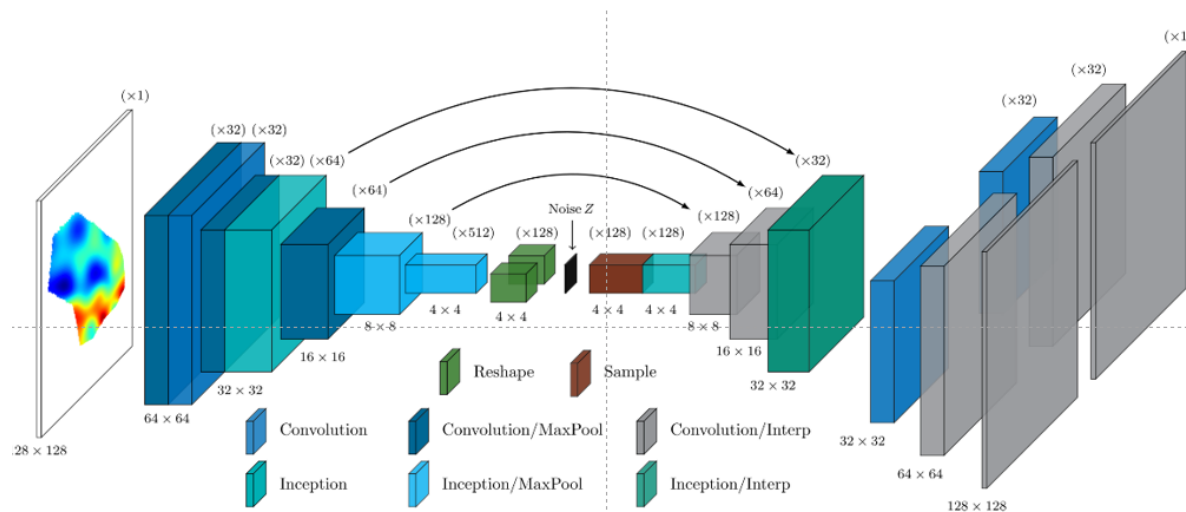


Towards Third Wave AI: Interpretable, Robust Trustworthy Machine Learning for Diverse Applications in Science and Engineering

Guang Lin,

Director, Data Science Consulting Service, Full Professor of Departments of Mathematics, Statistics, & Mechanical Engineering, Purdue University



MATH-IMS Joint Applied Mathematics Colloquium, Chinese University of Hong Kong, May 20, 2022

The Four Waves of AI

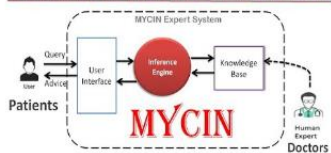
First Wave

c. 1970s - 1990s

Good at reasoning, but no ability to learn or generalize.

- GOFAI - "Good Old Fashioned AI."
- Symbolic, heuristic, rule based.
- Handcrafted knowledge, "expert systems."

ARTIFICIAL INTELLIGENCE



Second Wave

c. 2000s - present

Good at learning and perceiving, but minimal ability to reason or generalize.

- Statistical learning, "deep" neural nets, CNNs, RNNs.
- Advanced text, speech, language and vision processing.



Third Wave

est. 2020s - 2030s

Excellent at perceiving, learning and reasoning, and able to generalize.

- Contextual adaptation, able to explain decisions.
- Can converse in natural language.
- Requires far fewer data samples for training.
- Able to learn and function with minimal supervision.

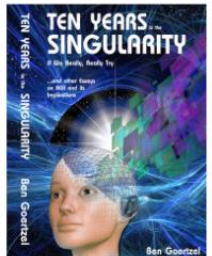
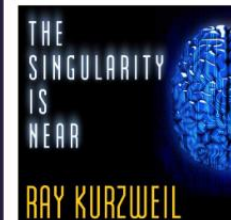
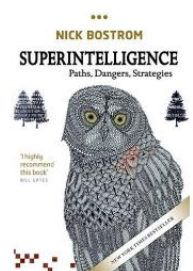


Fourth Wave

est. 2030s →

Able to perform any intellectual task that a human can.

- AGI (Artificial General Intelligence), possibly leading to ASI (Artificial Superintelligence) and the "Technological Singularity."



Six Kin Development (adapted from DARPA's "Three Waves of AI")

Artificial Intelligence: Image Coloring

100 year old pictures...



Credit to ColdFusion

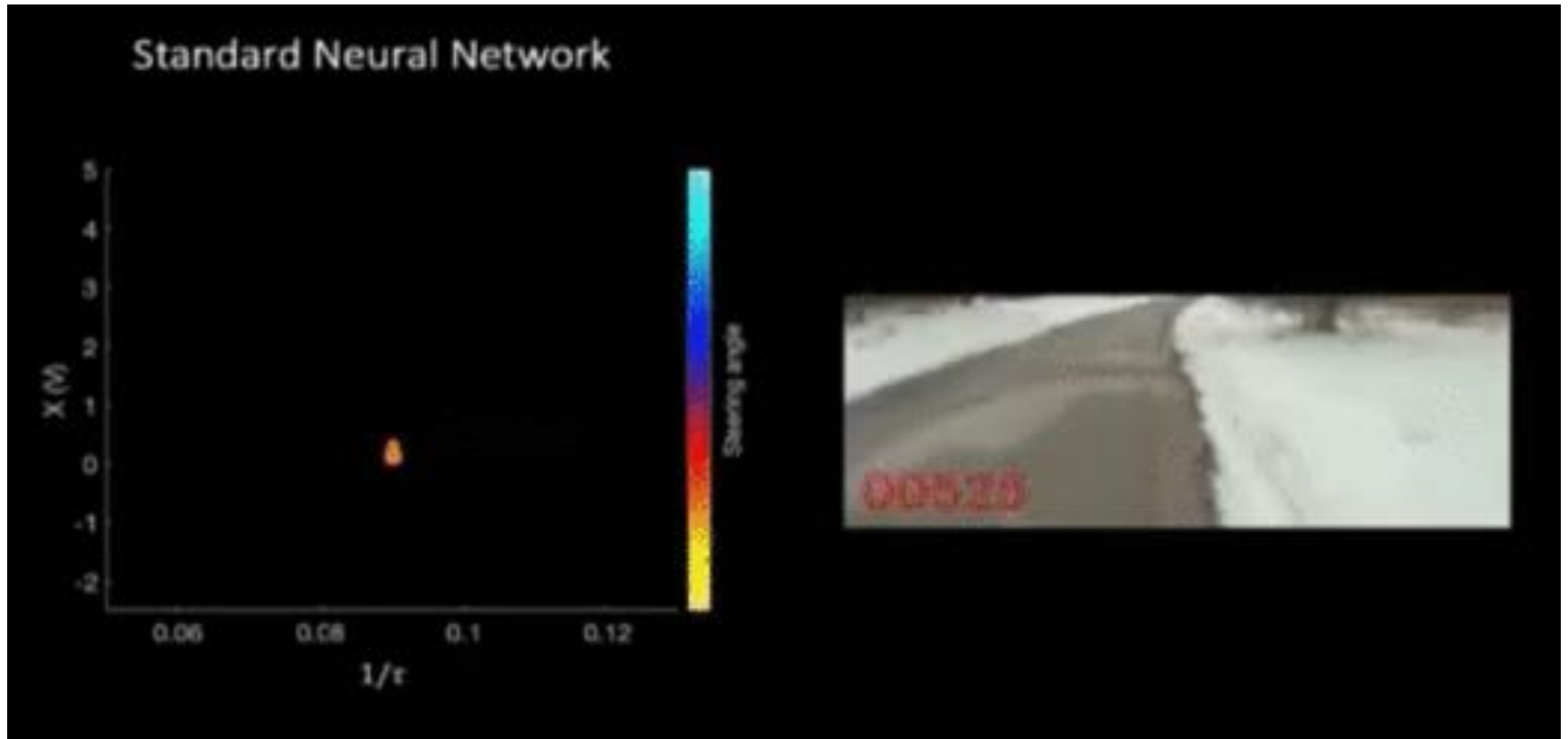
Artificial Intelligence: AlphaGo



Reinforcement Learning in AlphaGo

Credit to DeepMind

Artificial Intelligence: Autonomous Driving

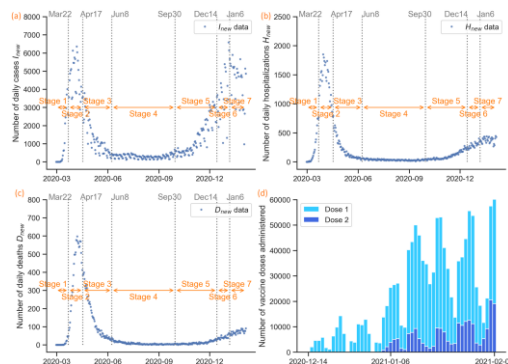


Credit to Ramin Hasani, MIT.

Guang Lin's Group's Main Research

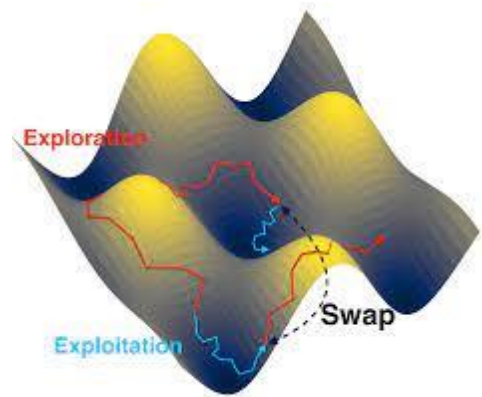
COVID 19 Pandemic Prediction

1. Nature Computational Science, 1-10, 2021
2. PLOS Computational Biology, 2021



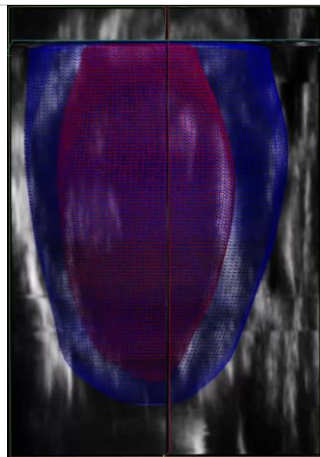
Uncertainty Quantification for Deep Learning

- Tier 1 AI conference: ICLR21, ICML21, WSDM21



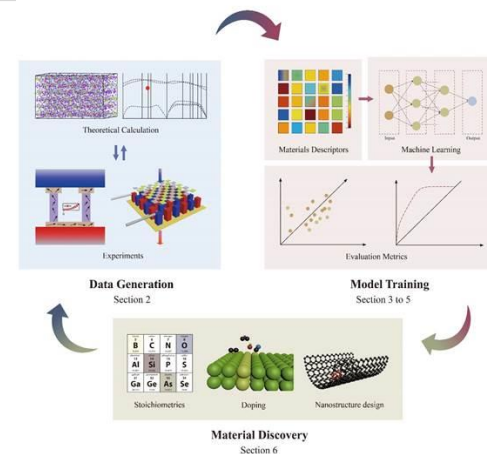
AI for Cardiac 4D Ultrasound

- Received Trask Innovation Award
- Patent submitted
- Applied Sciences 21



Machine Learning for Material Discovery

- ES Materials & Manufacturing 21



Outline:

- ❖ **Incorporate Physics Knowledge and AI to design new interpretable models**
- ❖ Incorporate Physics Knowledge into AI to predict multiscale problems: NH-PINN
- ❖ Interpretable AI enables data-driven scientific discovery with uncertainty quantification capability – ALZHEIMER's Disease Prediction
- ❖ Sparse Neural Architecture Design with quantified uncertainties
- ❖ Scalable training large-scale Deep Neural Network

How to incorporate Physics Knowledge and AI to design new interpretable models? - Interpretable AI for Science

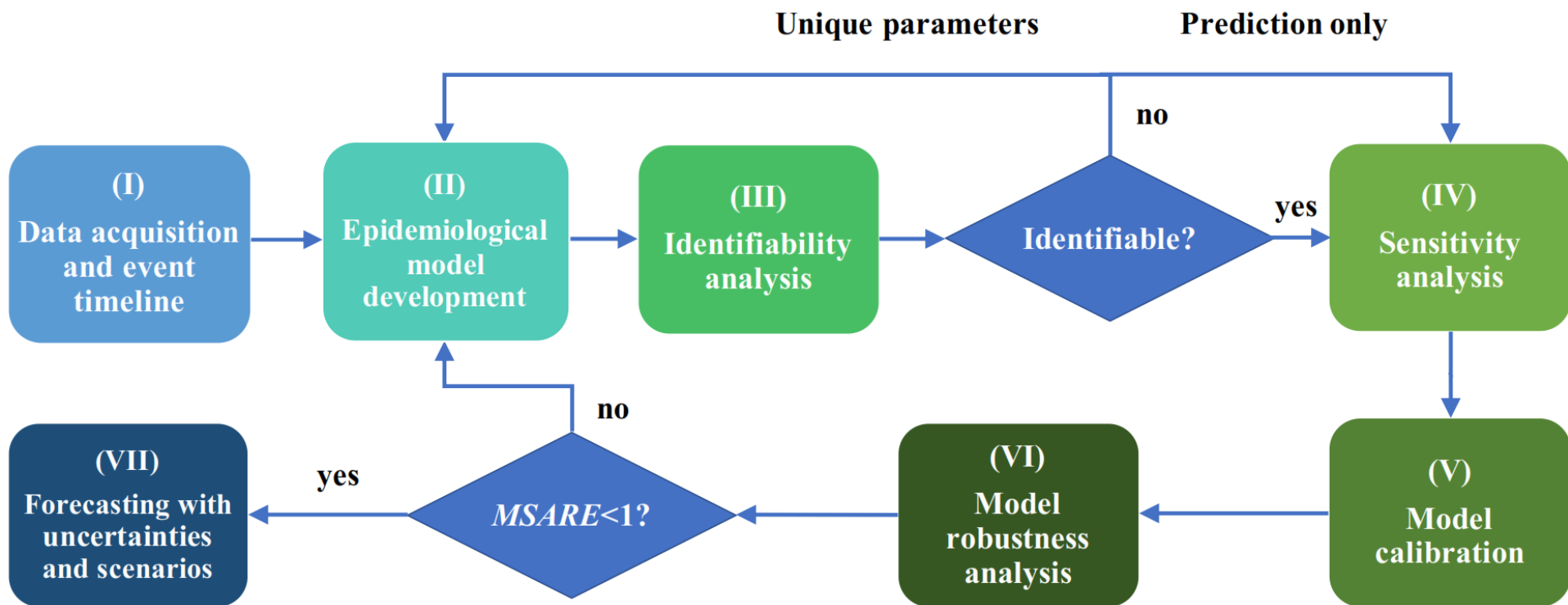
1. Ehsan Kharazmi, Min Cai, Xiaoning Zheng, Guang Lin, George Em Karniadakis, **Identifiability and predictability of integer- and fractional-order epidemiological models using physics-informed neural networks**, *Nature Computational Science*, **1**, 744-753, 2021
2. Sheng Zhang, Joan Ponce, Zhen Zhang, Guang Lin, George Karniadakis, **An integrated framework for building trustworthy data-driven epidemiological models: Application to the COVID-19 outbreak in New York City**, *PLoS Computational Biology* 17(9): e1009334. <https://doi.org/10.1371/journal.pcbi.1009334>

Predicting the COVID-19 pandemic with uncertainties using trustworthy data-driven epidemiological models

1. Ehsan Kharazmi, Min Cai, Xiaoning Zheng, Guang Lin, George Em Karniadakis, **Identifiability and predictability of integer- and fractional-order epidemiological models using physics-informed neural networks**, *Nature Computational Science*, 1, 744-753, 2021
2. Sheng Zhang, Joan Ponce, Zhen Zhang, Guang Lin, George Karniadakis, **An integrated framework for building trustworthy data-driven epidemiological models: Application to the COVID-19 outbreak in New York City**, *PLoS Computational Biology* 17(9): e1009334. <https://doi.org/10.1371/journal.pcbi.1009334>



A general framework for building a trustworthy data-driven epidemiological model

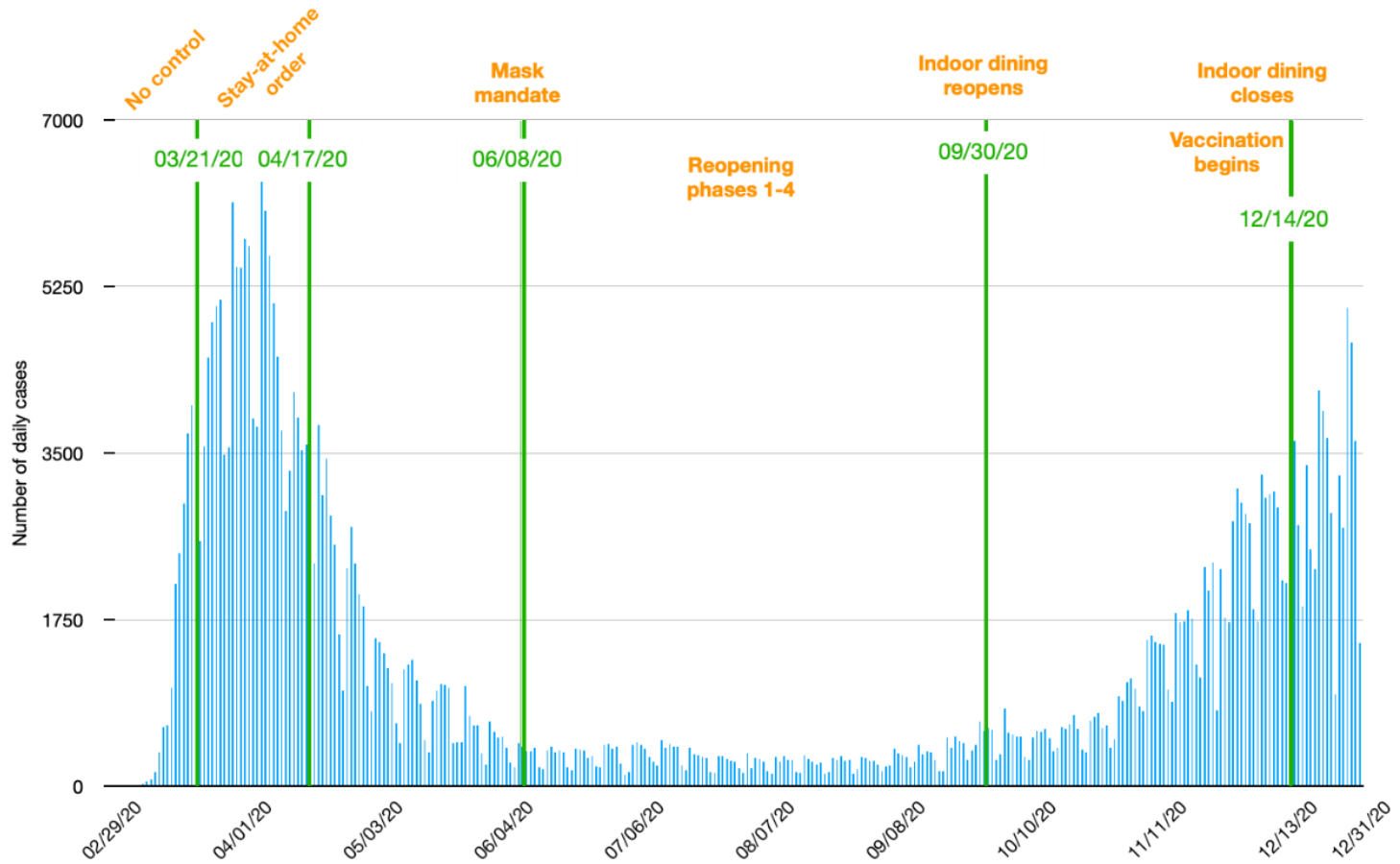


Sheng Zhang, Joan Ponce, Zhen Zhang, Guang Lin, George Karniadakis, **An integrated framework for building trustworthy data-driven epidemiological models: Application to the COVID-19 outbreak in New York City**, *PLoS Computational Biology* 17(9): e1009334.

