.
- ▷ Cost of training.
- ▷ Reliability?

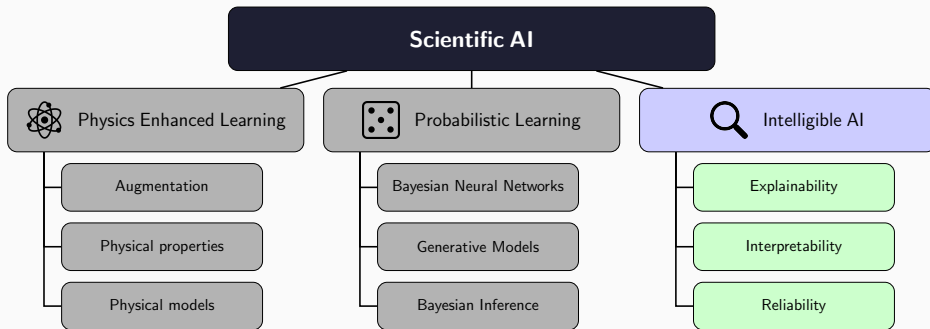


- ▷ Public datasets/benchmarks (e.g. BASE, IllustrisTNG, CAMELS, Quijote, CosmoGrid).
- ▷ Amortized inference (training **not** repeated for new observations).
- ▷ Integrate in Bayesian framework to provide statistical guarantees.
- ▷ Statistical validation (hold that thought... see upcoming Reliability section).

Intelligible AI

Machine learning methods that are able to be understood by humans.

(See Weld & Bansal 2018, Ras *et al.* 2020.)





Explainable ML techniques may or may not be interpretable themselves but their outputs can be explained to humans.



Poking the black box: may provide some explanation of outputs but humans still not able to comprehend underlying process.



Interpretable ML models are **white boxes** that can be understood by humans.



- ▶ Designed models limit flexibility.
- ▶ Availability and representativeness of training data.



- ▶ Benefits of designed models often outweigh (minimal) reduced flexibility.
- ▶ Public datasets/benchmarks (e.g. IllustrisTNG, CAMELS, Quijote, CosmoGrid).
- ▶ Transfer learning, self-supervised learning.



Reliability **critical for science** in order for humans to have confidence in results of ML models. Closely coupled with a **meaningful statistical distribution** of outputs.



- ▷ Uncertainties not always meaningful.
- ▷ Diversity of ML model often lacking.



- ▷ Integrate in statistical framework to inherit theoretical guarantees.
- ▷ Extensive validation tests (*e.g.* Hermans *et al.* 2022, Lemos *et al.* 2023).
- ▷ Meta sampling to recover distribution over manifold (*e.g.* Price *et al.* 2023).
- ▷ Well-posed frameworks (*e.g.* physics enhanced, probabilistic).