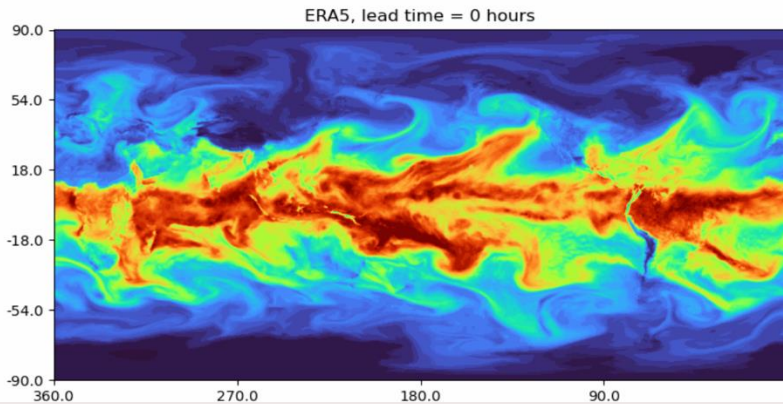
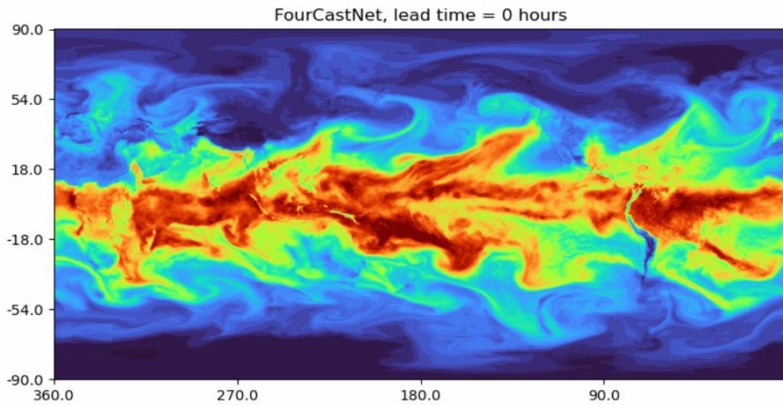


An introduction to SciML

VICTORITA DOLEAN

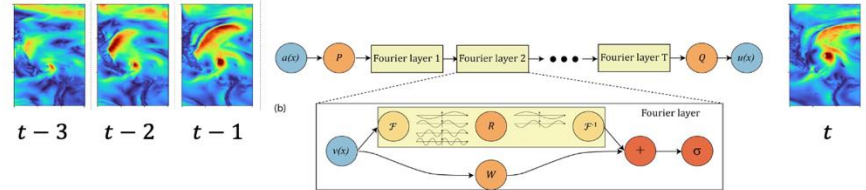


AI for science: a revolution?



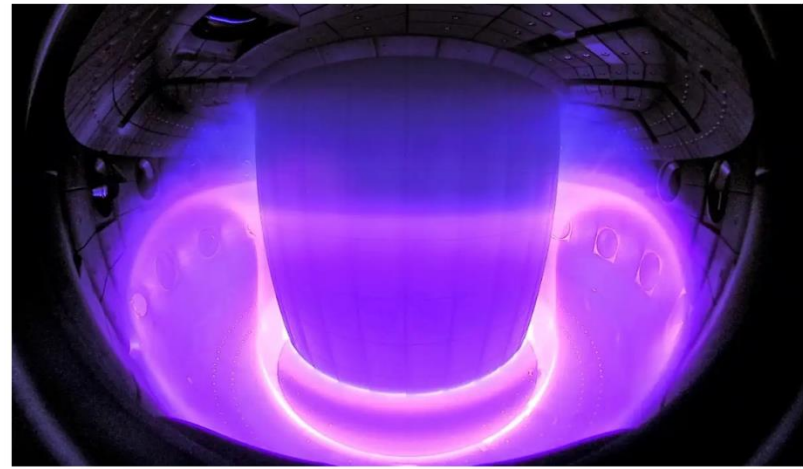
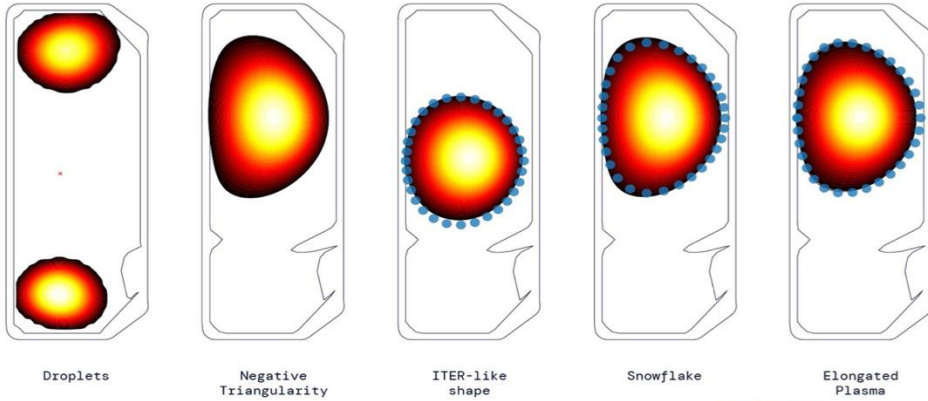
How Big Tech AI models nailed forecast for Hurricane Lee a week in advance