<https://doi.org/10.1371/journal.pcbi.1009334>

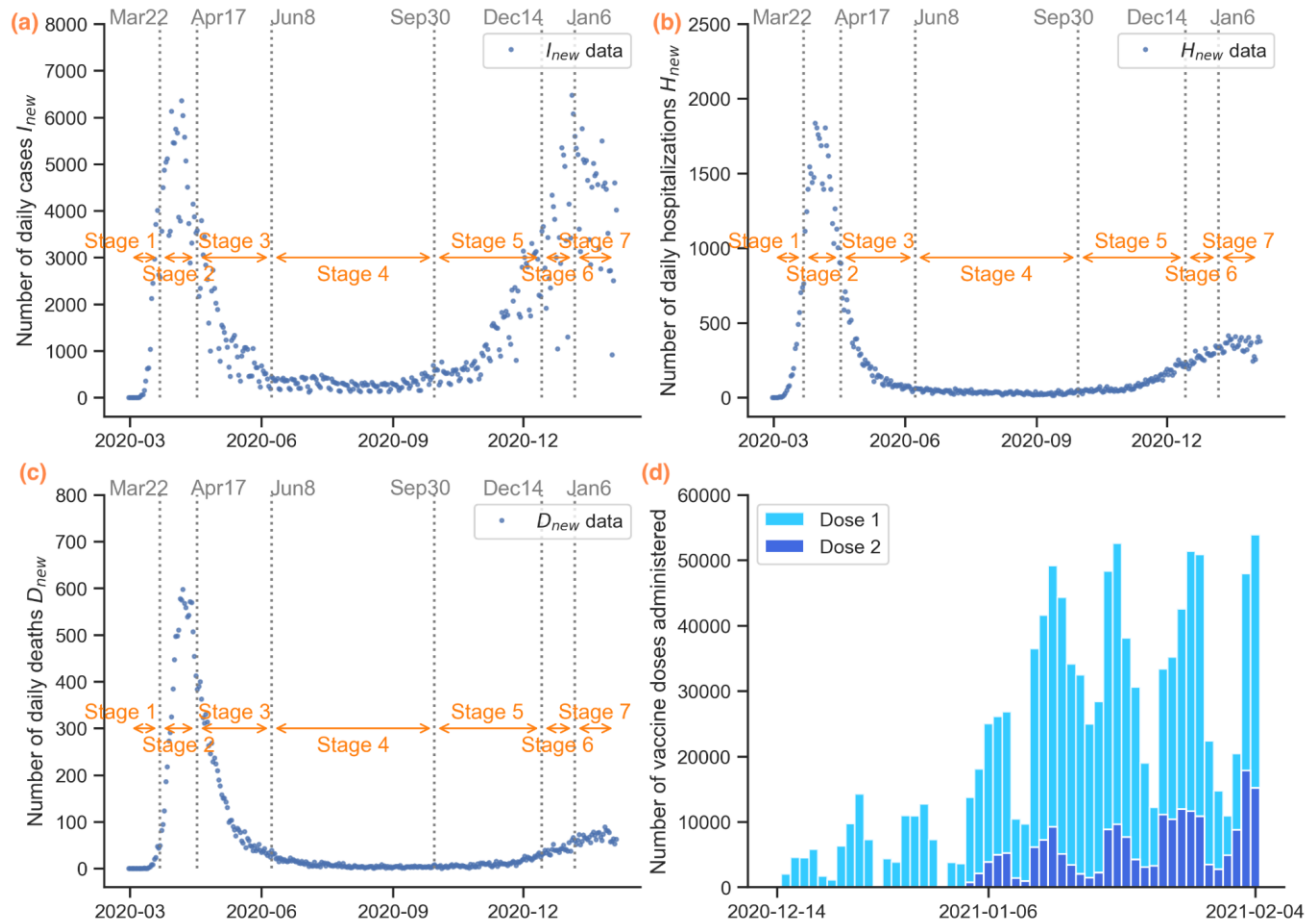
New York City COVID-19 related Event Timeline

Calibrate piecewise-constant model parameters to capture local epidemiological dynamics



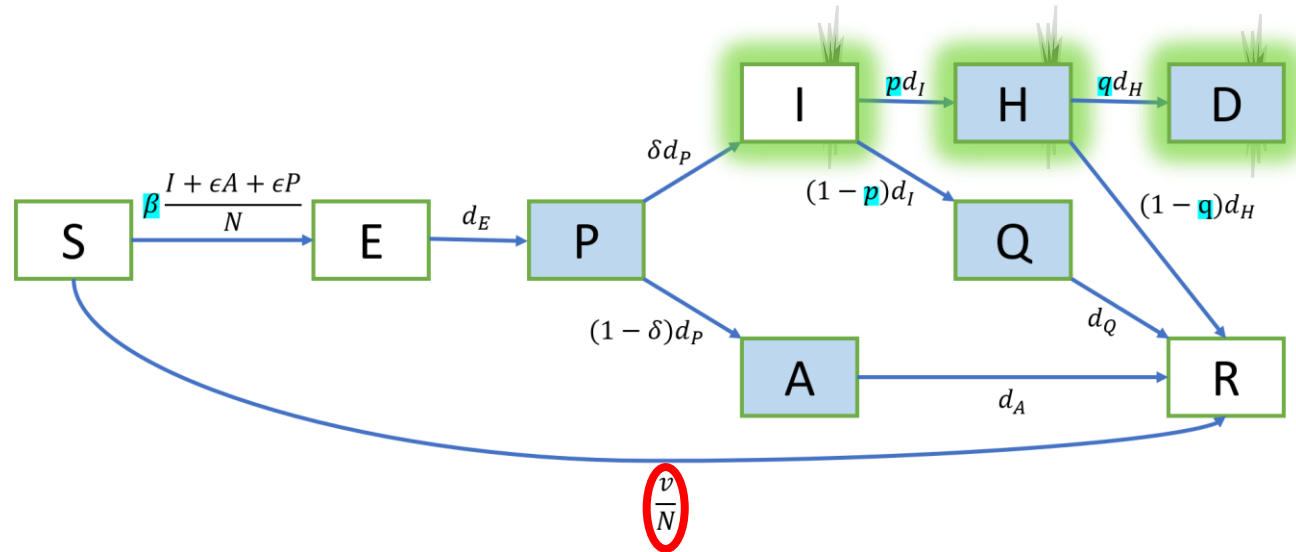
New York City COVID-19 related Event Timeline

Calibrate piecewise-constant model parameters to capture local epidemiological dynamics



Epidemiological Model Development

$$\left\{ \begin{array}{l} \frac{dS}{dt} = -\beta \frac{I + \epsilon A + \epsilon P}{N} S - \frac{v}{N} S \\ \frac{dE}{dt} = \beta \frac{I + \epsilon A + \epsilon P}{N} S - d_E E \\ \frac{dP}{dt} = d_E E - d_P P \\ \frac{dI}{dt} = \delta d_P P - d_I I \\ \frac{dA}{dt} = (1 - \delta) d_P P - d_A A \\ \frac{dH}{dt} = p d_I I - d_H H \\ \frac{dQ}{dt} = (1 - p) d_I I - d_Q Q \\ \frac{dD}{dt} = q d_H H \\ \frac{dR}{dt} = d_A A + (1 - q) d_H H + d_Q Q + \frac{v}{N} S. \end{array} \right.$$



Fixed
parameters:

eps = 0.75

delta = 0.6

d_E = 1/2.9

d_P = 1/2.3

d_I = 1/2.9

d_A = 1/7

d_H = 1/6.9

d_Q = 1/10

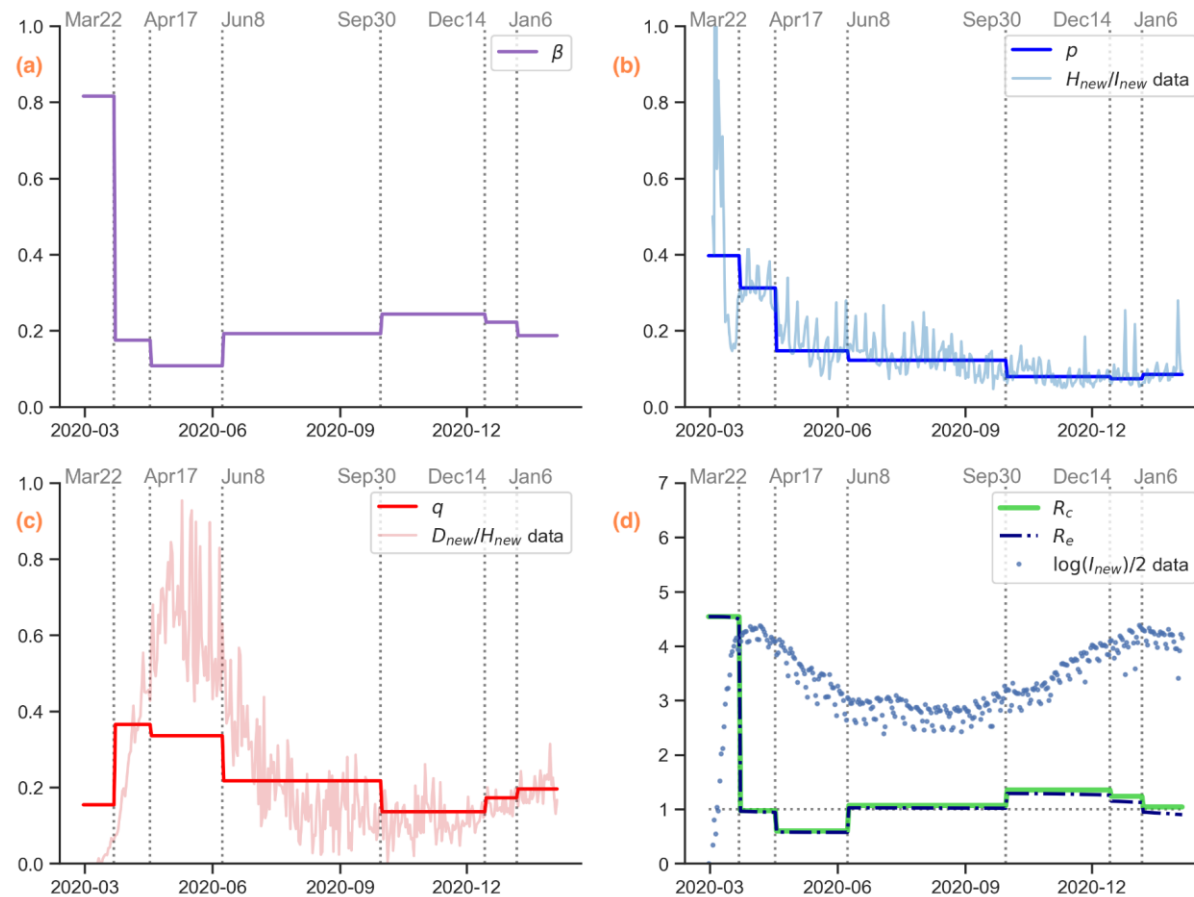
β : Transmission rate

p : Hospitalization rate

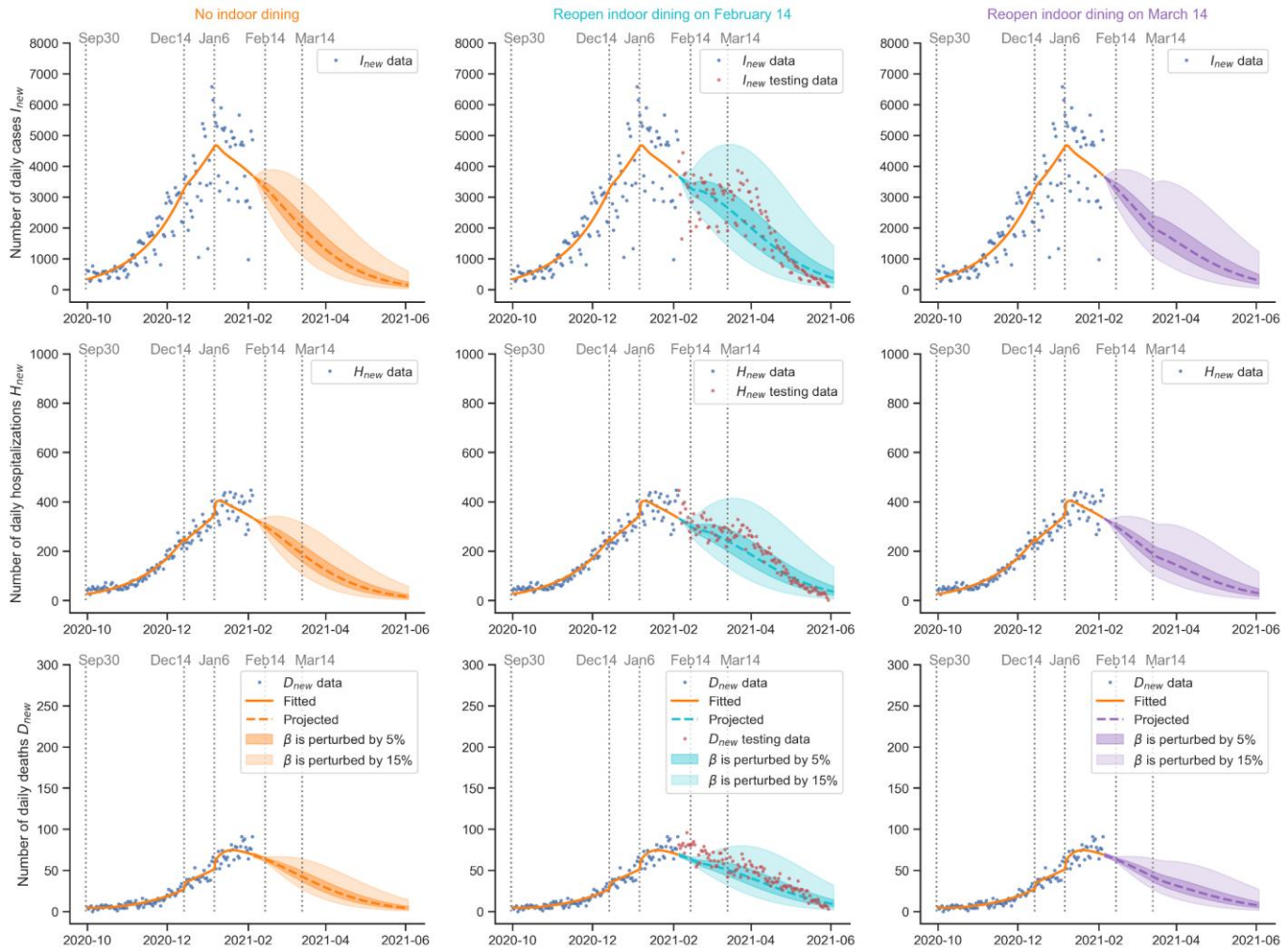
q : Death from hospital rate

Calibrated COVID-19 Transmission Rate for New York City

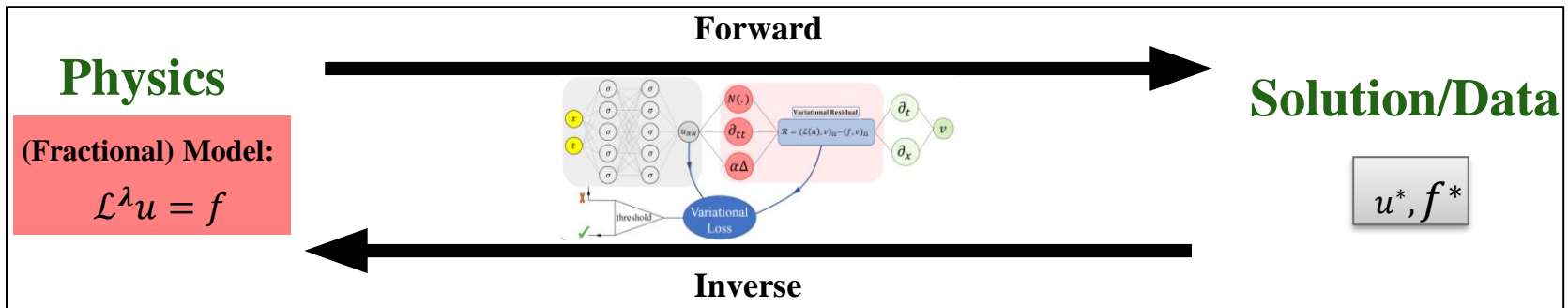
Calibrate piecewise-constant model parameters to capture local epidemiological dynamics



Forecasting with Uncertainties and Scenarios



Physics Informed Neural Networks (PINNs)



- \mathcal{L}^λ A (non-local) differential operator with parameters λ

A flexible **computational** tool to study **model uncertainty**

Incorporate **data** and **different models**

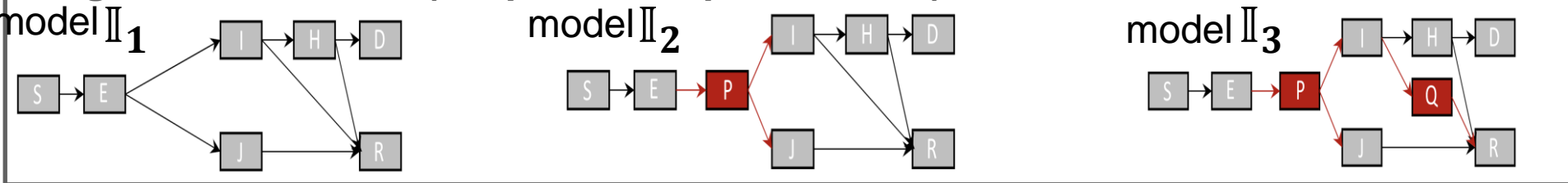
Accurate **fitting** to data

Inferring model parameters and **discovering** unobserved dynamics

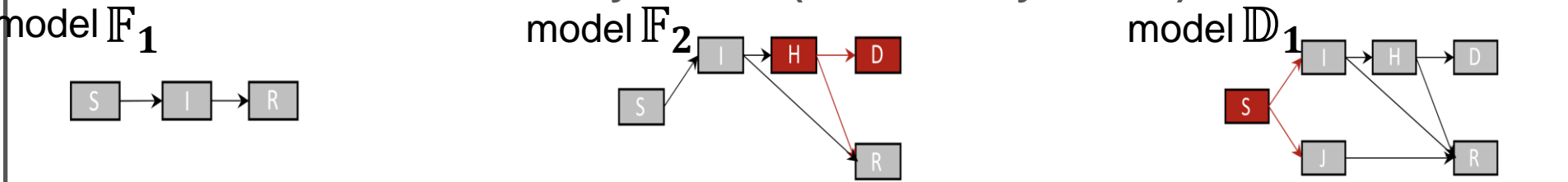
1. Ehsan Kharazmi, Min Cai, Xiaoning Zheng, Guang Lin, George Em Karniadakis, **Identifiability and predictability of integer- and fractional-order epidemiological models using physics-informed neural networks**, *Nature Computational Science*, 1, 744-753, 2021

Different Epidemiological Models

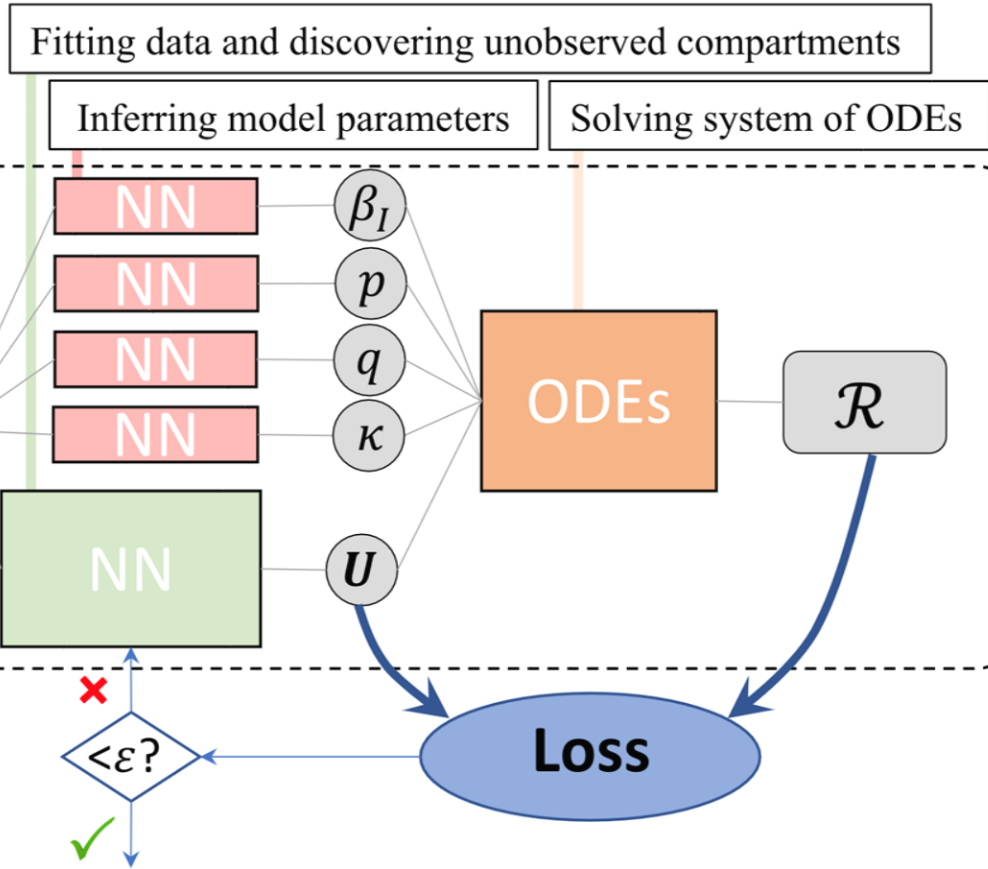
Integer-Order Models (simple to complex models)