U.S. and European weather agencies are escalating their engagement with artificial intelligence as the technology rapidly advances

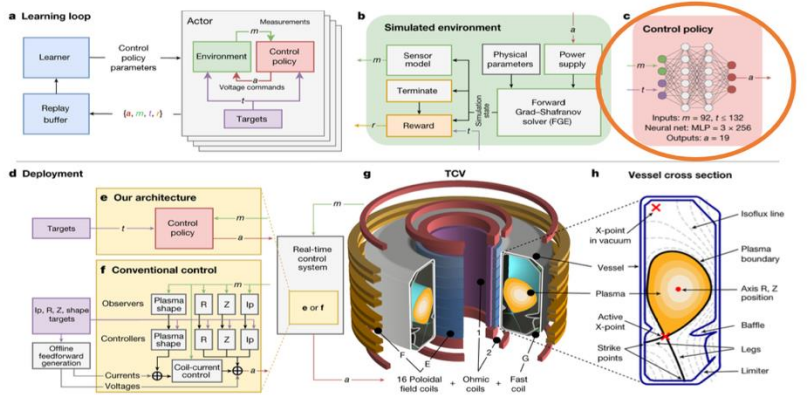


Pathak et al, FourCastNet: A Global Data-driven High-resolution Weather Model using Adaptive Fourier Neural Operators, ArXiv (2022)

AI for science: a revolution?



Variable Configuration Tokamak (TCV) in Lausanne, Switzerland
Source: DeepMind & SPC/EPFL



nature

Explore content ▾ About the journal ▾ Publish with us ▾

[nature](#) > [articles](#) > [article](#)

Article | [Open access](#) | Published: 16 February 2022

Magnetic control of tokamak plasmas through deep reinforcement learning

[Jonas Degraeve](#), [Federico Felici](#), [Jonas Buchli](#), [Michael Neunert](#), [Brendan Tracey](#), [Francesco Carpanese](#), [Timo Ewalds](#), [Roland Hafner](#), [Abbas Abdolmaleki](#), [Diego de las Casas](#), [Craig Donner](#), [Leslie Fritz](#), [Cristian Galperti](#), [Andrea Huber](#), [James Keeling](#), [Maria Tsimpoukelli](#), [Jackie Kay](#), [Antoine Merle](#), [Jean-Marc Moret](#), [Seb Noury](#), [Federico Pesamosca](#), [David Pfau](#), [Olivier Sauter](#), [Cristian Sommariva](#), ...
[Martin Riedmiller](#) + Show authors

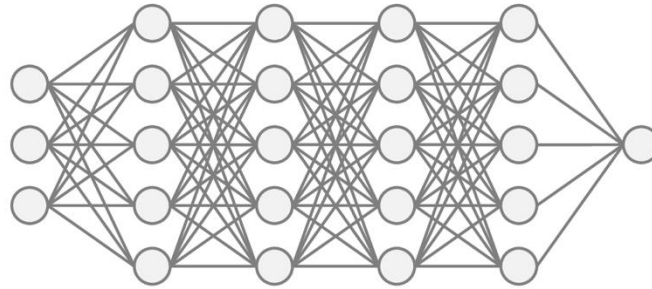
[Nature](#) **602**, 414–419 (2022) | [Cite this article](#)

206k Accesses | 223 Citations | 2430 Altmetric | [Metrics](#)

What is deep learning?



Input
 \mathbf{x}



Model
 $NN(\mathbf{x}, \boldsymbol{\theta})$

Probability(Dog) = 1

Output
 $\mathbf{y} = NN(\mathbf{x}, \boldsymbol{\theta})$

Deng et al,
ImageNet: A
large-scale
hierarchical
image
database,
CVPR (2009)



Neural networks are
simply **flexible functions**
fit to data

For example:

$$\mathbf{y} = W_2 \sigma(W_1 \mathbf{x} + \mathbf{b}_1) + \mathbf{b}_2$$

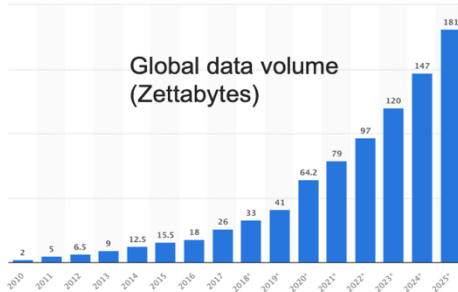
Trained ($\boldsymbol{\theta}$ learned) using:

- Many examples of inputs and outputs
- A loss function
- An optimisation algorithm, e.g. stochastic gradient descent

Why now?

Neural networks date back to the 1950's – so why is deep learning so popular today?

Rapidly increasing amounts of data



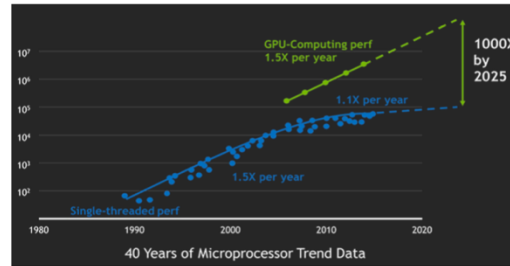
Source: Statista

IMAGENET



WIKIPEDIA
The Free Encyclopedia

Hardware improvements



Source: NVIDIA

- Graphical processing units (GPUs)
- Highly optimised for deep learning (massively parallel)

Software improvements

PyTorch



TensorFlow

Keras



- Mature deep learning frameworks
- Better training algorithms
- Deeper and more sophisticated architectures

Scientific machine learning (SciML)

Major problem

Naively using deep learning for scientific tasks usually leads to:

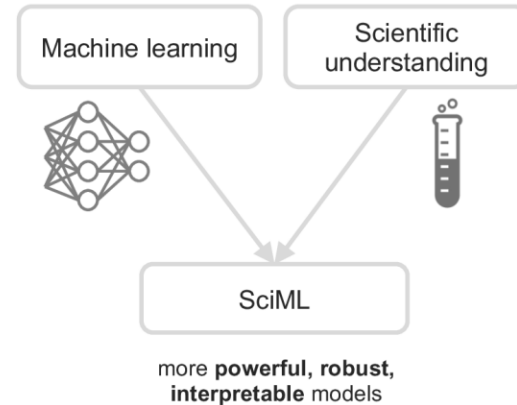
- Lack of interpretability
- Poor generalisation
- Lots of training data required

Do neural networks really “**understand**” the scientific tasks they are being applied to?

Traditional scientific method:

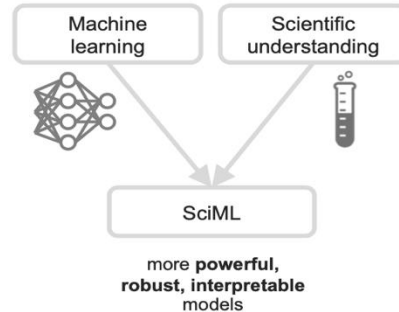
- Revolves around theory and experiment
- a good theory should be explainable and make **novel** predictions

Solution



Scientific machine learning

Hamiltonian neural networks
Learned sub-grid processes
Hidden physics models
Physics-informed neural networks
Solver-in-the-loop
Physics-constrained Gaussian processes
AI Feynman
DeepONets
PDE-Net
Algorithm unrolling
AlphaFold
Learned regularisation
Differentiable simulation
Physics-informed neural operators
Fourier neural operators
Encoding conservation laws
Encoding physical symmetries
Neural ODEs



A rapidly growing field

ICLR 2024 Workshop on
AI4DifferentialEquations in Science

Machine Learning and the Physical Sciences
Workshop at the 37th conference on Neural Information Processing Systems (NeurIPS)
December 15, 2023

Synergy of Scientific and Machine
Learning Modeling

ICML 2023 Workshop, July 28 2023, Room 320 of the Hawai'i Convention Center

AI for Science

NeurIPS 2021

ICML 2022

NeurIPS 2022

NeurIPS 2023



ICLR 2023 Workshop on Physics for Machine Learning

Physics4ML

The Symbiosis of Deep Learning and Differential Equations (DLDE)

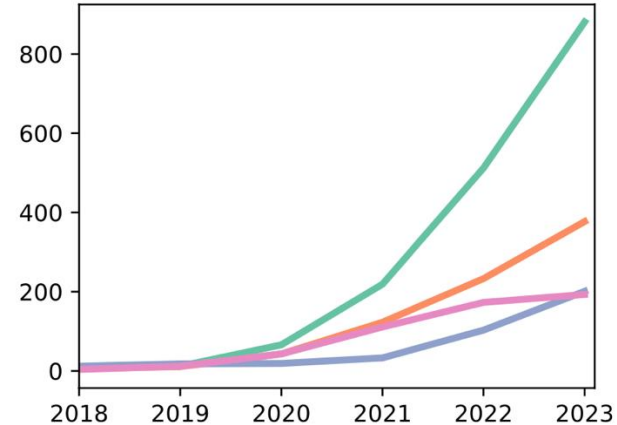
NeurIPS 2022 Workshop

AAAI-MLPS 2021

AAAI 2021 Spring Symposium on

Combining Artificial Intelligence and Machine Learning with Physics Sciences

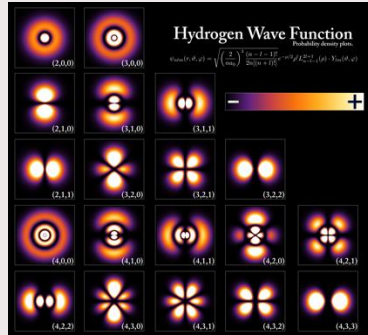
Number of publications each year



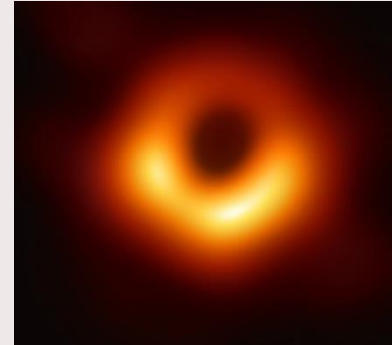
- physics-informed neural networks
- scientific machine learning / physics-informed ML / AI for science
- operator learning / neural operators
- differentiable physics / neural differential equations