Fractional-Order and Time-Delay Models (add memory effects)



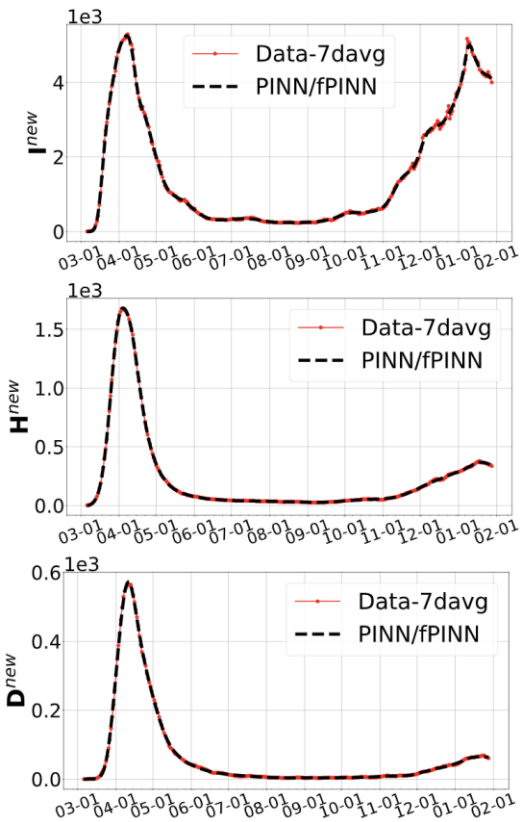
PINNs for (Fractional) Epidemiological Models



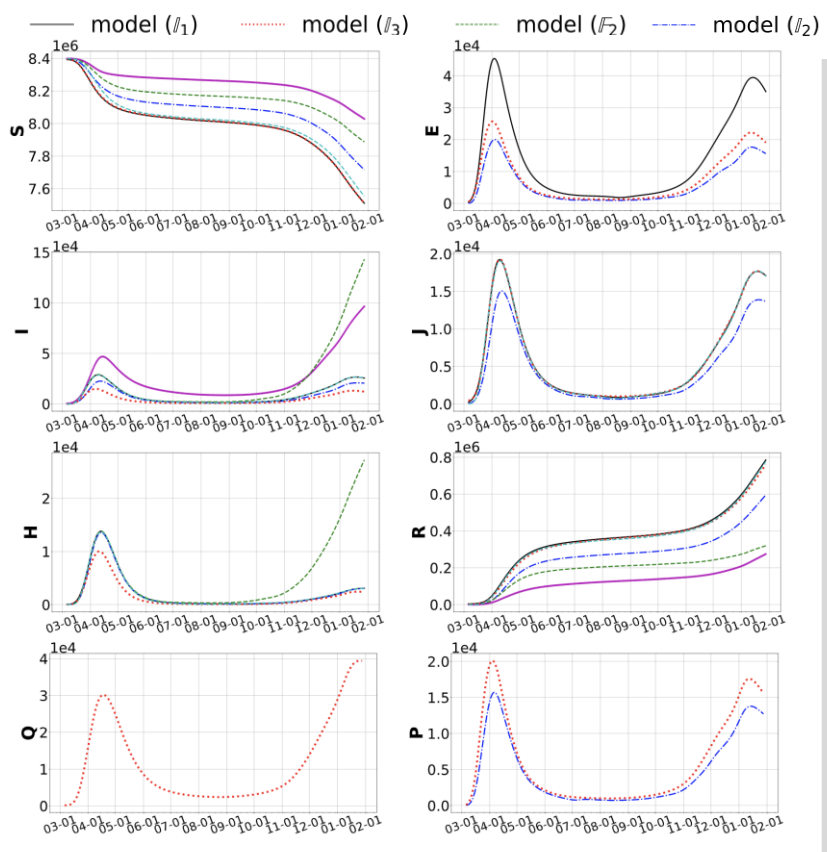
$$loss = \underbrace{\frac{1}{N_u} \sum_{i=1}^{N_u} |U(t_i; \theta) - data(t_i)|^2}_{\text{Loss } u \text{ (data/IC points)}} + \underbrace{\frac{1}{N_r} \sum_{j=1}^{N_r} |\mathcal{R}(t_j; \theta, \lambda)|^2}_{\text{Loss ODE (residual points)}}$$

PINN Results: Model Uncertainty based on NYC dataset

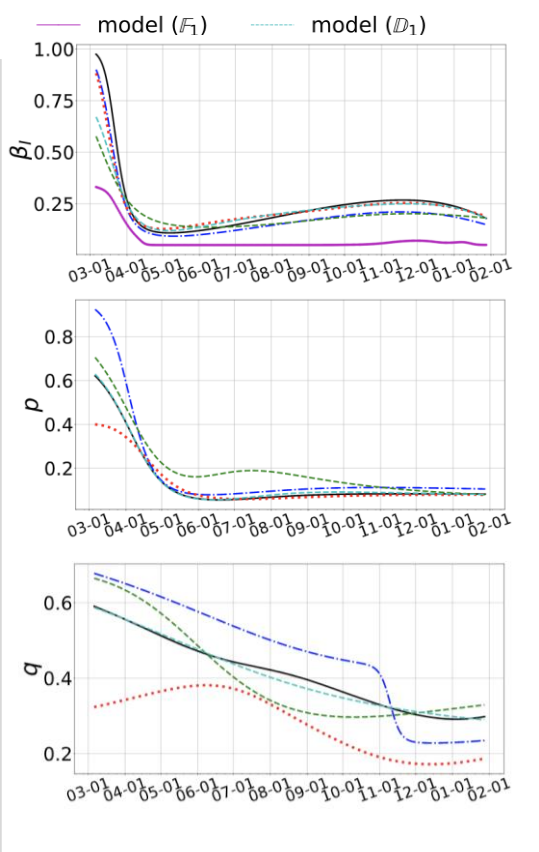
Fitting the data accurately



Discovering unobserved dynamics



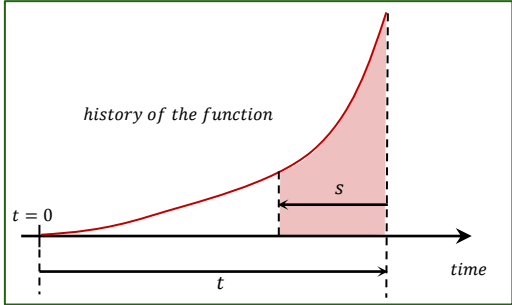
Inferring model parameters



Fractional Order Models Introduce Memory in the Dynamics

Caputo fractional derivative of order $\kappa \in (0,1)$: a **convolution** type **integro-differential** operator

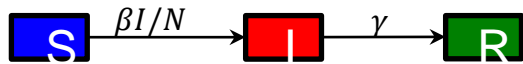
$$\frac{\partial^\kappa}{\partial t^\kappa} u(t) = {}^C_0\mathcal{D}_t^\kappa u(t) = \frac{1}{\Gamma(1-\kappa)} \int_0^t \frac{1}{(t-s)^\kappa} \frac{du(s)}{ds} ds$$



Memory: The derivative at time t depends on the weighted values of the function **from initial point $t = 0$ up to current time t .**

- Fractional order κ is the notion of memory effect
- Smaller κ can induce a delay in the dynamics
- $\kappa = \kappa(t)$ can be time varying

Different Compartments May Have Different Memory Effects! model \mathbb{F}_1

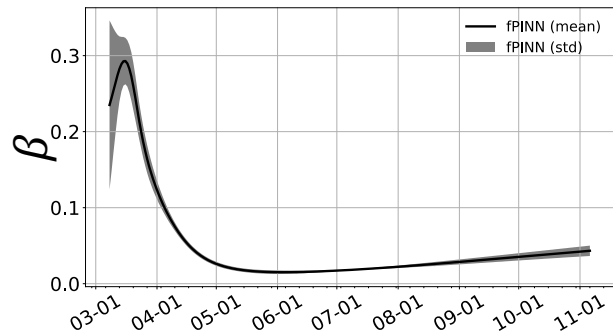
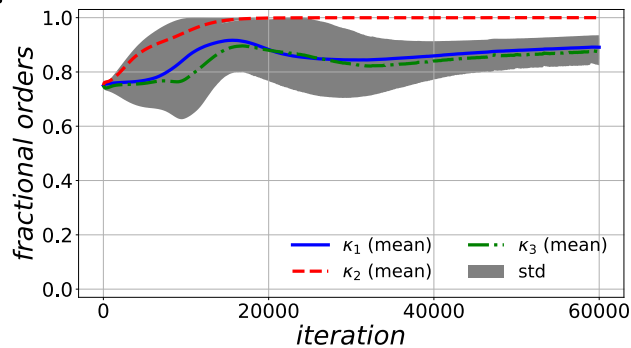


$${}^C_0\mathcal{D}_t^{\kappa_1} S(t) = -\frac{\beta}{N} I(t) S(t),$$

$${}^C_0\mathcal{D}_t^{\kappa_2} I(t) = \frac{\beta}{N} I(t) S(t) - \gamma I(t),$$

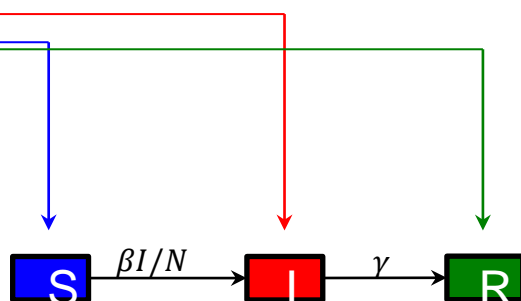
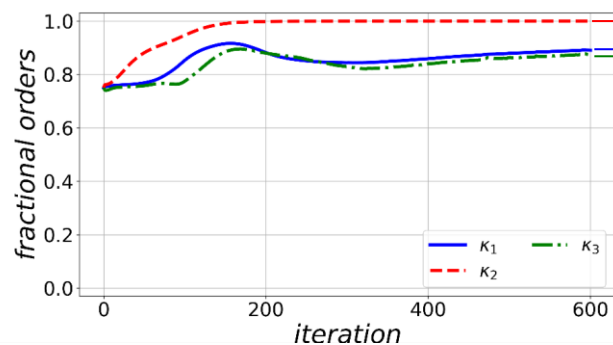
$${}^C_0\mathcal{D}_t^{\kappa_3} R(t) = \gamma I(t),$$

$${}^C_0\mathcal{D}_t^{\kappa_2} I^c(t) = \frac{\beta}{N} I(t) S(t).$$

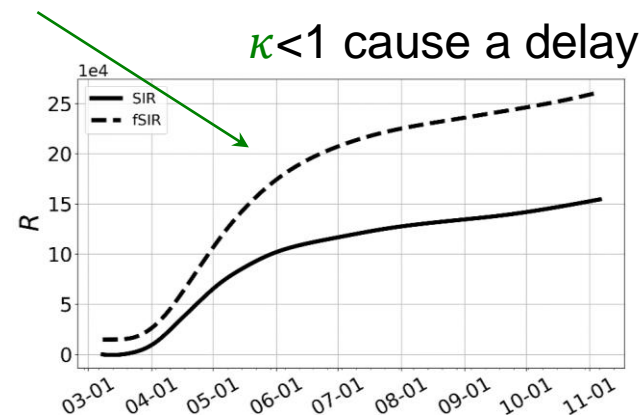
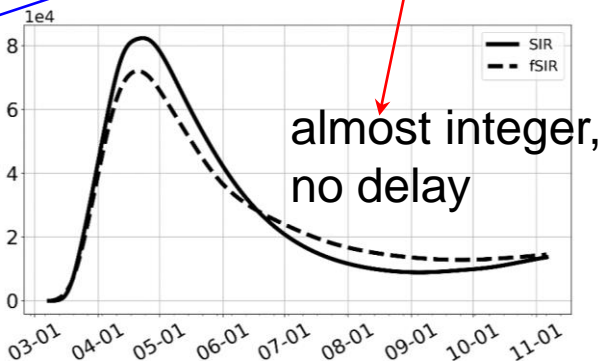
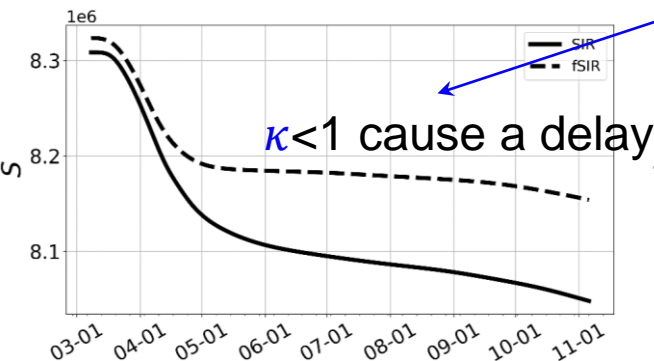


model \mathbb{F}_1

Fractional Order SIR V.S. Integer Order SIR



$$\kappa_1 = 0.89 \quad \kappa_2 = 0.99 \quad \kappa_3 = 0.87$$



Summary

This is the first work to employ structural and practical identifiability tools to study COVID-19 model identifiability based on the available data.

A general data-driven epidemiological modeling framework is developed, which seamlessly integrates model identifiability, model sensitivity analysis, model calibration, model prediction with confidence intervals, and evaluating control strategies under uncertainties.

We treat β (transmission rate), p (proportion of isolated individuals), and q (proportion of disease-related deaths) as time-dependent piece-wise model parameters and calibrate them using the available New York City COVID-19 dataset.

The developed COVID-19 model is employed to evaluate the effects of vaccination deployment scenarios.

We developed a flexible computational framework using physics-informed neural networks (PINNs) to study model uncertainty and discover time-dependent parameters.

Outline:

- ❖ Incorporate Physics Knowledge and AI to design new interpretable models
- ❖ **Incorporate Physics Knowledge into AI to predict multiscale problems: NH-PINN**
- ❖ Interpretable AI enables data-driven scientific discovery with uncertainty quantification capability – ALZHEIMER's Disease Prediction
- ❖ Sparse Neural Architecture Design with quantified uncertainties
- ❖ Scalable training large-scale Deep Neural Network

NH-PINN: Neural homogenizationbased the physics-informed neural network for the multiscale problems

Wing Tat Leung, Guang Lin, Zecheng Zhang, NH-PINN: Neural homogenization based Physics-informed Neural Network for Multiscale Problems, 2022, <https://arxiv.org/abs/2108.12942>

Content

- ➊ Physics-informed neural network (PINN)
- ➋ Homogenization
- ➌ Neural homogenization based PINN (NH-PINN)
- ➍ Numerical examples

Wing Tat Leung, Guang Lin, Zecheng Zhang. NH-PINN: Neural homogenization based the physics-informed neural network for the multiscale problems.
arXiv:2108.12942.

Physics-informed neural network (PINN)

$$\begin{aligned}\mathcal{L}(u) &= f \text{ in } \Omega \\ \mathcal{B}(u) &= b \text{ on } \partial\Omega\end{aligned}$$

where \mathcal{L} is a differential operator and \mathcal{B} is the boundary condition operator. f is the given source term, b is the given boundary condition.

$$\min_{\beta} \frac{w_1}{N_f} \sum_{i=1}^{N_f} |\mathcal{L}(\mathcal{F}_{\beta}(p_i)) - f(p_i)|^2 + \frac{w_2}{N_b} \sum_{i=1}^{N_b} |\mathcal{B}(\mathcal{F}_{\beta}(q_i)) - b(q_i)|^2, \quad (1)$$

where $w_1 + w_2 = 1$ are the positive weights; $\{p_i\} \subset \Omega$, $\{q_i\} \subset \partial\Omega$ and N_f, N_b are the number of points used in discretizing the domain and boundary respectively, $\mathcal{F}_{\beta}(\cdot)$ is the network and β is the parameters associated with the network.

Motivation

$$\begin{aligned} -\nabla \cdot (\kappa \nabla u(x)) &= f, x \in \Omega, \\ u &= 0, x \in \partial\Omega, \end{aligned}$$

where $\Omega = [0, \pi]$ and $f = \sin(x)$, $\kappa(x) = 0.5 \sin(2\pi x/\epsilon) + 2$, where $\epsilon = \frac{1}{8}$.

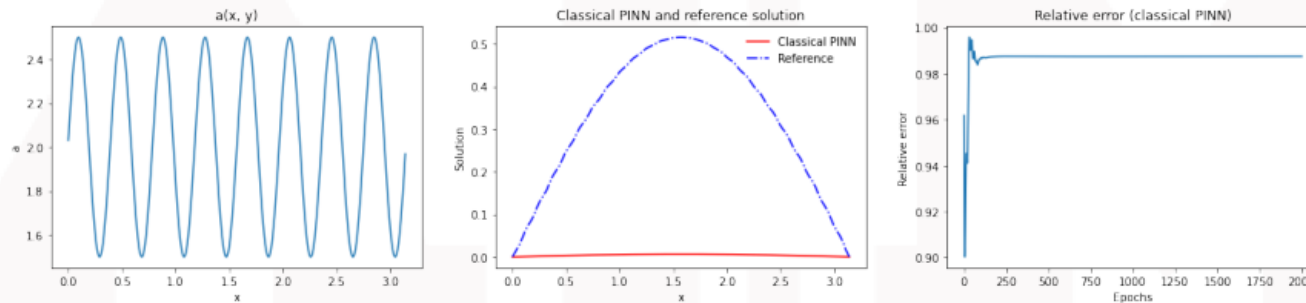


Figure: Left: demonstration of the permeability $\kappa(x) = 0.5 \sin(2\pi x/\epsilon) + 2$, where $\epsilon = 1/8$. Middle: learnt solution by the classical PINN vs the reference solution. Right: relative error as a function of the training epochs. The average error of the last 500 epochs is 0.987471.

Homogenization

Consider we are solving:

$$-\frac{\partial}{\partial x_i} \left(a_{ij}(x/\epsilon) \frac{\partial}{\partial x_j} u_\epsilon(x) \right) = f(x), x \in \Omega \quad (2)$$

with $u_\epsilon(x) = 0$ on $\partial\Omega$. The asymptotic expansion:

$$u_\epsilon = u_0 + \epsilon u_1 + \epsilon^2 u_2 + \dots \quad (3)$$

① Solve the cell problem:

$$-\frac{\partial}{\partial y_i} \left(a_{ij}(y) \frac{\partial}{\partial y_j} \right) \chi_j = \frac{\partial}{\partial y_i} a_{ij}(y), \quad (4)$$

$$\chi_j \text{ is periodic in } y \text{ with mean } 0. \quad (5)$$

② Solve the homogenized equation:

$$-\frac{\partial}{\partial x_i} \left(a_{ij}^* \frac{\partial}{\partial x_j} \right) u_0 = f, \quad (6)$$

$$\text{where } a_{ij}^* = \int_Y \left(a_{ij} + a_{ik} \frac{\chi_j}{\partial y_k} \right) dy$$

Neural homogenization based PINN (NH-PINN)

Proposed NH-PINN:

- 1 Solve the cell problems using PINN.
- 2 Evaluate the homogenized coefficients.
- 3 Solve the homogenized equation using PINN.

Oversampling (to improve PINN accuracy for solving periodic problems):

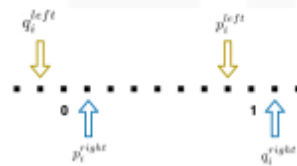


Figure: Oversampling mesh demonstration.

Suppose $w_1 + w_2 + w_3 = 1$ are positive constants, the new loss function:

$$\min_{\beta} \left\{ \frac{w_1}{N_f} \sum_{i=1}^{N_f} |\mathcal{L}(\mathcal{F}_{\beta}(p_i)) - f(p_i)|^2 + w_2 |(\mathcal{F}_{\beta}(q_1) - \mathcal{F}_{\beta}(q_2))|^2 \right. \\ \left. + \frac{w_3}{N_o} \sum_{i=1}^{N_o} \left(|(\mathcal{F}_{\beta}(q_i^{left}) - \mathcal{F}_{\beta}(p_i^{left}))|^2 + |(\mathcal{F}_{\beta}(q_i^{right}) - \mathcal{F}_{\beta}(p_i^{right}))|^2 \right) \right\}.$$

Notations for the numerical examples

Solution notation	Cell problems solver	Homogenized equation solver
$p(x)$	PINN	PINN
$v(x)$	FEM	FEM
$w(x)$	PINN	FEM

Table: Notations of the solutions. u_ϵ which is not listed here is the reference solution which solves the multiscale PDE directly by fine scale finite element methods.

Relative errors:

$$e_1 = \frac{\|p(x) - u_\epsilon(x)\|}{\|u_\epsilon\|}, e_2 = \frac{\|w(x) - u_\epsilon(x)\|}{\|u_\epsilon(x)\|},$$
$$e_3 = \frac{\|p(x) - w(x)\|}{\|w(x)\|}, e_4 = \frac{\|v(x) - u_\epsilon(x)\|}{\|u_\epsilon\|},$$

where $\|\cdot\|$ is the L_2 norm.

Numerical examples

We consider the following 2D elliptic equation:

$$-\frac{\partial}{\partial x_i} \left(a\left(\frac{x}{\epsilon}\right) \frac{\partial}{\partial x_i} u_\epsilon(x) \right) = f(x), x \in \Omega, \quad (7)$$

$$u_\epsilon(x) = 0, x \in \partial\Omega. \quad (8)$$

In our examples, $\Omega = [0, 1]^2$ and the permeability $a(x/\epsilon) = 2 + \sin(2\pi x_1/\epsilon) \cos(2\pi x_2/\epsilon)$ and $\epsilon = \frac{1}{8}$.

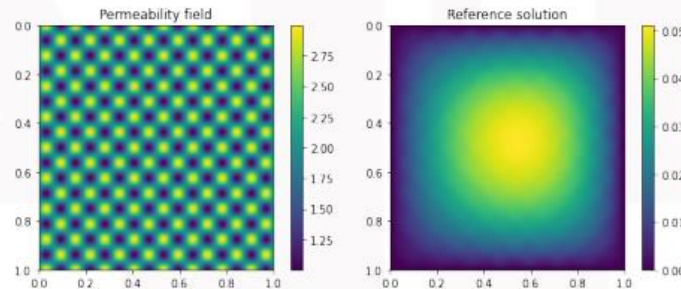


Figure: Left: permeability of the 2D elliptic problem. Note that $\epsilon = 1/8$. Right: the reference solution.

Oversampling and cell problems

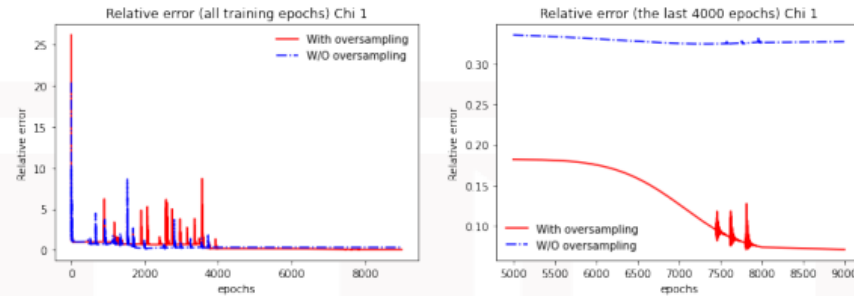


Figure: $2D$ elliptic cell problem χ_1 . Relative error as a function of the training epochs for χ_1 . Left: all training epochs; right: the last 4000 epochs. The average relative errors of the last 300 epochs for the oversampling and without oversampling are 0.072034 and 0.326963 respectively.

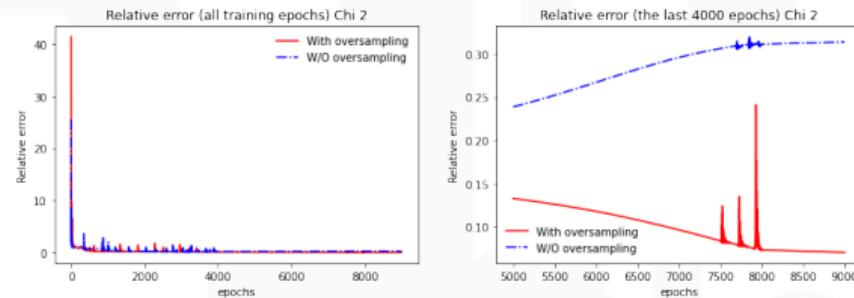


Figure: $2D$ elliptic cell problem χ_2 . Relative error as a function of the training epochs for χ_2 . Left: all training epochs; right: the last 4000 epochs. The average relative errors of the last 300 epochs for the oversampling and without oversampling are 0.071396 and 0.312961 respectively.

Convergence results

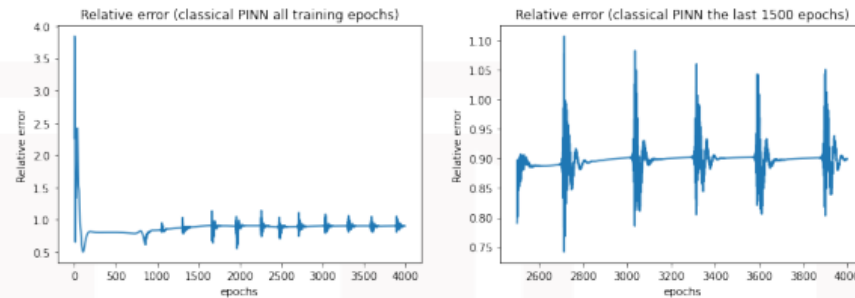


Figure: 2D elliptic problem by the classical PINN. Relative error as a function of the training epochs, the entire history (left), the last 1500 epochs (right).

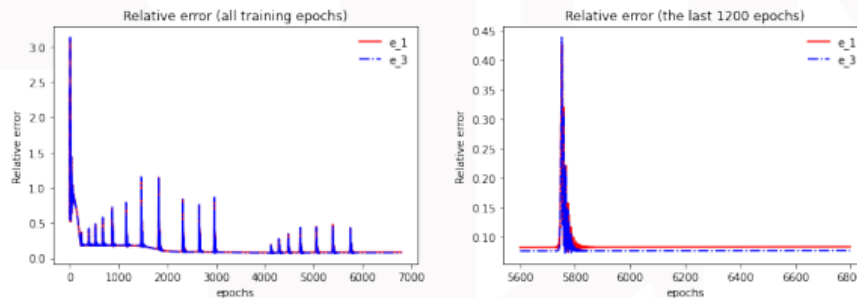


Figure: 2D elliptic cell problem e_1 and e_3 relative errors (NH-PINN) with respect to the training epochs. Left: history of all training epochs; right: history of the last 1200 epochs.

Analysis and transfer learning

Errors:

e_1	e_2	e_3	e_4
0.082994	0.0212534	0.076604	0.021316

Table: Relative errors for the 2D elliptic problem.

Transfer learning: We use the trained NH-PINN network to initialize the PINN. The failure shows that the minimizers of two loss functions are not closed. Modify the PINN loss function should be the future.

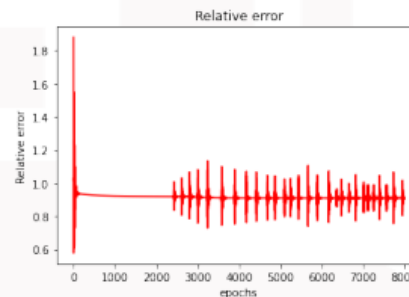


Figure: 2D elliptic transfer learning of classical PINN.

Conclusion

- ① We find that PINN accuracy on solving multiscale problems degenerates. We propose a 3-step approach, neural homogenization based PINN (NH-PINN). We first apply PINN to solve the cell problems which are used to derive the homogenized equation; the homogenized equation can then be easily solved by PINN.
- ② We propose an oversampling strategy to solve the periodic PDE by PINN; this method greatly improves the accuracy of the PINN when solving the high dimensional periodic problems.
- ③ We also observe that NH-PINN can improve the homogenization accuracy. If we apply PINN to implement homogenization, the solution may be more accurate than the traditional numerical methods. PINN may be a potential alternative of implementing the homogenization.

Outline:

- ❖ Incorporate Physics Knowledge and AI to design new interpretable models
- ❖ Incorporate Physics Knowledge into AI to predict multiscale problems: NH-PINN
- ❖ **Interpretable AI enables data-driven scientific discovery with uncertainty quantification capability – ALZHEIMER's Disease Prediction**
- ❖ Sparse Neural Architecture Design with quantified uncertainties
- ❖ Scalable training large-scale Deep Neural Network

Interpretable AI:

Question: Can we use available observation data to discover the physical laws?

Goal: Enable Data-driven Scientific Discovery?

S. Zhang, **G. Lin**, Robust data-driven discovery of governing physical laws with error bars, Proceedings of the Royal Society of London. Series A, mathematical, physical and engineering sciences, in press, 2018.

Jiuhai Chen, Lulu Kang, Guang Lin, Gaussian process assisted active learning of physical laws, Technometrics, in press, 2020.

<https://doi.org/10.1080/00401706.2020.1817790>

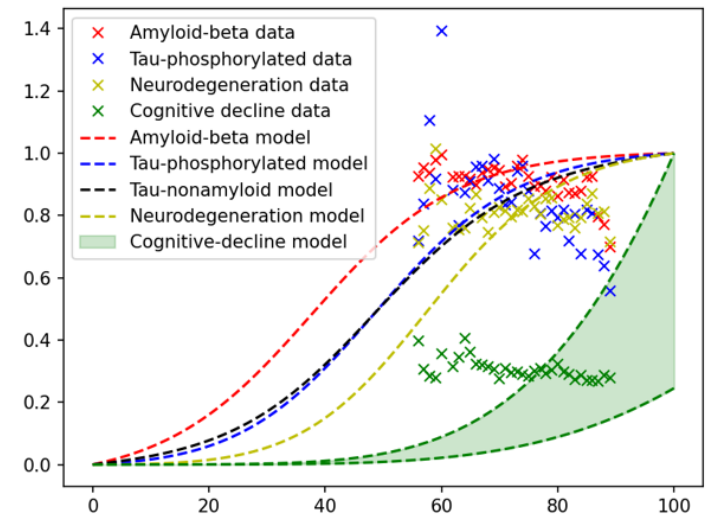
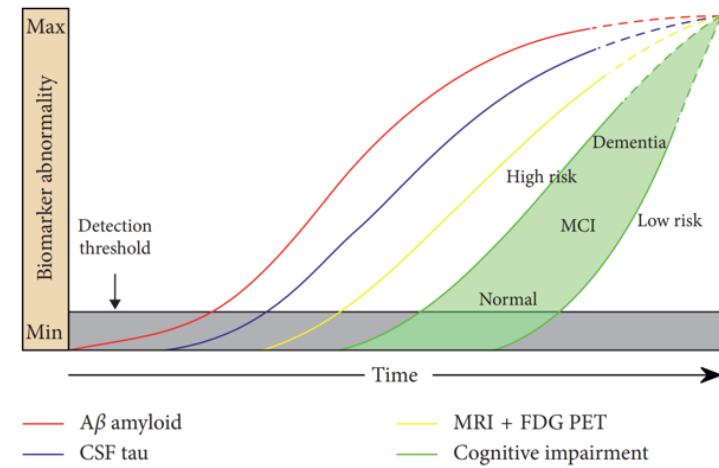
Sheng Zhang, Guang Lin, Robust subsampling-based threshold sparse Bayesian regression to tackle high noise and outliers for data-driven discovery of differential equations, Journal of Computational Physics, 428: 109962, 2021.

ALZHEIMER'S DISEASE PREDICTION

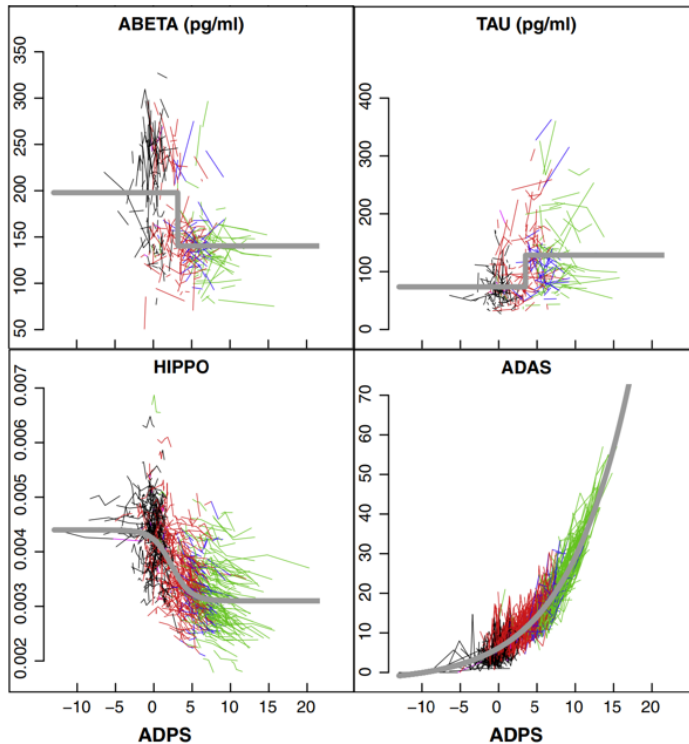
Haoyang Zheng, Jeffrey Petrella, P.Murali Doraiswamy, Guang Lin, Wenrui Hao, Nature Medicine, in review, 2022

AD prediction

How to apply patient data (ADNI dataset) to optimize ODEs?



Dataset

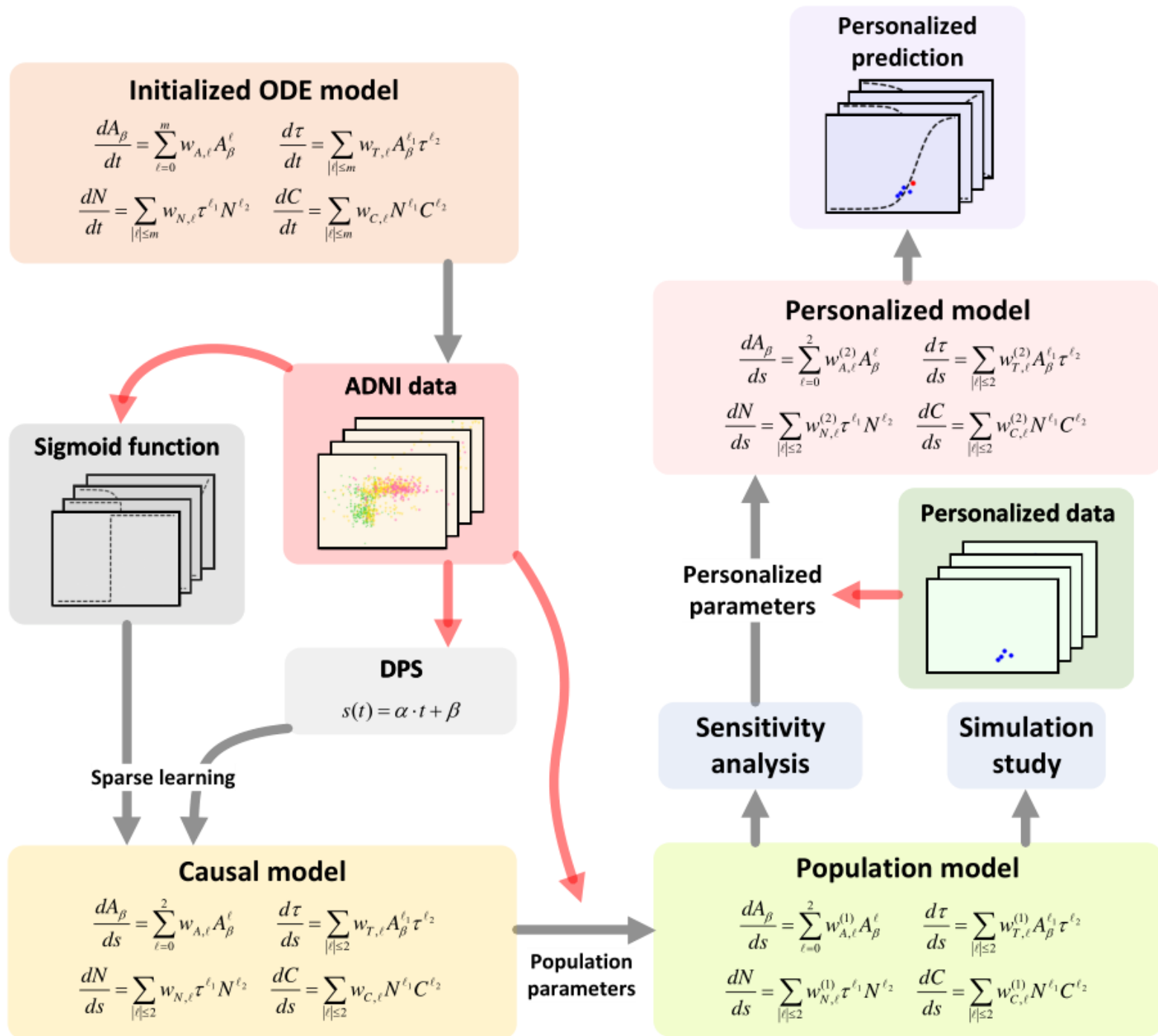


Progression of ADNI biomarkers
as function of the Alzheimer's
Disease Progression Score (ADPS)

— Normal–Normal
— MCI–MCI
— AD–AD
— MCI–AD
— MCI–N
— sigmoids

Disease progression score
(DPS)

- Patients' disease progression differs in their age of onset and rate of progression.
- Each biomarker follows a sigmoid shaped curve.