Source: Scopus keyword search (Feb 2024)

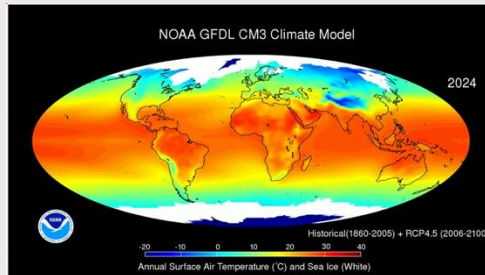
Importance of PDEs



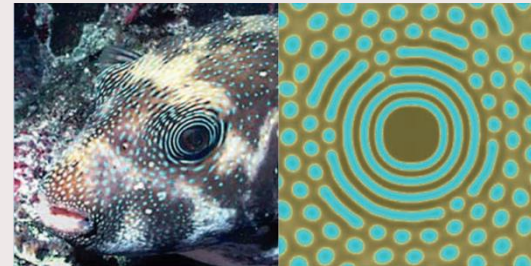
Source: Wikipedia
Schrödinger equation



Source: The Event Horizon Telescope (2019)
Einstein field equations



Source: NOAA
Navier-Stokes equations



Source: Kondo and Miura, Science (2010)
Reaction-diffusion equation

Key scientific tasks: simulation

$$b = F(a)$$

a = set of input conditions

F = physical model of the system (usually a PDE)

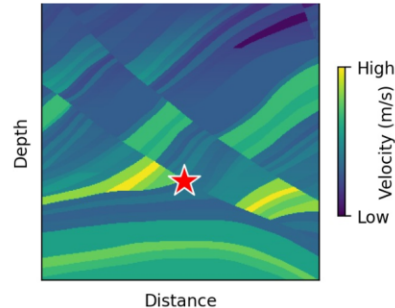
b = resulting properties given F and a

Key scientific tasks: simulation



Simulation is:

- Crucial for practically all domains of science
- Essential for understanding the behaviour of complex phenomena
- Usually used as a starting point for other tasks (e.g. inverse / control / design problems)

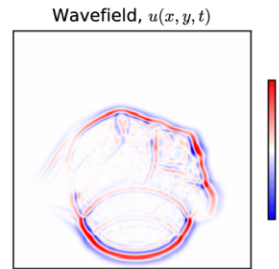
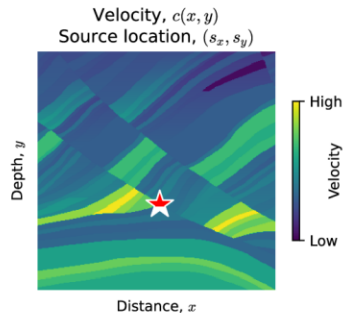


$$\nabla^2 u - \frac{1}{c(x)^2} \frac{\partial^2 u}{\partial t^2} = s$$

Wave equation

Key scientific tasks: simulation

$$b = F(a)$$



a = velocity model, $c(x, y)$,
source location, s

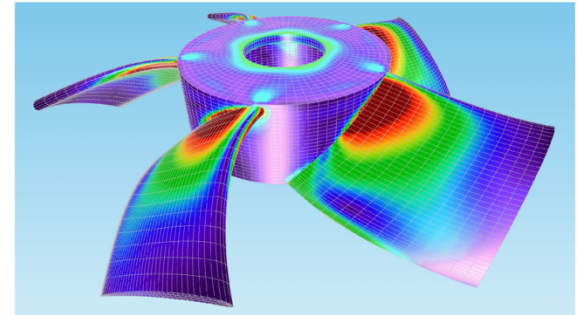
b = wavefield, $u(x, t)$

$$\begin{aligned}\nabla^2 u - \frac{1}{c(x)^2} \frac{\partial^2 u}{\partial t^2} &= \delta(x = s, t = 0) \\ u(x, t = 0) &= 0 \\ u'(x, t = 0) &= 0\end{aligned}$$

F = A method for solving the wave equation (e.g. finite difference simulation)

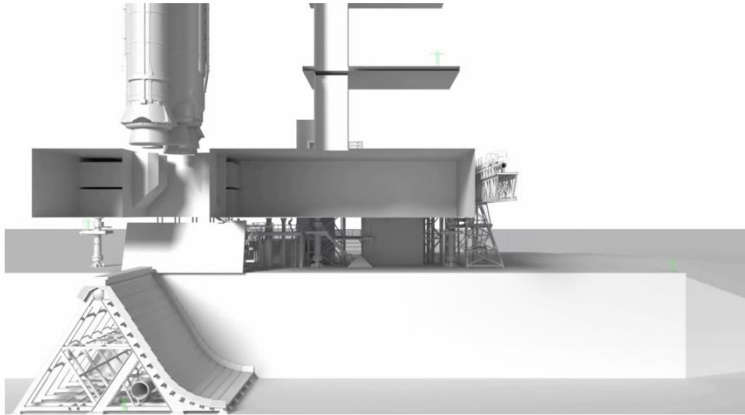
How can we solve simulation problems?