Fit ODEs with the use of patient data

Results

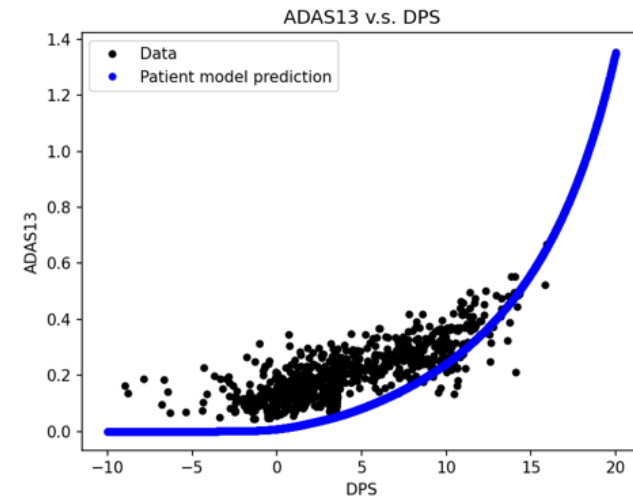
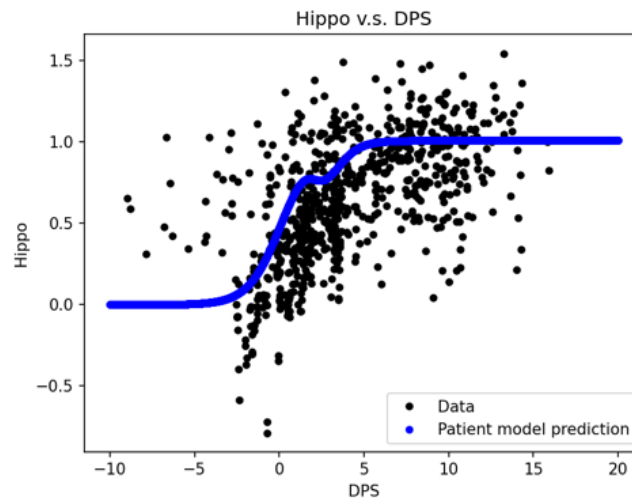
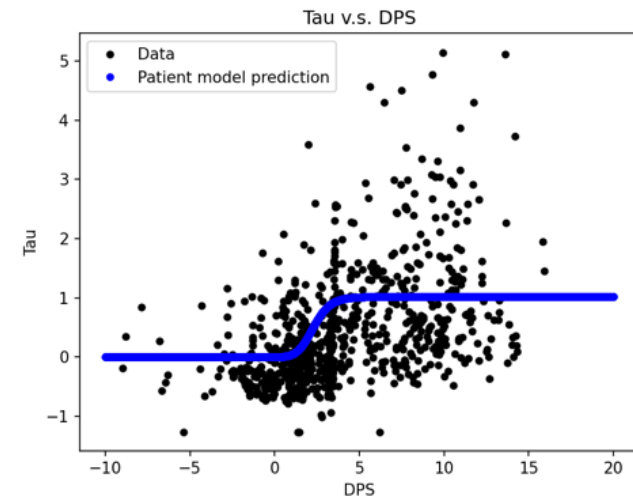
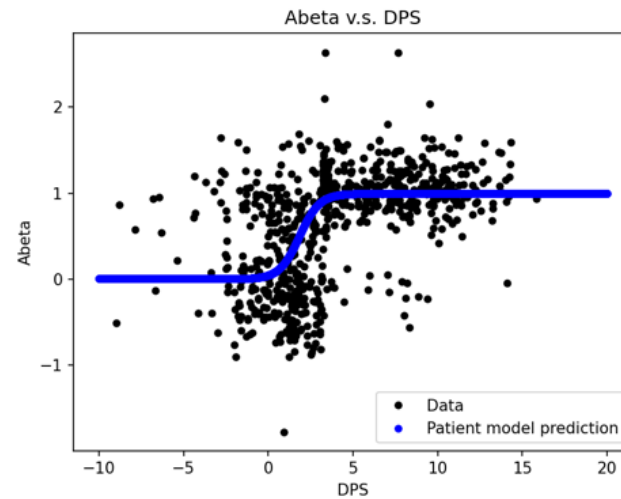
$$\frac{dA_{\beta}(t)}{dt} = -1.7112 \cdot (A_{\beta}(t) + 0.005) \cdot (A_{\beta}(t) - 1) + 0.0045;$$

$$\frac{d\tau_{\rho}(t)}{dt} = -0.0276 \cdot A_{\beta}(t) + 3.6918 \cdot A_{\beta}^2(t) - 0.2498 \cdot \tau_{\rho}(t)[\tau_{\rho}(t) - 1 + 14.2932 \cdot A_{\beta}(t)];$$

$$\frac{dN(t)}{dt} = -0.5106 \cdot \tau_{\rho}(t) + 1.0708 \cdot \tau_{\rho}^2(t) - 1.0257 \cdot N(t)[N(t) - 1 + 0.5533 \cdot \tau_{\rho}(t)]$$

$$\frac{dC(t)}{dt} = -0.0017 \cdot N(t) + 5.9 \cdot 10^{-7} \cdot N^2(t) - 0.0372 \cdot C(t)[C(t) - 2 - 0.0118 \cdot N(t)]$$

Results



Outline:

- ❖ Incorporate Physics Knowledge and AI to design new interpretable models
- ❖ Incorporate Physics Knowledge into AI to predict multiscale problems: NH-PINN
- ❖ Interpretable AI enables data-driven scientific discovery with uncertainty quantification capability – ALZHEIMER's Disease Prediction
- ❖ **Sparse Neural Architecture Design with quantified uncertainties**
- ❖ Scalable training large-scale Deep Neural Network

Sparse Neural Architecture Design with quantified uncertainties:

Question: Can we develop a fast, small & accurate deep neural network with better interpretability and less demanding on the computational resource?

Goal: Enable Fast Interpretable Nonlinear Data-driven Scientific Discovery.

W. Deng, X. Zhang, F. Liang, **G. Lin**, An adaptive empirical Bayesian method for sparse deep learning, **2019 Conference on Neural Information Processing Systems (NIPS)**, Dec. 8 – Dec. 14, 2019, Vancouver, Canada.

NeurIPS'19, NeurIPS'20, ICML'20, ICLR'21, JCP'20, JCP'21a, JCP'21b

Bayesian Sparse learning with preconditioned stochastic gradient MCMC and its applications

Guang Lin¹

Joint work with Yating Wang, Wei Deng, Xiao Zhang, Faming Liang

¹Departments of Mathematics, Statistics & School of Mechanical Engineering

Purdue University

NeurIPS'19, JCP'20

Outline

- 1 Introduction
- 2 A New Class of Adaptive Stochastic Gradient MCMC
- 3 Preconditioned SGLD and hierarchical Bayesian model
- 4 Numerical tests
- 5 Summary

Motivation and Objective

DNN challenges:

- Over-parameterized DNN requires heavy memory and computation power, and may cause overfitting
- Cost function is difficult to optimize and a good local minima is hard to obtain

Objective:

- Sparse learning to enhance efficiency, robustness and interpretability
- Bayesian approach to capture uncertainty
- Compute posterior expectation to obtain more robust result
- Escape "shallow" local optimas and saddle points to achieve better point estimate

Notations

- The entire data: $\mathcal{D} = \{d_i\}_{i=1}^N$, a mini-batch of data \mathcal{B}
- Model parameters: $\beta \in \mathbb{R}^D$
- Log posterior density $L(\beta) = \log(p(\beta|\mathcal{D}))$ and true gradient

$$\nabla_{\beta} L(\beta) = \nabla_{\beta} \log p(\beta) + \sum_{i=1}^N \nabla_{\beta} \log p(d_i|\beta)$$

- Stochastic gradient using a mini-batch is $\nabla \tilde{L}(\beta)$.

$$\nabla_{\beta} \tilde{L}(\beta) = \nabla_{\beta} \log p(\beta) + \frac{N}{n} \sum_{i=1}^n \nabla_{\beta} \log p(d_i|\beta)$$

Key Components:

- **Sparse learning** - Enhance efficiency, robustness and interpretability
- **Bayesian approach** - Capture uncertainty
- **Empirical Bayesian method** - Learn a class of hierarchical Bayes models, yield data-driven adaptive penalties
- **Adaptive Stochastic Gradient Langevin Dynamics (SGLD) or SG-MCMC** - Capture parameter uncertainty and avoid overfitting, escaping "shallow" local optima and saddle points, and compute posterior expectation to obtain more robust result
- **Preconditioned SGLD (PSGLD)** - Update parameters with different step size, adaptive to local geometric and resulting in faster convergence for components of β have different scales
- **Stochastic Approximation (SA)** - Optimize latent variables in prior (SGLD-SA, PSGLD-SA) to converge to the asymptotically correct distribution with a controllable bias
- **Pruning Strategy** - Enable Sparse Deep Neural Network; Sparsity is ensured, resulting in less usage in memory and computational power

SGD

SGD update the parameters using

$$\beta_{k+1} = \beta_k + \epsilon_k \nabla_{\beta} \tilde{L}(\beta_k)$$

- Find MAP for model parameter through stochastic optimization
- Do not capture parameter uncertainty and can potentially overfit data.

Langevin Dynamics and SGLD

- A Langevin diffusion with stationary distribution $p(\beta)$ can be described by the SDE

$$d\beta(t) = \nabla_{\beta} L(\beta) dt + \sqrt{2} dW(t) \quad (1)$$

where $W(t)$ is a Brownian motion.

- Euler discretization of (1), and approximate the true gradient by the stochastic gradient:

$$\beta_{k+1} = \beta_k + \epsilon_k \nabla_{\beta} \tilde{L}(\beta_k) + N(0, 2\epsilon_k \tau^{-1} I) \quad (2)$$

(2) asymptotically converges to $\pi(\beta|\mathcal{D}) \propto e^{\tau L(\beta)}$. As τ increases, ϵ_k decreases, the solution tends to the global optima with a higher probability.

Advantages of SGLD

- Sample DNN posterior to model uncertainty
- Compute posterior expectation to obtain more robust result
- Escape local optimas to achieve better point estimate

SGLD updates all parameters with the same step size, this may cause slow mixing when components of β have different scales.

A class of Adaptive Stochastic Gradient MCMC

Stochastic gradient Langevin Dynamics (SGLD), the first order SG-MCMC algorithm, is a sampling algorithm in DNN which asymptotically converges to a stationary distribution of $e^{\tau L(\beta)}$:

$$\beta^{(k+1)} = \beta^{(k)} + \epsilon^{(k)} \nabla_{\beta} \tilde{L}(\beta^{(k)}) + \mathcal{N}(0, 2\epsilon^{(k)}\tau^{-1}), \quad (4)$$

For the adaptive SGLD, we introduce an auxiliary variable θ and formulate an inhomogeneous Markov Chain as follows:

$$\begin{aligned} \beta^{(k+1)} &= \beta^{(k)} + \epsilon^{(k)} \nabla_{\beta} \tilde{L}(\beta^{(k)}, \theta^{(k)}) + \tilde{\eta}^{(k)}, \\ \theta^{(k+1)} &= \theta^{(k)} + \omega^{(k+1)} H(\theta^{(k)}, \beta^{(k+1)}), \end{aligned} \quad (5)$$

where $\tilde{\eta}^{(k)} \sim \mathcal{N}(0, 2\epsilon^{(k)}/\tau \mathbf{I})$, $H(\theta, \beta)$ is a **biased but asymptotically unbiased estimator** of the mean field function $h(\theta)$. The interpretation of this algorithm is that we sample $\beta^{(k+1)}$ from $\tilde{L}(\beta^{(k)}, \theta^{(k)})$ and adaptively optimize θ such that $\langle \theta - \theta^*, h(\theta) \rangle = 0$.

Pruning Strategy - Enable Sparse Deep Neural Network

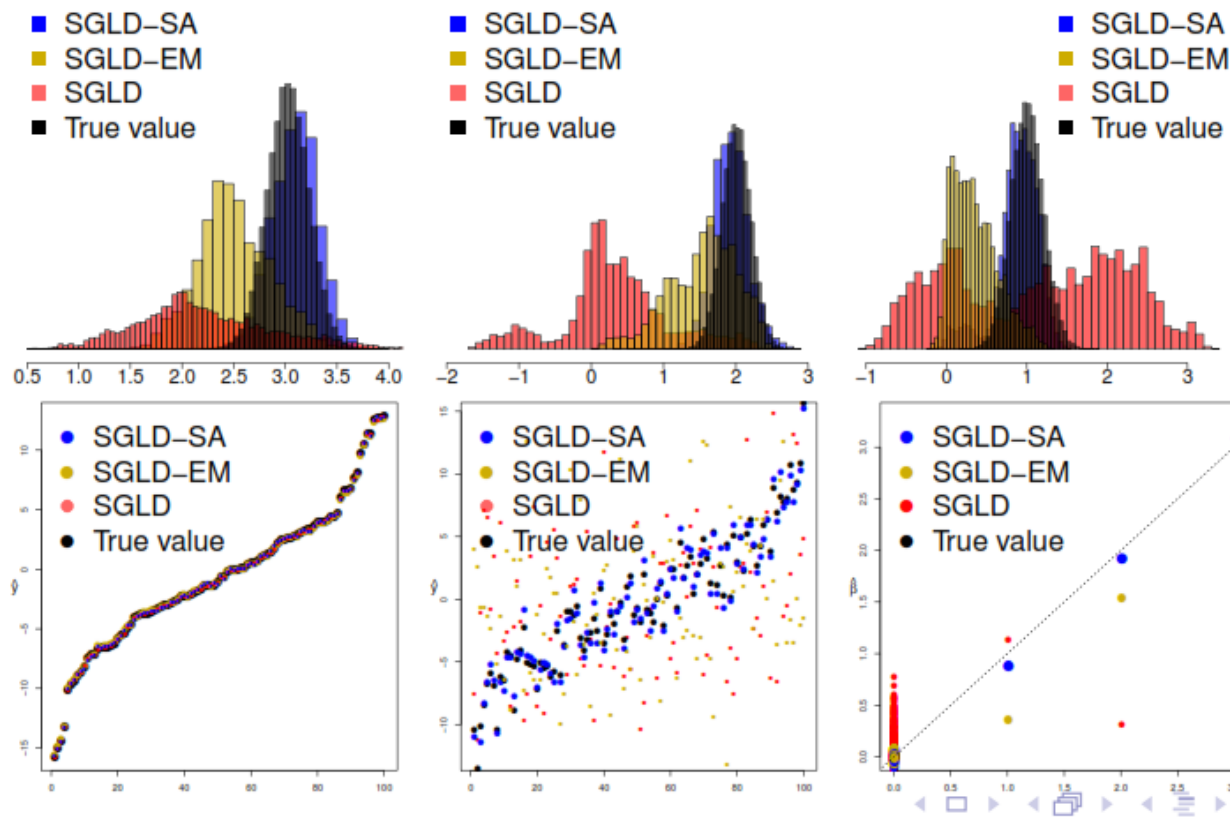
Although the magnitude-based unit pruning shows more computational savings, it doesn't demonstrate robustness under coarser pruning. We instead apply the magnitude-based weight-pruning to our compression experiments. The weight pruning can also be viewed as a greedy L^0 regularization in optimization.

Preconditioned SGLD (PSGLD)

- Updating parameters with different step size, adaptive to local geometric and resulting in faster convergence for components of β have different scales

Simulation of Large-p-Small-n Regression

Dataset: $n = 100$ and $p = 1000$. $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\eta}$ where $\mathbf{X} \sim \mathcal{N}_p(\mathbf{0}, \boldsymbol{\Sigma})$,
 $\boldsymbol{\beta} = (\beta_1, \beta_2, \beta_3, 0, 0, \dots, 0)'$, $\boldsymbol{\eta} \sim \mathcal{N}_n(\mathbf{0}, 3\mathbf{I}_n)$, $\beta_1 \sim \mathcal{N}(3, \sigma_c^2)$,
 $\beta_2 \sim \mathcal{N}(2, \sigma_c^2)$, $\beta_3 \sim \mathcal{N}(1, \sigma_c^2)$, $\sigma_c = 0.2$.



Residual Network Compression

Table 1: Resnet20 Compression on CIFAR10. $L_2 = 1 \times 10^{-4}$ means we apply weight decay 1×10^{-4} to the sparse layers with target sparse rate \mathbb{S} .

PENALTY \ \mathbb{S}	30%	50%	70%	90%
$L_1 = 1 \times 10^{-3}$	92.88	92.75	92.62	89.95
$L_1 = 1 \times 10^{-2}$	89.50	89.79	90.07	89.83
$L_2 = 1 \times 10^{-2}$	94.17	93.82	93.17	90.11
$L_2 = 1 \times 10^{-1}$	94.02	93.96	93.50	91.20
SGHMC-SA	94.23	94.27	93.74	91.68

Most notably, **91.68% accuracy based on 27K parameters (90% sparsity) in Resnet20 is the besting existing result.** By contrast, targeted dropout (2018) achieved 91.48% accuracy based on 47K parameters (90% sparsity) of Resnet32, BC-GHS (2017) achieved 91.0% accuracy based on 8M parameters (94.5% sparsity) of VGG models.

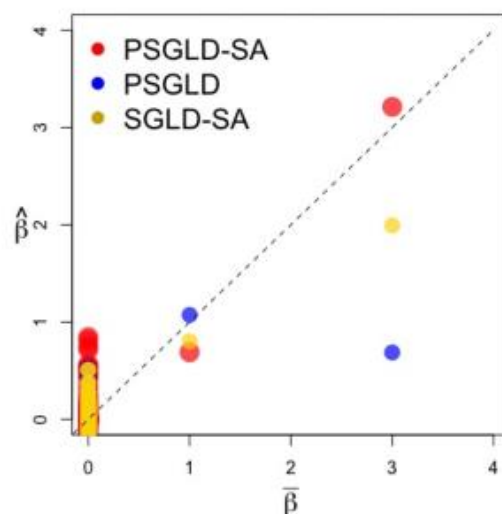
Small n large p problem

Number of observations $n = 100$, Number of predictors $p = 200$.

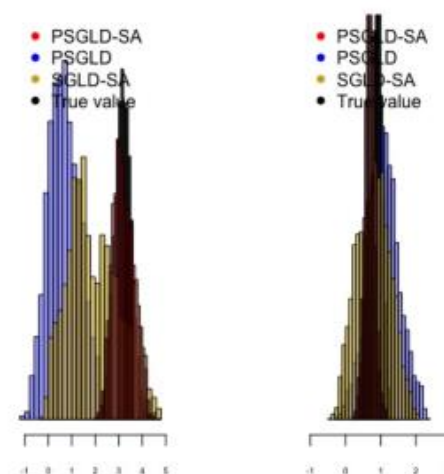
Predictors: $X \in \mathcal{N}(0, \Sigma)$, $\Sigma_{ij} = 0.6^{|i-j|}$, $X[:, 1] * 0.3$.

Model parameters: $\beta_1 = 3, \beta_2 = 1, \beta_j = 0$, for $j = 1, \dots, p$.

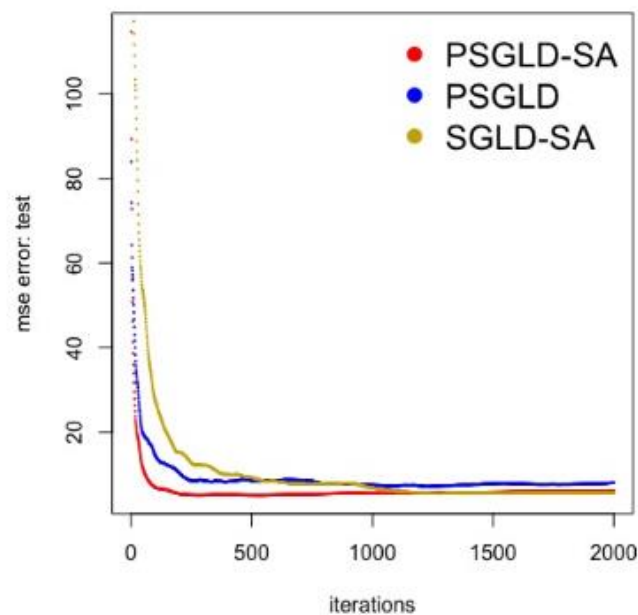
Responses: $y = X\beta + \epsilon$, $\epsilon \sim \mathcal{N}_n(0, 3I_n)$.



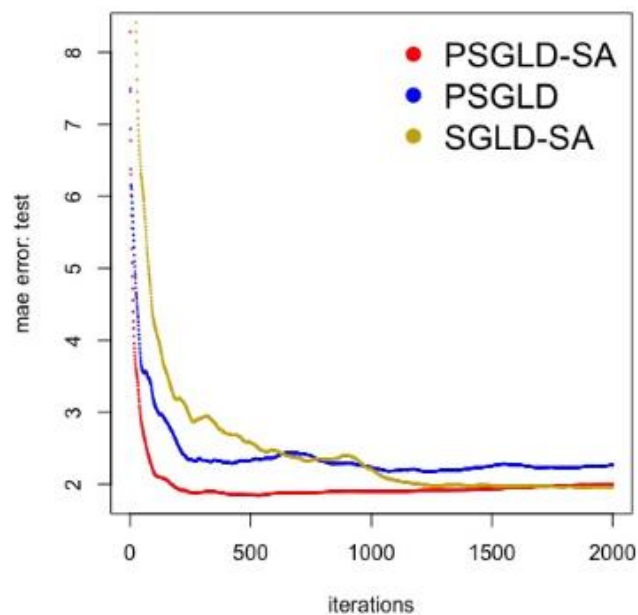
(a) Posterior mean vs true



(b) Posterior estimation of β_1, β_2



(c) Testing MSE error history



(d) Testing MAE error history

Multiscale flow problem

Flow equation:

$$\begin{aligned}\kappa^{-1}u + \nabla p &= 0 && \text{in } D \\ \operatorname{div}(u) &= f && \text{in } D \\ u \cdot n &= 0 && \text{on } \partial D\end{aligned}$$

κ : heterogeneous permeability field.

On fine grid \mathcal{T}_h , using Mixed finite element method:

$$\begin{bmatrix} A_h(\kappa) & B_h^T \\ B_h & 0 \end{bmatrix} \begin{bmatrix} u_h \\ p_h \end{bmatrix} = \begin{bmatrix} 0 \\ -F \end{bmatrix} \quad (9)$$

Velocity approximation: RT_0 (the lowest order Raviart-Thomas element) .

Pressure approximation: P_0 (piecewise constant element).

Difficulties: $A_h(\kappa)$ is large, and depends on κ .