- Usually, analytical solutions do not exist, and we must resort to **numerical** modelling
- Many approaches exist, the most appropriate choice highly depends on the specific problem and scientific domain
- Popular methods (for systems modelled by PDEs) include:
 - finite difference methods, finite element methods, finite volume methods, spectral methods, domain decomposition, mesh-free methods, ...



Mesh for finite element method
Source: COMSOL

Challenges of simulation



Angel et al, Predicting SLS Launch Environment using a Novel Multiphase Formulation (NASA) SC22 (2022) Source: NASA

Required: 500 million grid cells, ran for several weeks on 8,000 cores, generating 400 TB (!)

- Typically, **computational cost** is the main challenge
- Simulation can require **elaborate** parallel software implementations (especially for multi-scale, multi-physics systems) with 10,000s of lines of code
- Significant **human effort** is often required, e.g. in defining high-quality meshes for finite element simulations

How to Contribute to OpenFOAM

OpenFOAM is a large piece of software (of the order of 1 million lines of code) in a complex area of scientific application. Since its open source release in 2004, it has become the CFD software of choice for many thousands of people from industry, government laboratories, academic institutions, etc., who download OpenFOAM and use it for free. With such a large user base, sometimes working on mission-critical applications, we have a responsibility to maintain OpenFOAM as a robust, efficient, and scalable CFD software package for a wide range of applications.

Source: OpenFOAM

Key scientific tasks: inverse problems

$$b = F(a)$$

a = set of input conditions

F = physical model of the system (usually a PDE)

b = resulting properties given F and a

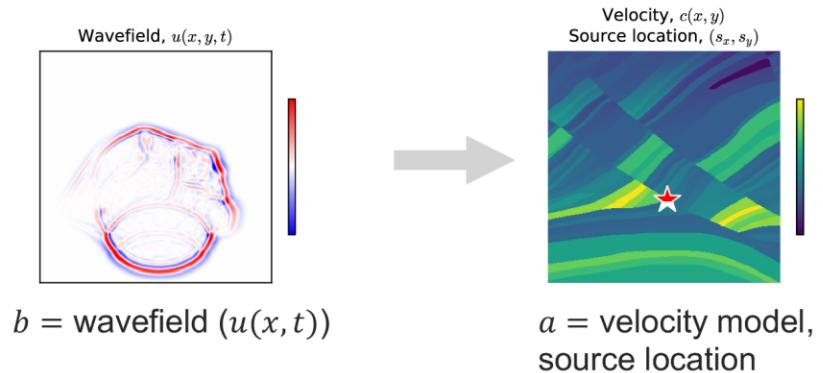
Key scientific tasks: inverse problems

$$b = F(a)$$

Inverse problems are **pervasive** across all domains of science and solving them is essential for many real-world tasks

Example inverse problems:

- Seismic imaging
- Magnetic resonance imaging
- Image denoising
- Estimating infection rates
- Design problems
- ...



$F =$ A method for solving the PDE (e.g. finite difference simulation)

How can we solve inverse problems?

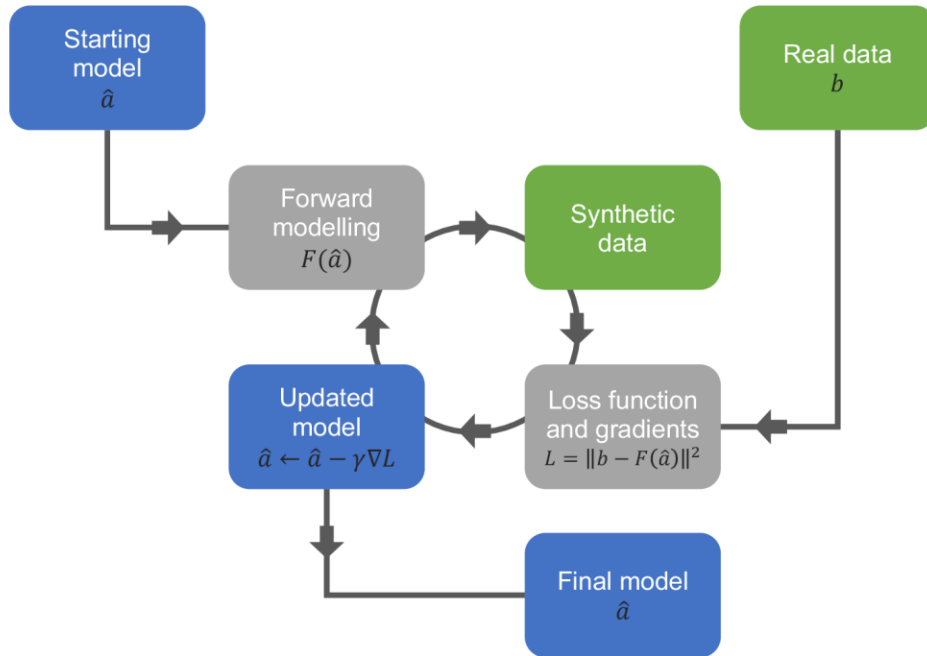
$$b = F(a)$$

- Fundamentally, inverse problems are **search** problem
- It is often useful to frame them as an optimisation problem, for example:

$$\min_{\hat{a}} \|b - F(\hat{a})\|^2$$

- If F is differentiable, one option is to use gradient-based methods (e.g. **gradient descent**)
- Otherwise, we can use gradient-free methods (e.g. evolutionary algorithms, Bayesian optimisation, brute-force search, ...)

Solving inverse problems with gradient descent



$$\min_{\hat{a}} \|b - F(\hat{a})\|^2$$

Loss function is:

$$L(\hat{a}) = \|b - F(\hat{a})\|^2$$

Gradient descent step:

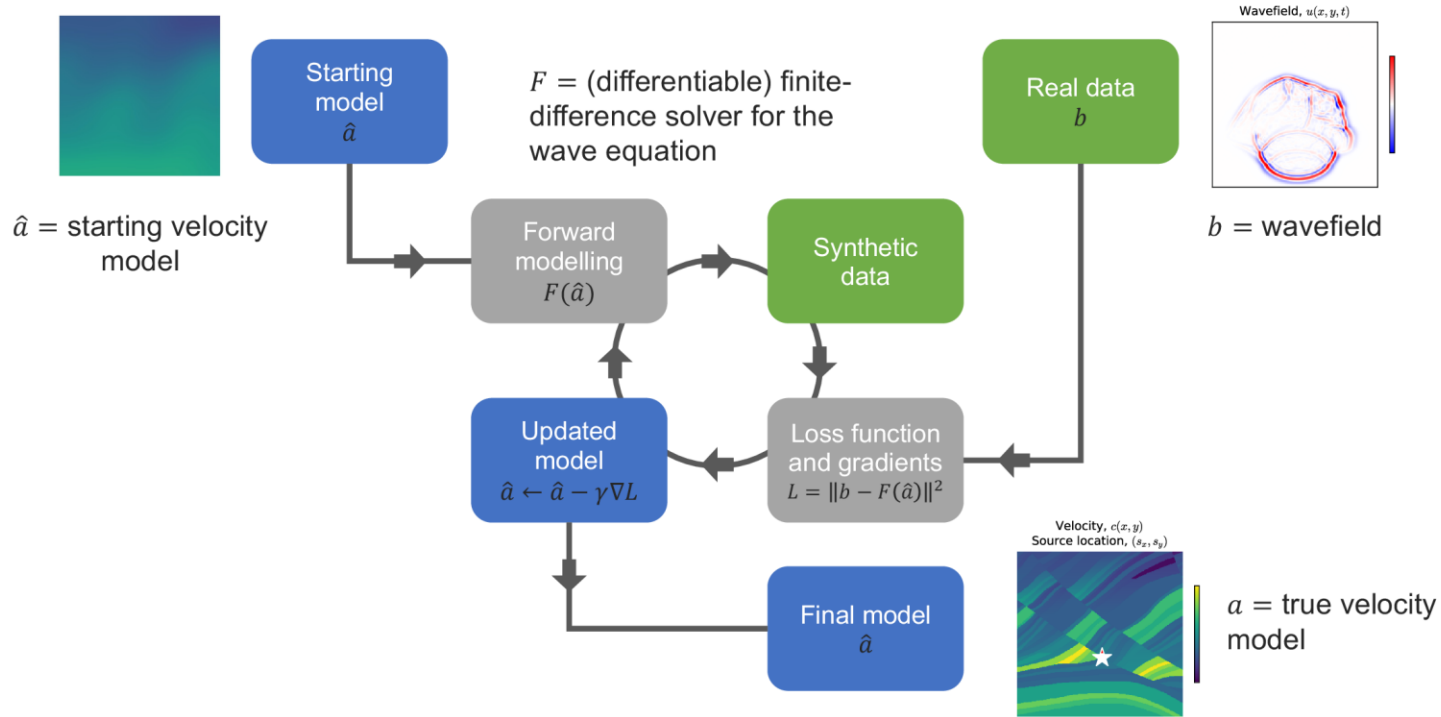
$$\hat{a} \leftarrow \hat{a} - \gamma \nabla L$$

- Requires F to be differentiable

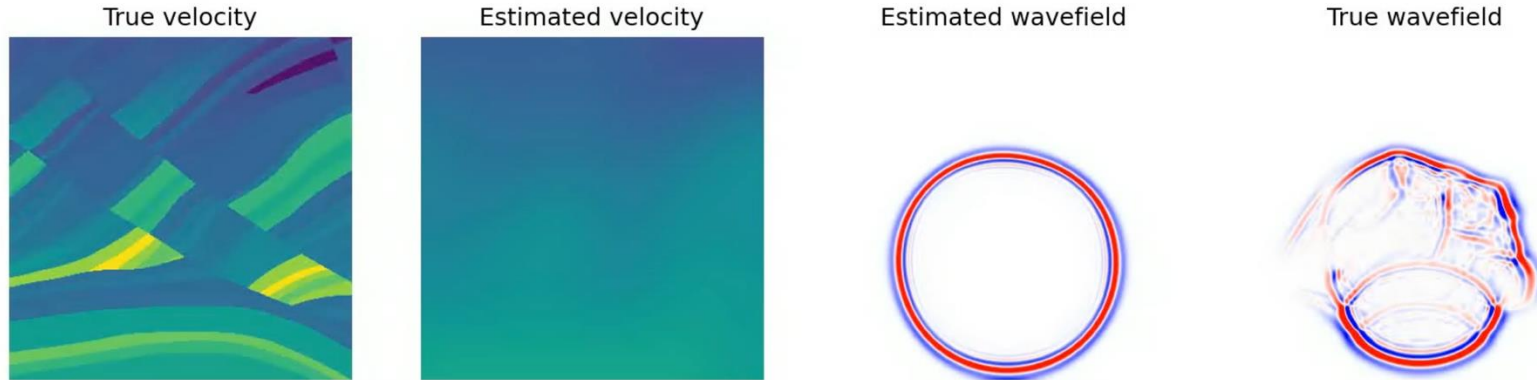


For later: note similarity to training deep neural networks

Solving inverse problems with gradient descent



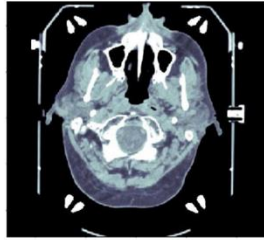
Solving inverse problems with gradient descent



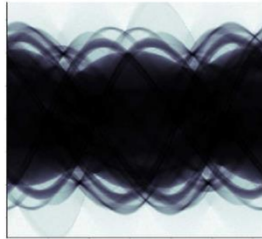
$$\min_{\hat{a}} \|b - F(\hat{a})\|^2$$

- In geophysics, this inverse problem is known as **full waveform inversion**

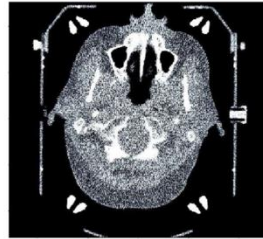
Challenges of inverse problems



Ground truth computed tomography image



Resulting tomographic data (sinogram)



Result of inverse algorithm (filtered back-projection)

$a(x)$

$$F(a)(l) = I_0 \exp\left(-\int_l a(x) dx\right)$$

$\hat{a}(x)$



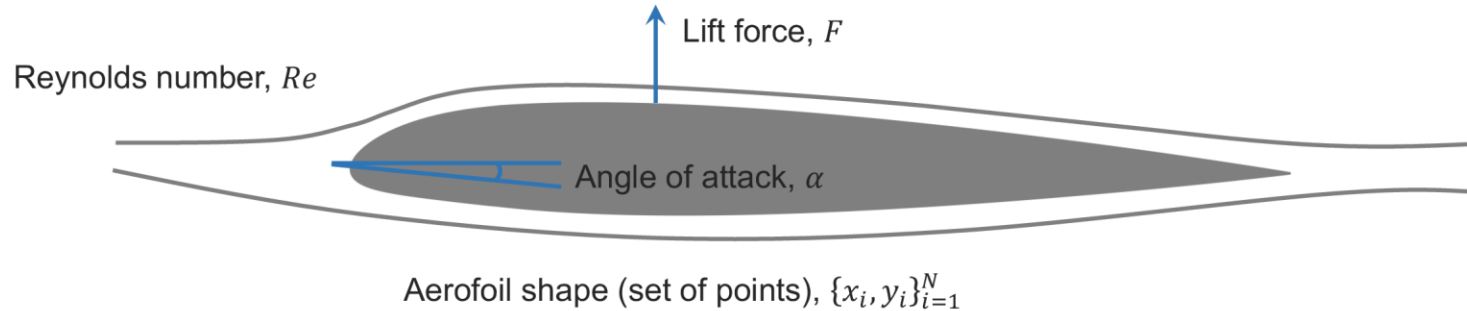
Adler et al, Solving ill-posed inverse problems using iterative deep neural networks, Inverse Problems (2017)

Typically, inverse problems are **incredibly challenging** to solve because:

- They are usually **ill-posed** (not enough information for a unique solution)
- Observed real-world data is usually **noisy** and **sparse**
- Often require forward modelling to be carried out thousands of times – making them extremely **computationally demanding**

Key scientific tasks: design

- Also a type of inverse problem



Design task:

Find $\{x_i, y_i\}_{i=1}^N, \alpha$ which maximize F

Key scientific tasks: discovering unknown physics

$$b = F(a)$$

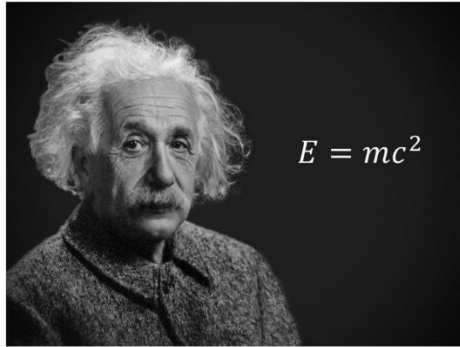
a = set of input conditions

F = physical model of the system (usually a PDE)

b = resulting properties given F and a

Challenges of discovering unknown physics

$$b = F(a)$$



- Historically, F (= laws of physics) has been found through remarkable human **intuition**
- Constant interplay between theory and experiment
- From a computational standpoint, discovering physics is like solving an inverse problem (trying to fit a model to observed data)
- But the model should be **explainable**, **generalisable** and make **novel** predictions

How can deep learning help?