Network architecture

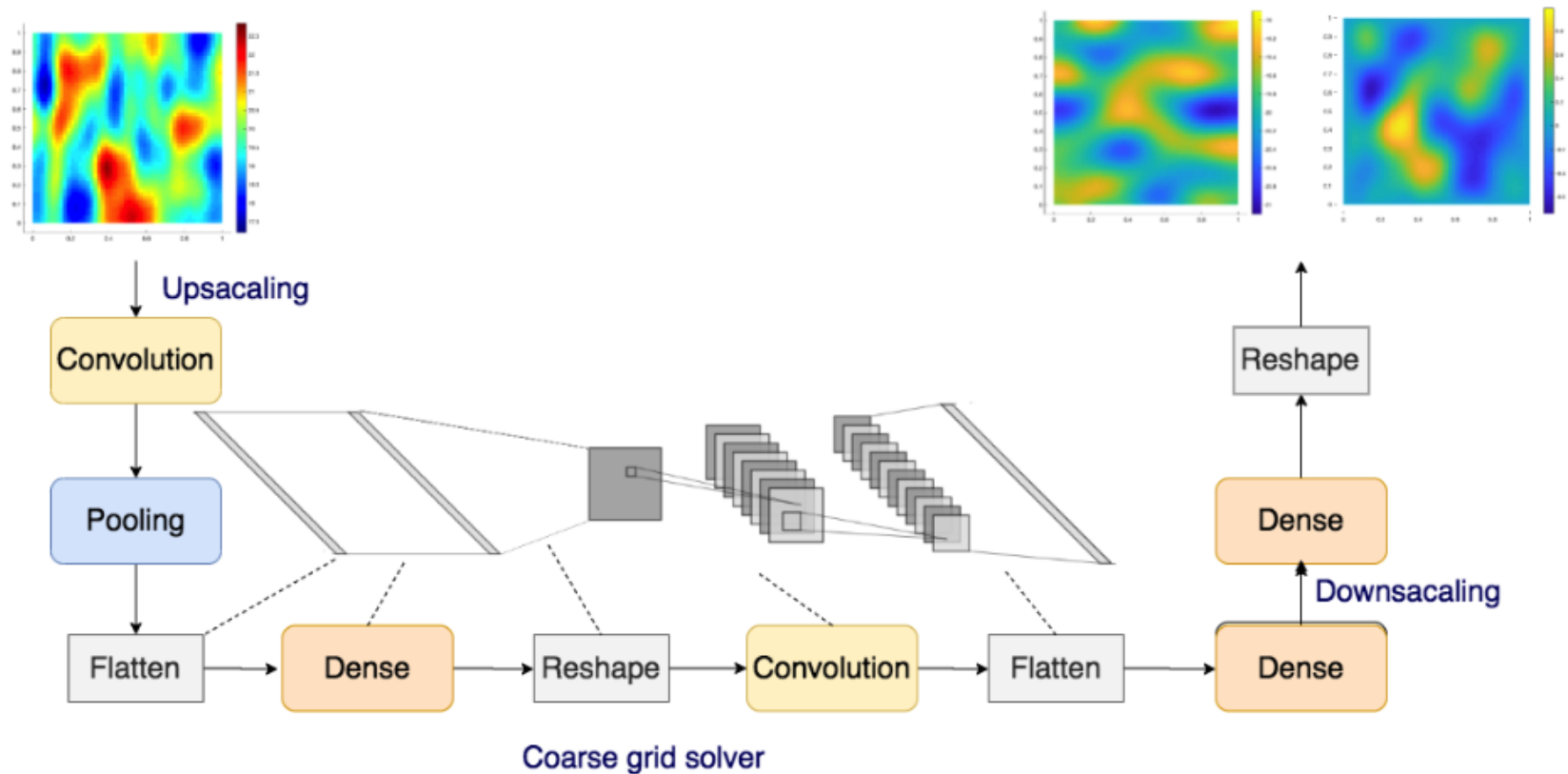
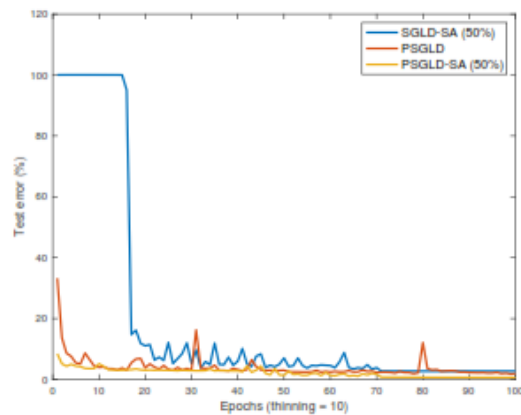


Figure 2: An illustration of the network architecture for flow approximation.

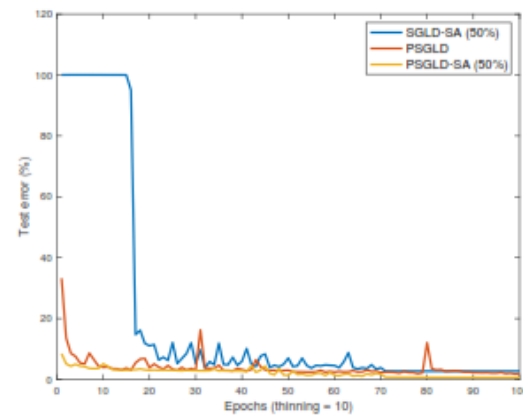
Numerical results

Dense (6,110,624 weight parameters)		
	PSGLD (e_1/ e_2 %)	SGLD(e_1/ e_2 %)
KLE32	0.75/0.57	2.37 /2.17
KLE64	0.82/0.63	2.38 /2.25
KLE128	2.13 /1.93	2.90 /2.60
	PSGLD-SA (e_1/ e_2 %)	SGLD-SA (e_1/ e_2 %)
Sparse rate 50% (2,326,049 weight parameters)		
KLE32	0.59/0.56	2.67 /2.35
KLE64	0.78 /0.58	2.68 /2.41
KLE128	1.60 /1.31	3.47 /3.00
Sparse rate 70% (758,738 weight parameters)		
KLE32	0.58/ 0.51	2.28 /2.10
KLE64	0.76 /0.61	2.40 /2.97
KLE128	1.79 /1.60	3.51/3.02

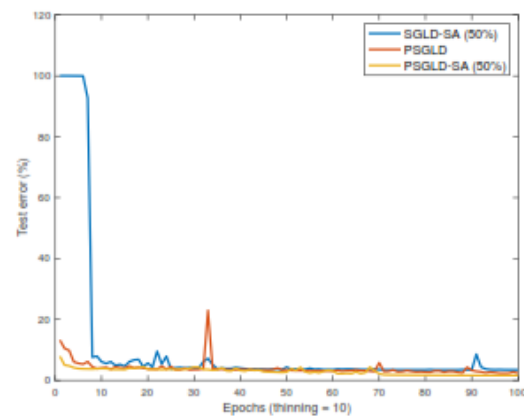
Table 2: Errors between the true velocity and predicted velocity (from trained neural networks) using SGLD, PSGLD, SGLD-SA, and proposed PSGLD-SA. Mean errors of 300 testing cases.



(a) KLE 32, learning curves



(b) KLE 64, learning curves



(c) KLE 128, learning curves

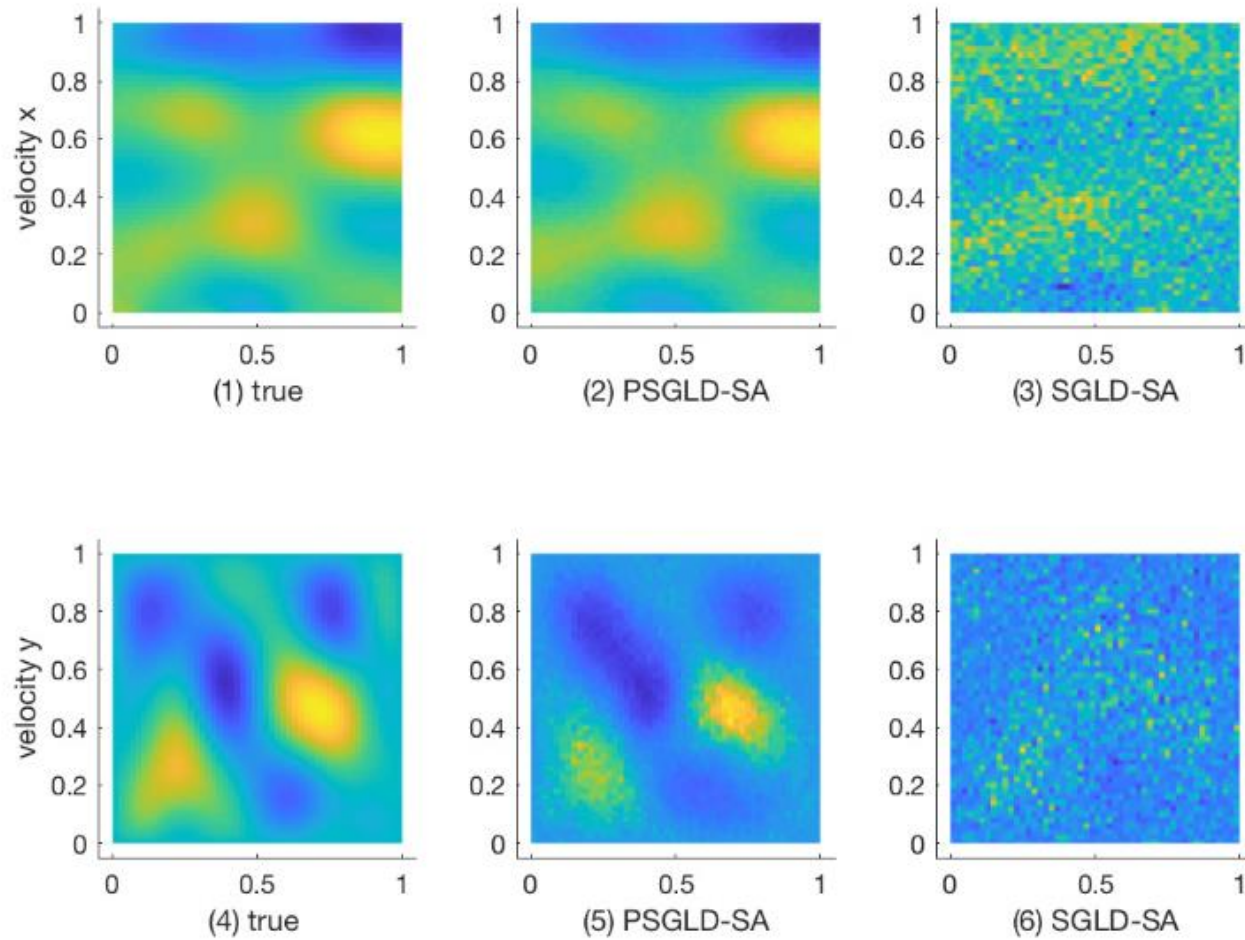


Figure 3: KLE 32. True and prediction solutions.

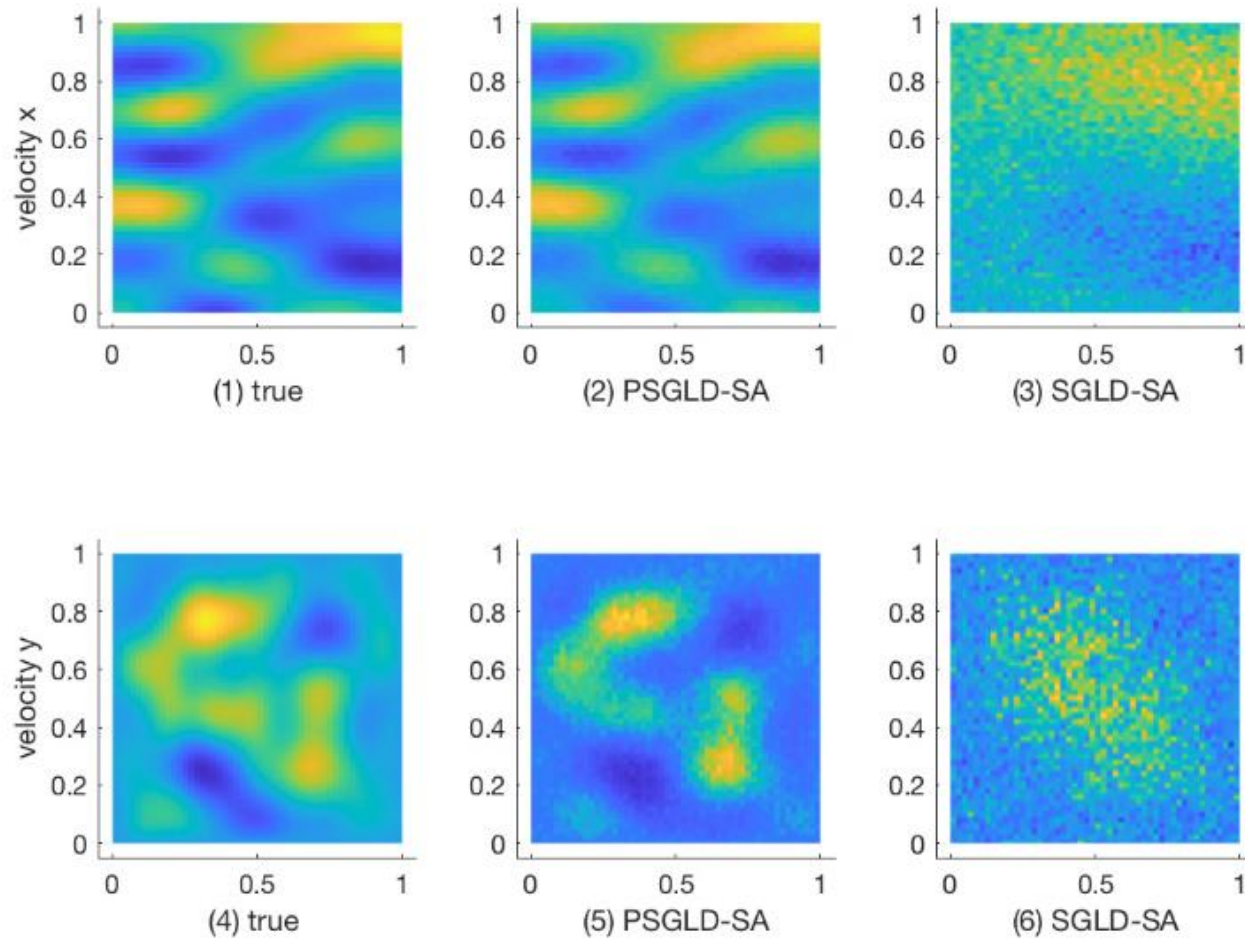


Figure 4: KLE 64. True and prediction solutions.

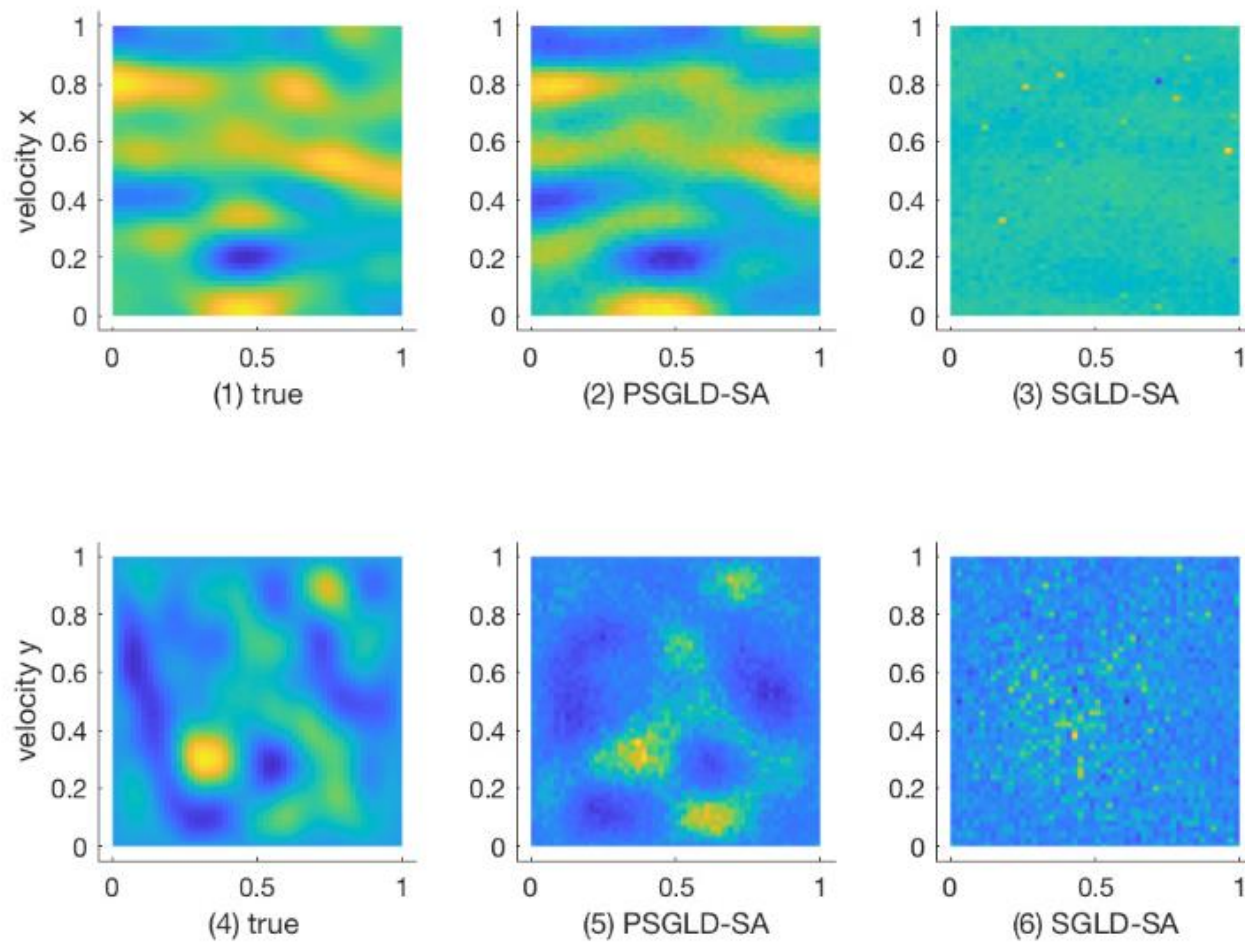
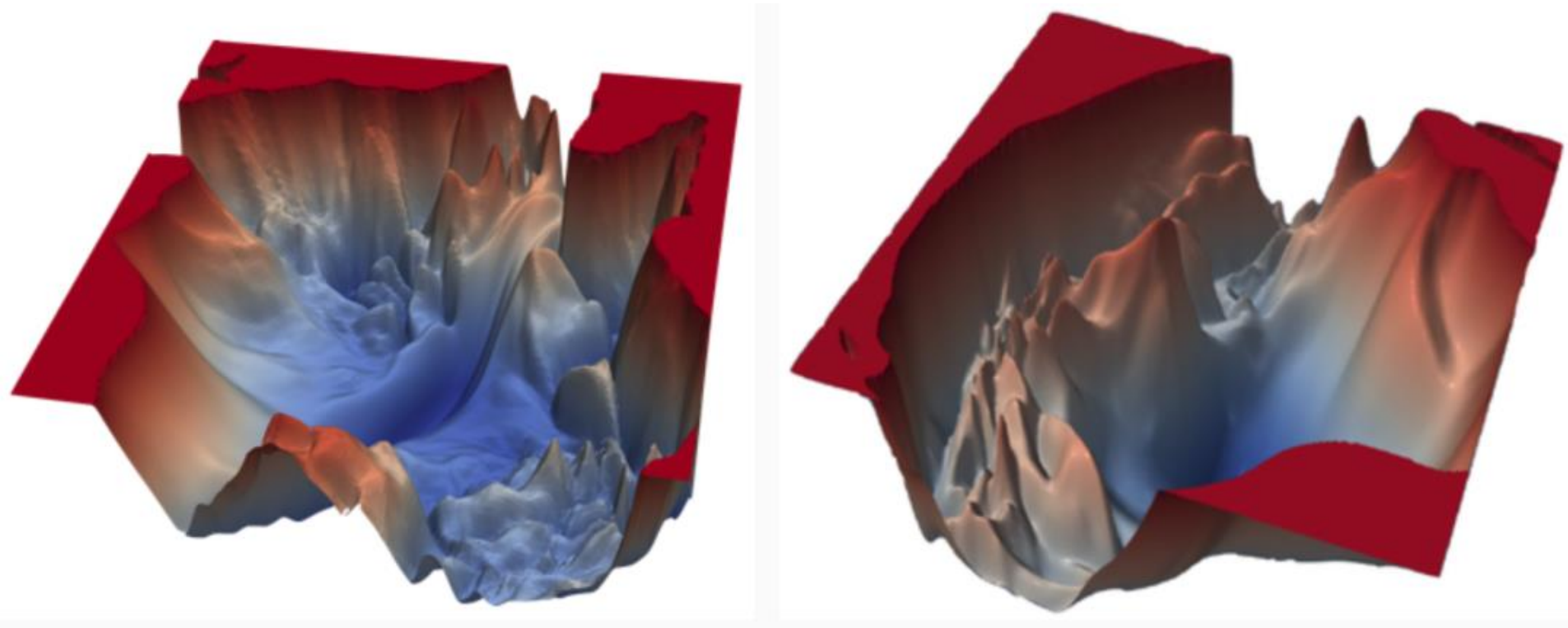


Figure 5: KLE 128. True and prediction solutions.

Summary

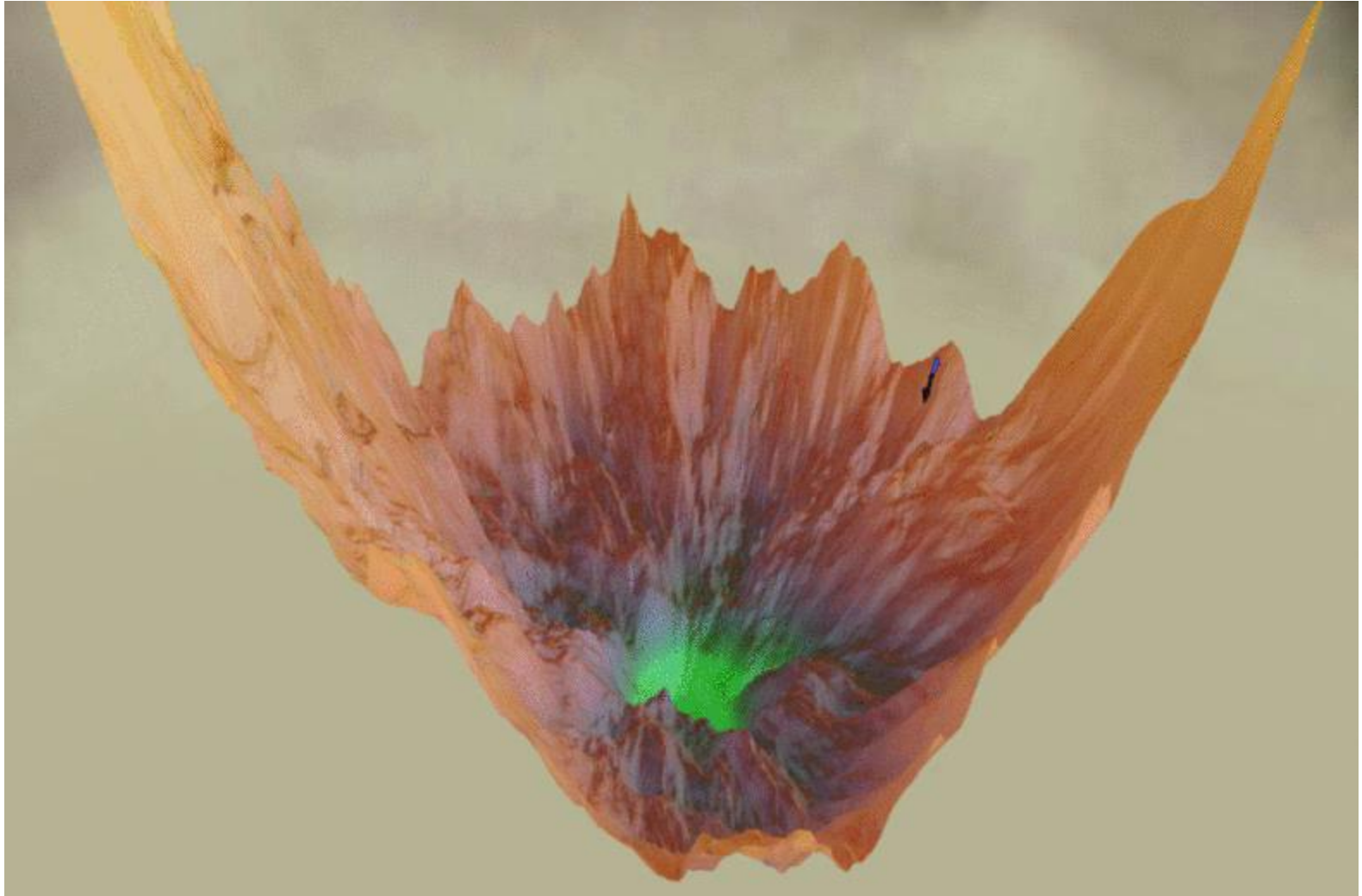
- Propose a class of adaptive stochastic gradient MCMC framework with wide applications and proved the convergence under mild assumptions.
- It can be applied to the empirical Bayesian method to learn a class of hierarchical Bayes models, yielding data-driven adaptive penalties.
- Propose an adaptive Bayesian sparse deep learning algorithm for regression problems
- Optimize SSSL priors through stochastic approximation. Sparsity is ensured, resulting in less usage in memory and computational power
- Sample with preconditioner SGLD, which is adaptive to local geometric and results in faster convergence
- Converge to the asymptotically correct distribution with a controllable bias introduced by SA
- It achieves the state-of-the-art compression performance on Resnet20, which outperforms the existing methods by a large margin.

Visualization the Loss Landscape of Deep Neural Nets



The loss landscape of modern deep neural nets [Li et al., 2018]

Gradient Descent Fails



Credit to losslandscape.com

Outline:

- ❖ Incorporate Physics Knowledge and AI to design new interpretable models
- ❖ Incorporate Physics Knowledge into AI to predict multiscale problems: NH-PINN
- ❖ Interpretable AI enables data-driven scientific discovery with uncertainty quantification capability – ALZHEIMER's Disease Prediction
- ❖ Sparse Neural Architecture Design with quantified uncertainties
- ❖ **Scalable training large-scale Deep Neural Network**

Scalable training large-scale Deep Neural Network:

Question: How can we design efficient optimization/sampling algorithms to train large-scale deep neural networks?

Goal: Enable Fast training large-scale DNN.

W. Deng, X. Zhang, F. Liang, **G. Lin**, An adaptive empirical Bayesian method for sparse deep learning, **2019 Conference on Neural Information Processing Systems (NIPS)**, Dec. 8 – Dec. 14, 2019, Vancouver, Canada.

NeurIPS'19, NeurIPS'20, ICML'20, ICLR'21, JCP'20, JCP'21a, JCP'21b

Scalable algorithms for Bayesian deep learning via Stochastic Gradient Monte Carlo and Beyond

Guang Lin ¹

Joint work with **W. Deng, Y. Wang, Q. Feng, L. Gao, G. Karagiannis, F. Liang**

August 13, 2021

¹Departments of Mathematics & School of Mechanical Engineering, Purdue University

NeurIPS'19, NeurIPS'20, ICML'20, ICLR'21, JCP'20, JCP'21a, JCP'21b

Markov chain Monte Carlo

Uncertainty quantification is crucial for AI safety problems and reinforcement learning, which draws our attention to **Markov chain Monte Carlo (MCMC)**, which is known for

- Multi-modal sampling → Accurate predictive **confidence interval**
- Non-convex optimization → Better **point estimate**

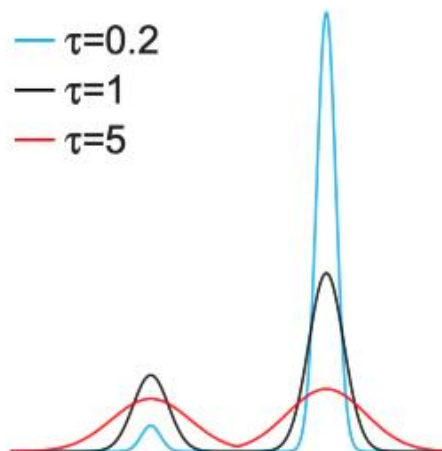
Langevin diffusion

A famous sampling algorithm is called Langevin diffusion.

$$d\beta_t = -\nabla U(\beta_t)dt + \sqrt{2\tau}d\mathbf{W}_t,$$

where β_t is the parameter at time t , $U(\cdot)$ is the energy function, \mathbf{W}_t is a Brownian motion and τ is the temperature.

As $t \rightarrow \infty$, β_t converges to the stationary Gibbs distribution $Ce^{-\frac{U(\beta)}{\tau}}$.



(a) Gibbs measures at three different temperatures τ .

Stochastic gradient Langevin dynamics

However, evaluating gradient in big data problems is **too costly**.

To tackle this issue, Max Welling, etc [Welling and Teh, 2011] proposed the stochastic gradient Langevin dynamics algorithm (SGLD)

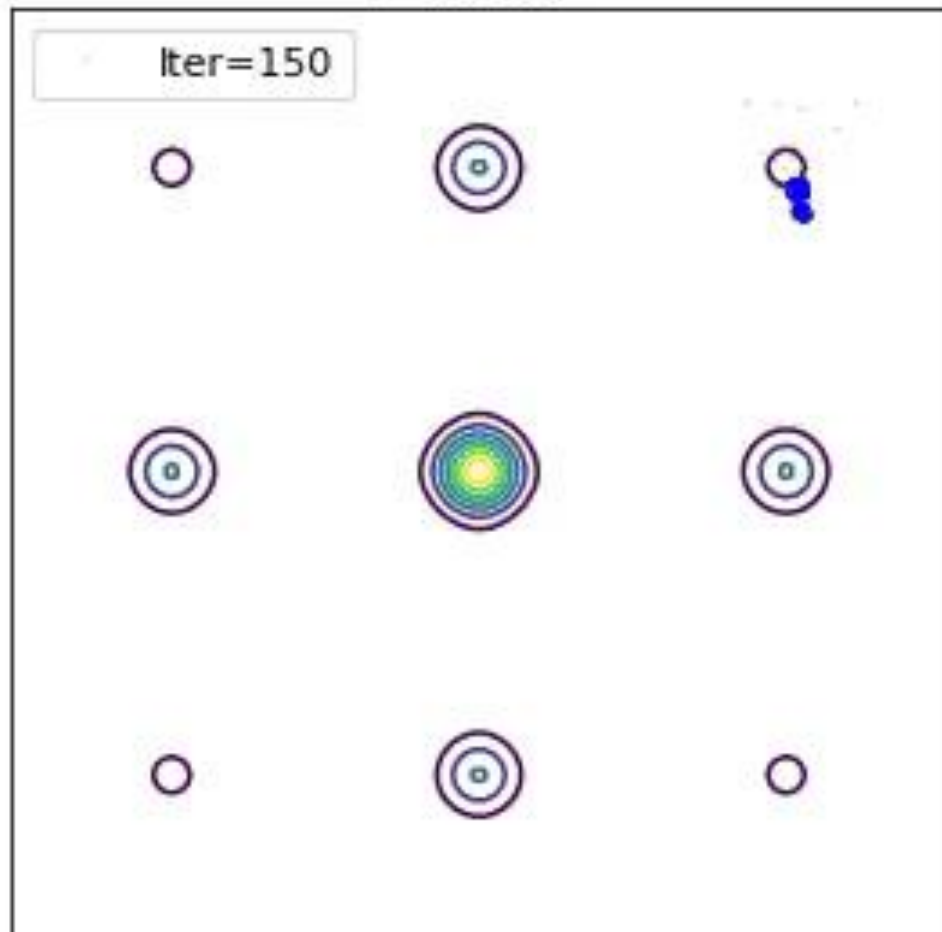
$$\beta_{k+1} = \beta_k - \eta \nabla \tilde{U}(\beta_k) + \mathcal{N}(0, 2\eta\tau \mathbf{I}). \quad (1)$$

As $t \rightarrow \infty$ and $\eta \rightarrow 0$, β_t converges weakly to the stationary Gibbs distribution $Ce^{-\frac{U(\beta)}{\tau}}$.

Stochastic gradient Langevin dynamics

Sample from a multi-modal distribution

SGLD



Acceleration strategies for MCMC

Most popular strategies to *accelerate* MCMC:

- Simulated annealing [Kirkpatrick et al., 1983]
- **Replica exchange MCMC [Swendsen and Wang, 1986]**

Replica Exchange SGLD

Wei Deng, et al., ICML 2020

Replica exchange Langevin diffusion

Consider two Langevin diffusion processes with $\tau_1 > \tau_2$

$$\begin{aligned}d\beta_t^{(1)} &= -\nabla U(\beta_t^{(1)})dt + \sqrt{2\tau_1}d\mathbf{W}_t^{(1)} \\d\beta_t^{(2)} &= -\nabla U(\beta_t^{(2)})dt + \sqrt{2\tau_2}d\mathbf{W}_t^{(2)},\end{aligned}$$

Moreover, the positions of the two particles swap with a probability

$$S(\beta_t^{(1)}, \beta_t^{(2)}) := e^{\left(\frac{1}{\tau_1} - \frac{1}{\tau_2}\right)(U(\beta_t^{(1)}) - U(\beta_t^{(2)}))}$$

In other words, a jump process is included in a Markov process

$$\begin{aligned}\mathbb{P}(\beta_{t+dt} = (\beta_t^{(2)}, \beta_t^{(1)}) | \beta_t = (\beta_t^{(1)}, \beta_t^{(2)})) &= rS(\beta_t^{(1)}, \beta_t^{(2)})dt \\ \mathbb{P}(\beta_{t+dt} = (\beta_t^{(1)}, \beta_t^{(2)}) | \beta_t = (\beta_t^{(1)}, \beta_t^{(2)})) &= 1 - rS(\beta_t^{(1)}, \beta_t^{(2)})dt\end{aligned}$$

A demo

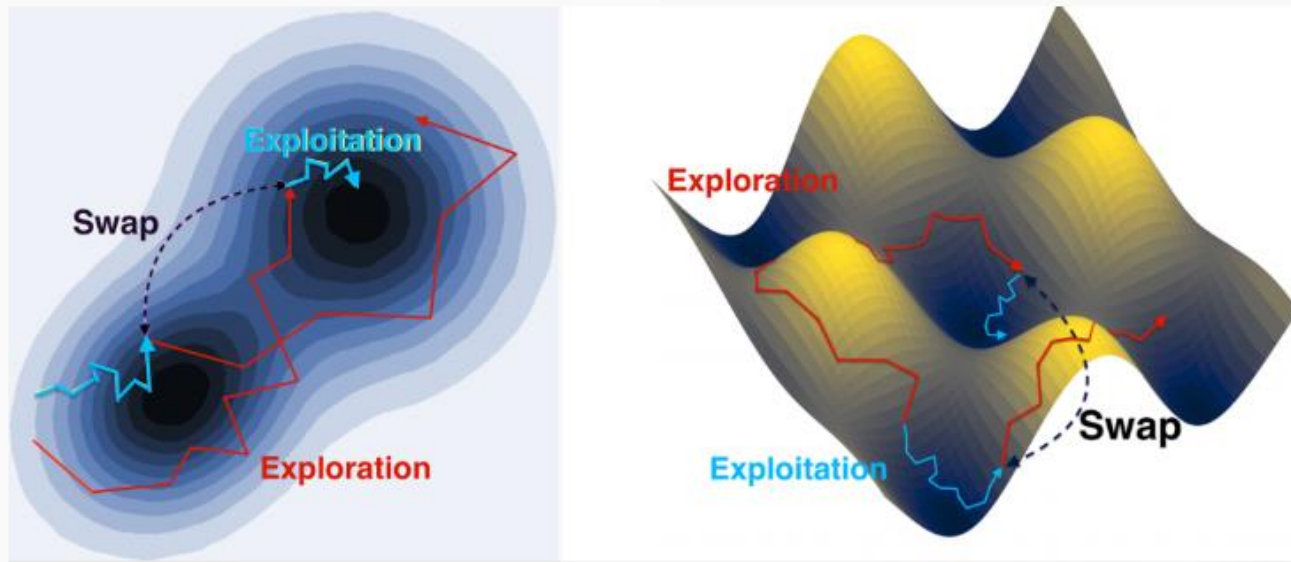


Figure 1: Trajectory plot for replica exchange Langevin diffusion.

Why the naïve numerical algorithm fails

Consider the scalable stochastic gradient Langevin dynamics algorithm [Welling and Teh, 2011]

$$\begin{aligned}\tilde{\beta}_{k+1}^{(1)} &= \tilde{\beta}_k^{(1)} - \eta_k \nabla \tilde{L}(\tilde{\beta}_k^{(1)}) + \sqrt{2\eta_k \tau_1} \xi_k^{(1)} \\ \tilde{\beta}_{k+1}^{(2)} &= \tilde{\beta}_k^{(2)} - \eta_k \nabla \tilde{L}(\tilde{\beta}_k^{(2)}) + \sqrt{2\eta_k \tau_2} \xi_k^{(2)}.\end{aligned}$$

Swap the chains with a **naïve** swapping rate $r\mathbb{S}(\tilde{\beta}_{k+1}^{(1)}, \tilde{\beta}_{k+1}^{(2)})\eta_k$ [§]:

$$\mathbb{S}(\tilde{\beta}_{k+1}^{(1)}, \tilde{\beta}_{k+1}^{(2)}) = e^{\left(\frac{1}{\tau_1} - \frac{1}{\tau_2}\right) \left(\tilde{L}(\tilde{\beta}_{k+1}^{(1)}) - \tilde{L}(\tilde{\beta}_{k+1}^{(2)})\right)}. \quad (2)$$

Exponentiating the unbiased estimators $\tilde{L}(\tilde{\beta}_{k+1}^{(\cdot)})$ leads to a **large bias**.

[§]In the implementations, we fix $r\eta_k = 1$ by default.

A corrected algorithm

Assume $\tilde{L}(\theta) \sim \mathcal{N}(L(\theta), \sigma^2)$ and consider the **geometric Brownian motion** of $\{\tilde{S}_t\}_{t \in [0,1]}$ in each swap as a Martingale

$$\begin{aligned}\tilde{S}_t &= e^{\left(\frac{1}{\tau_1} - \frac{1}{\tau_2}\right) \left(\tilde{L}(\tilde{\beta}^{(1)}) - \tilde{L}(\tilde{\beta}^{(2)}) - \left(\frac{1}{\tau_1} - \frac{1}{\tau_2}\right) \sigma^2 t\right)} \\ &= e^{\left(\frac{1}{\tau_1} - \frac{1}{\tau_2}\right) \left(L(\tilde{\beta}^{(1)}) - L(\tilde{\beta}^{(2)}) - \left(\frac{1}{\tau_1} - \frac{1}{\tau_2}\right) \sigma^2 t + \sqrt{2} \sigma W_t\right)}.\end{aligned}\tag{3}$$

Taking the derivative of \tilde{S}_t with respect to t and W_t , Itô's lemma gives,

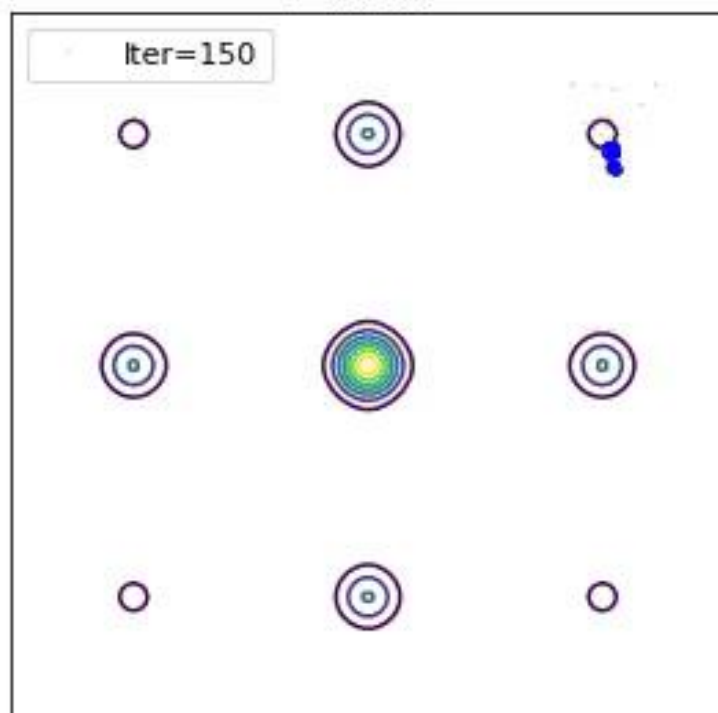
$$d\tilde{S}_t = \left(\frac{d\tilde{S}_t}{dt} + \frac{1}{2} \frac{d^2\tilde{S}_t}{dW_t^2} \right) dt + \frac{d\tilde{S}_t}{dW_t} dW_t = \sqrt{2} \left(\frac{1}{\tau_1} - \frac{1}{\tau_2} \right) \sigma \tilde{S}_t dW_t.$$

By fixing $t = 1$ in (3), we have the **suggested unbiased swapping rate**

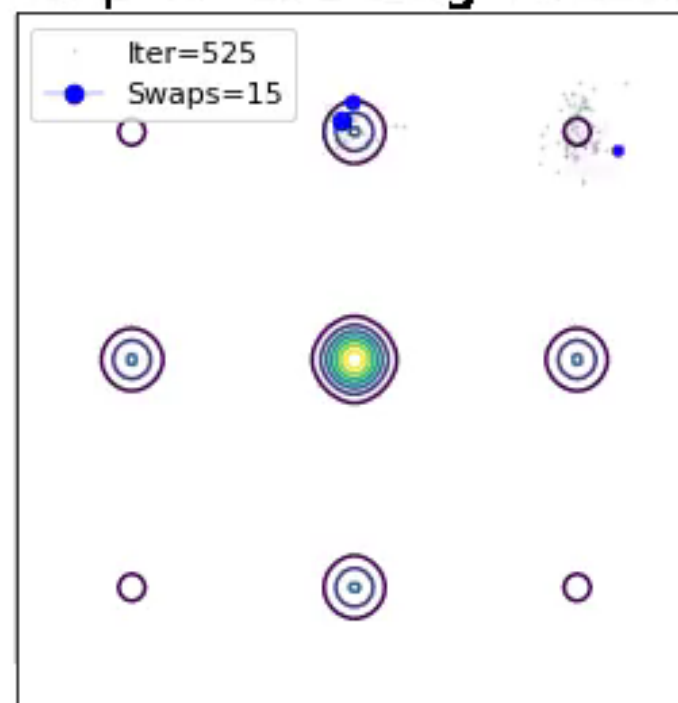
$$\tilde{S}_1 = e^{\left(\frac{1}{\tau_1} - \frac{1}{\tau_2}\right) \left(\tilde{L}(\tilde{\beta}^{(1)}) - \tilde{L}(\tilde{\beta}^{(2)}) - \left(\frac{1}{\tau_1} - \frac{1}{\tau_2}\right) \sigma^2\right)}.$$

Replica exchange Stochastic gradient Langevin dynamics

SGLD



Replica exchange SGLD



Acceleration via replica exchange

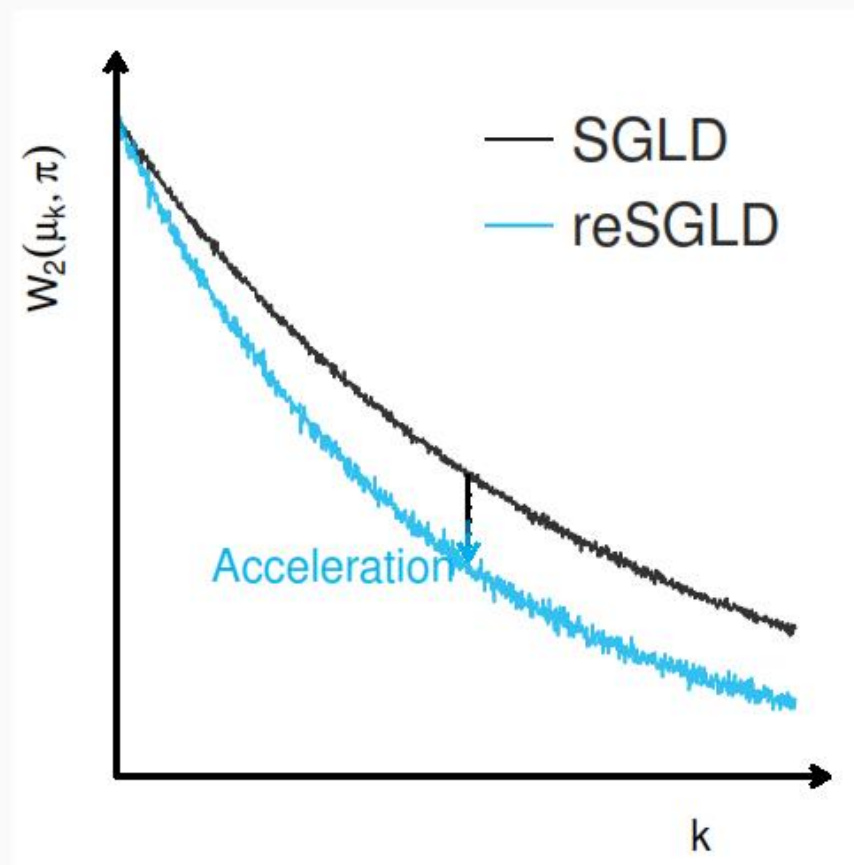


Figure 2: Acceleration via replica exchange (swaps/ interactions)

Accelerating convergence via variance reduction

Can we do better?

Exponential acceleration via variance reduction

Wei Deng et al., ICLR 2021

Accelerating convergence via variance reduction

The desire to obtain more effective swaps drives us to design more efficient energy estimators.

To reduce the variance of the noisy energy estimator

$L(B|\beta^{(h)}) = \frac{N}{n} \sum_{i \in B} L(\mathbf{x}_i|\beta^{(h)})$ for $h \in \{1, 2\}$, we consider an unbiased estimator $L(B|\hat{\beta}^{(h)})$ for $\sum_{i=1}^N L(\mathbf{x}_i|\hat{\beta}^{(h)})$ and a constant c , we see that a new estimator $\tilde{L}(B|\beta^{(h)})$, which follows

$$\tilde{L}(B|\beta^{(h)}) = L(B|\beta^{(h)}) + c \left(L(B|\hat{\beta}^{(h)}) - \sum_{i=1}^N L(\mathbf{x}_i|\hat{\beta}^{(h)}) \right), \quad (4)$$

is still the unbiased estimator for $\sum_{i=1}^N L(\mathbf{x}_i|\beta^{(h)})$.

Accelerating convergence via variance reduction

By decomposing the variance, we have

$$\text{Var}(\tilde{L}(B|\beta^{(h)})) = \text{Var}\left(L(B|\beta^{(h)})\right) + c^2 \text{Var}\left(L(B|\hat{\beta}^{(h)})\right) + 2c \text{Cov}\left(L(B|\beta^{(h)}), L(B|\hat{\beta}^{(h)})\right).$$

In such a case, $\text{Var}(\tilde{L}(B|\beta^{(h)}))$ achieves the minimum variance

$(1 - \rho^2) \text{Var}(L(B|\beta^{(h)}))$ given $c^* := \frac{-\text{Cov}(L(B|\beta^{(h)}), L(B|\hat{\beta}^{(h)}))}{\text{Var}(L(B|\beta^{(h)}))}$, where $\text{Cov}(\cdot, \cdot)$ denotes the covariance and ρ is the correlation coefficient.

Accelerating convergence via variance reduction

To make variance reduction work, it requires two crucial components.

- To propose a correlated control variate $\hat{\beta}$
 - Update $\hat{\beta}^{(h)} = \beta_{m \lfloor \frac{k}{m} \rfloor}^{(h)}$ every m iterations
- The optimal c is unknown.
 - Set $c = -1$ for highly correlated energy estimators.
 - Set adaptive c for the less correlated.

Reduction of Variance

VR-reSGLD may lead to a more efficient energy estimator with a much smaller variance.

Lemma (Variance-reduced energy estimator)

Under the smoothness and dissipativity assumptions, the variance of the variance-reduced energy estimator $\tilde{L}(B|\beta^{(h)})$, where $h \in \{1, 2\}$, is upper bounded by

$$\text{Var} \left(\tilde{L}(B|\beta^{(h)}) \right) \leq \min \left\{ \mathcal{O} \left(\frac{m^2 \eta}{n} \right), \text{Var} \left(\frac{N}{n} \sum_{i \in B} L(\mathbf{x}_i | \beta^{(1)}) \right) \right\},$$

where the detailed $\mathcal{O}(\cdot)$ constants is shown in the appendix [Deng et al., 2021].

A smaller variance implies more effective swaps

The variance-reduced energy estimator $\tilde{L}(B|\beta^{(h)})$ doesn't directly affect $\mathbb{E}[\tilde{S}_{\eta,m,n}]$ within the support $[0, \infty]$. However, the unbounded support is not appropriate for numerical algorithms, and only the truncated swapping rate $S_{\eta,m,n} = \min\{1, \tilde{S}_{\eta,m,n}\}$ is considered. As such, the truncated swapping rate becomes significantly smaller.

Lemma (Variance reduction for larger swapping rates)

Given a large enough batch size n , the variance-reduced energy estimator $\tilde{L}(B_k|\beta_k^{(h)})$ yields a truncated swapping rate that satisfies

$$\mathbb{E}[S_{\eta,m,n}] \approx \min \left\{ 1, S(\beta^{(1)}, \beta^{(2)}) \left(\mathcal{O}\left(\frac{1}{n^2}\right) + e^{-\mathcal{O}\left(\frac{m^2\eta}{n} + \frac{1}{n^2}\right)} \right) \right\}. \quad (5)$$

Acceleration via variance-reduced replica exchange

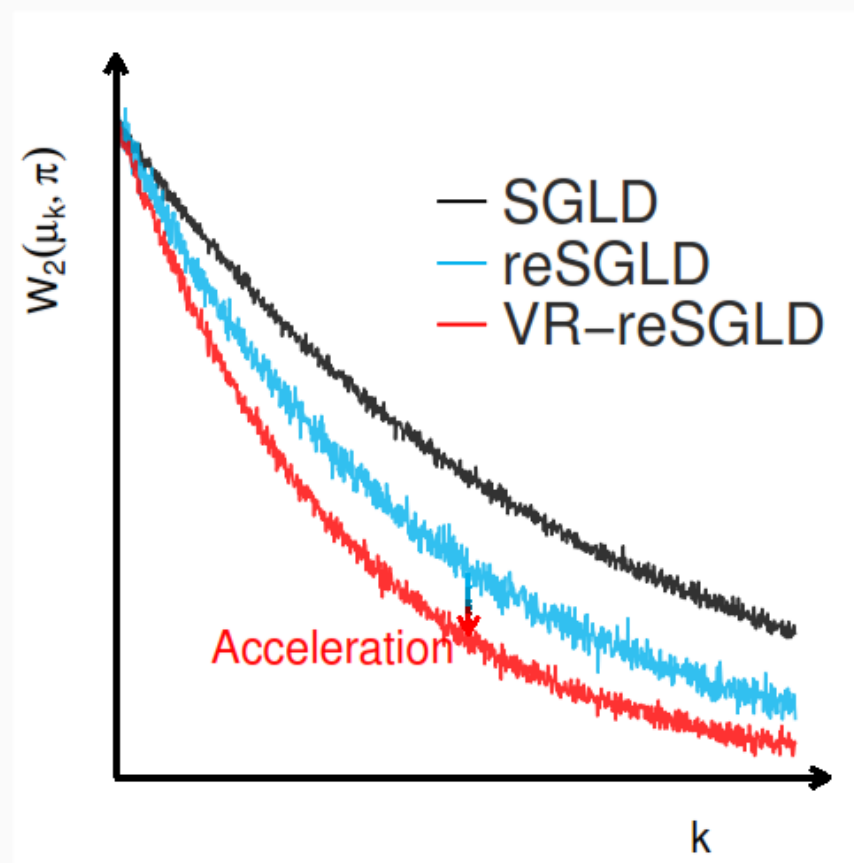


Figure 3: Acceleration via variance-reduced replica exchange.

1D simulation of Gaussian mixture

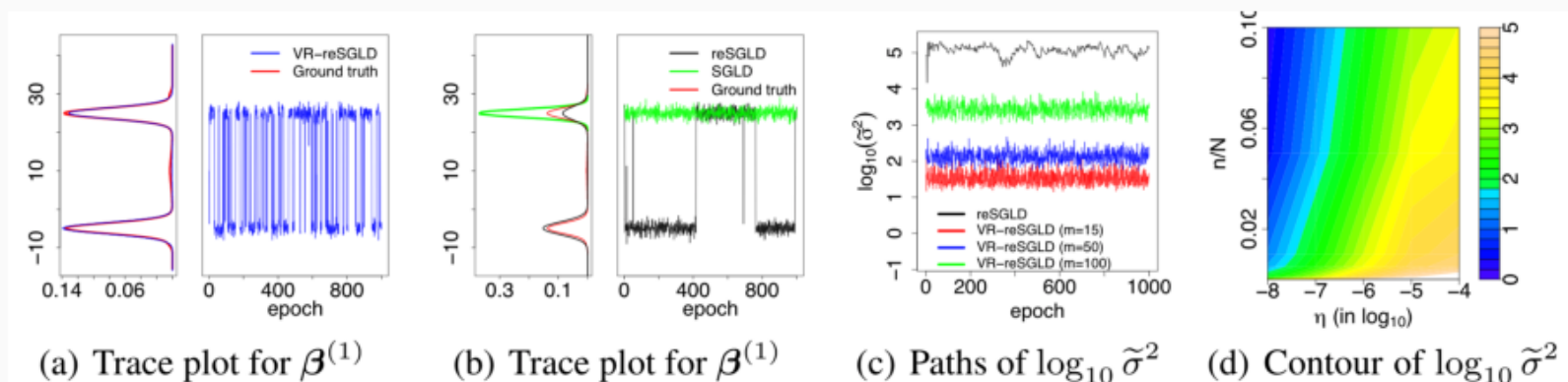


Figure 2: Trace plots, KDEs of $\beta^{(1)}$, and sensitivity study of $\tilde{\sigma}^2$ with respect to m , η and n .

Non-convex optimization on CIFAR10 and CIFAR100

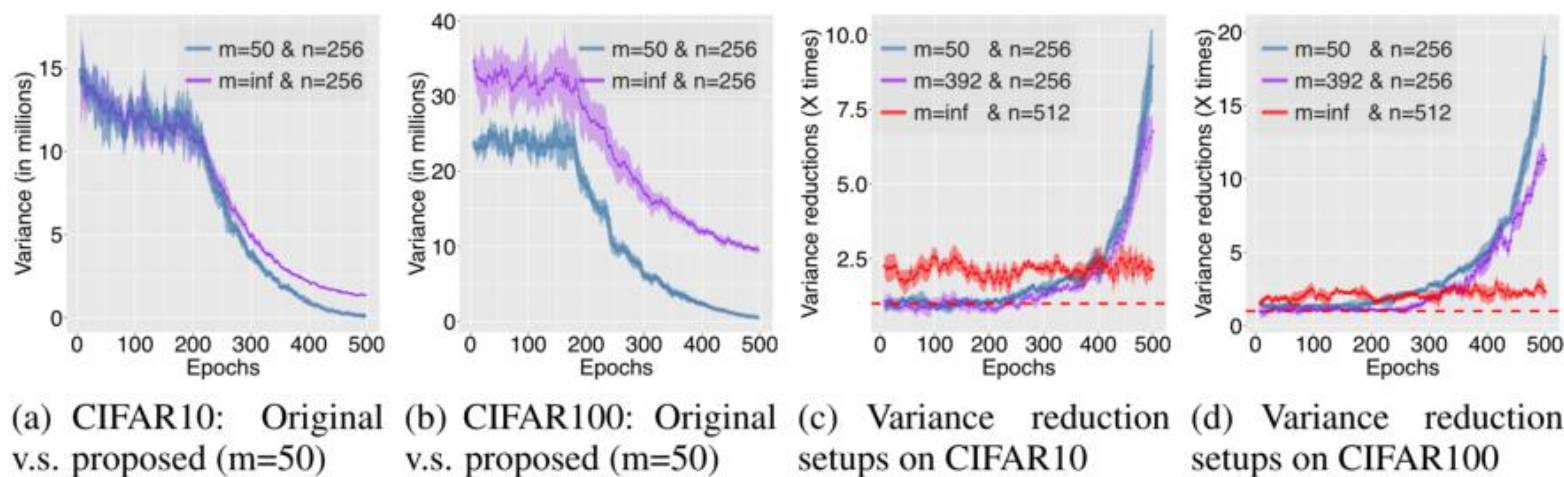


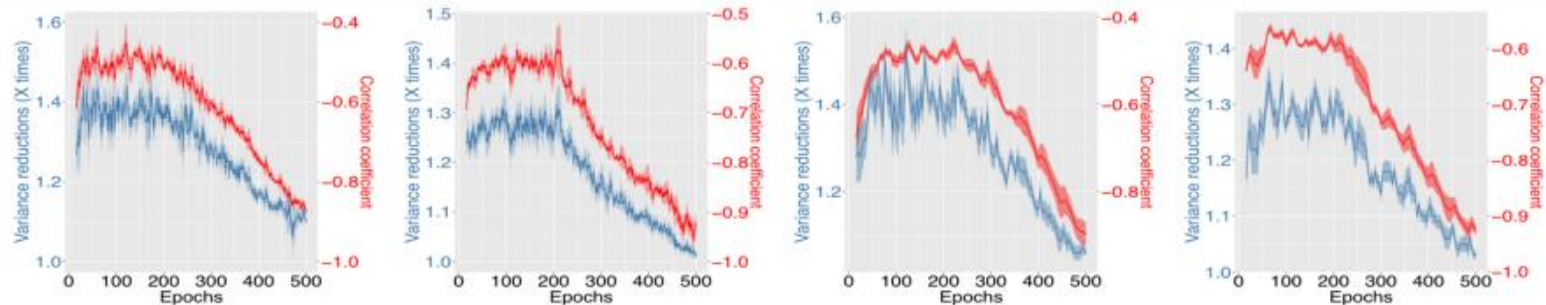
Figure 3: Variance reduction on the noisy energy estimators on CIFAR10 & CIFAR100 datasets.

Non-convex optimization on CIFAR10 and CIFAR100

TABLE 1: PREDICTION ACCURACIES (%) BASED ON BAYESIAN MODEL AVERAGING.

METHOD	CIFAR10			CIFAR100		
	RESNET20	RESNET32	RESNET56	RESNET20	RESNET32	RESNET56
M-SGD	94.07 \pm 0.11	95.11 \pm 0.07	96.05 \pm 0.21	71.93 \pm 0.13	74.65 \pm 0.20	78.76 \pm 0.24
SGHMC	94.16 \pm 0.13	95.17 \pm 0.08	96.04 \pm 0.18	72.09 \pm 0.14	74.80 \pm 0.19	78.95 \pm 0.22
reSGHMC	94.56 \pm 0.23	95.44 \pm 0.16	96.15 \pm 0.17	73.94 \pm 0.34	76.38 \pm 0.23	79.86 \pm 0.26
VR-reSGHMC	94.84\pm0.11	95.62\pm0.09	96.32\pm0.15	74.83\pm0.18	77.40\pm0.27	80.62\pm0.22
cycSGHMC	94.61 \pm 0.15	95.56 \pm 0.12	96.19 \pm 0.17	74.21 \pm 0.22	76.60 \pm 0.25	80.39 \pm 0.21
cVR-reSGHMC	94.91\pm0.10	95.64\pm0.13	96.36\pm0.16	75.02\pm0.19	77.58\pm0.21	80.50\pm0.25

Non-convex optimization on CIFAR10 and CIFAR100



(a) CIFAR10 & $m=50$ (b) CIFAR100 & $m=50$ (c) CIFAR10 & $m=392$ (d) CIFAR100 & $m=392$

Figure 5: A study of variance reduction techniques using adaptive coefficient and non-adaptive coefficient on CIFAR10 & CIFAR100 datasets.

Summary

- Replica exchange stochastic gradient MCMC shows a potential in exponentially accelerating the convergence in non-convex learning. [Deng et al., 2020]
- Variance reduction of energy estimators yields exponential more effective swaps, which further accelerates the exponential convergence in non-convex learning. [Deng et al., 2021]
- This is the first work to do variance reduction on energy estimators in deep learning, which paves the road for accelerating advanced stochastic gradient MCMC algorithms in non-convex learning.

References i



Deng, W., Feng, Q., Gao, L., Liang, F., and Lin, G. (2020).
Non-Convex Learning via Replica Exchange Stochastic Gradient MCMC.

In Proc. of the International Conference on Machine Learning (ICML).



Deng, W., Feng, Q., Karagiannis, G., Lin, G., and Liang, F. (2021).
Accelerating Convergence of Replica Exchange Stochastic Gradient MCMC via Variance Reduction.

In Proc. of the International Conference on Learning Representation (ICLR).



Kirkpatrick, S., Jr, D. G., and Vecchi, M. P. (1983).
Optimization by Simulated Annealing.

Science, 220(4598):671–680.

Towards Third Wave AI: Interpretable, Robust Trustworthy Machine Learning for Diverse Applications in Science and Engineering



“...Because I had worked in the closest possible ways with physicists and engineers, I knew that our data can never be precise...”

Norbert Wiener