How can deep learning help?

B What are the benefits of using deep learning in science?



Deep learning, a subset of machine learning, has revolutionized various scientific fields by offering numerous benefits. Some of the benefits of using deep learning in science include:

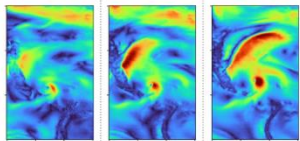
1. **Improved accuracy:** Deep learning algorithms can learn patterns in large amounts of data and make highly accurate predictions. This is particularly useful in fields like medical imaging where human-level accuracy is required to detect diseases and anomalies.
2. **Automation:** Deep learning algorithms can automate tasks that were previously performed manually, freeing up scientists' time to focus on other tasks. For example, in astronomy, deep learning algorithms can be used to automatically identify objects in images, saving researchers hours of manual analysis.
3. **Enhanced understanding:** By revealing patterns and relationships in data, deep learning algorithms can help scientists gain new insights into their fields. This is particularly useful in areas such as genomics where large amounts of data need to be analyzed to uncover important biological insights.
4. **New applications:** Deep learning can be applied to new areas that were previously difficult to tackle using traditional machine learning methods. For example, it has been used in chemistry to predict the properties of new chemical compounds and accelerate drug discovery.
5. **Increased efficiency:** By reducing the time required to analyze data and perform experiments, deep learning can increase the overall efficiency of scientific research. This can lead to new discoveries being made more quickly and cost-effectively.

Overall, deep learning has the potential to significantly enhance the speed and accuracy of scientific research, enabling scientists to make new discoveries and progress more quickly.

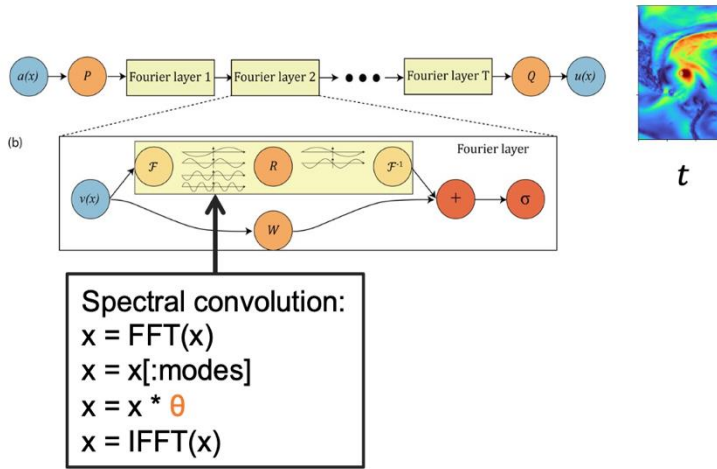
- Accuracy
- Efficiency
- Automation
- Discovery

How Deep Learning can help?

Fourier Neural Operator

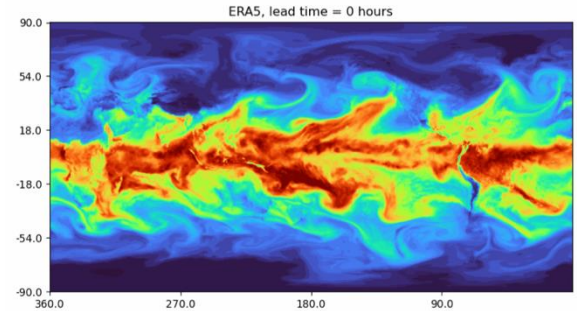
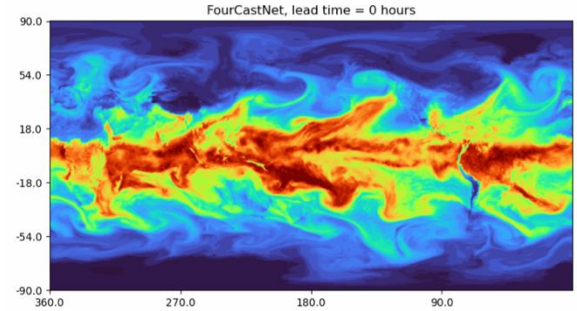


$t - 3, t - 2, t - 1$



Idea: use deep neural network to predict next timestep of multiple atmospheric variables, given previous timesteps

- Use a Fourier neural operator architecture
- **4-5 orders of magnitude** faster than numerical simulation

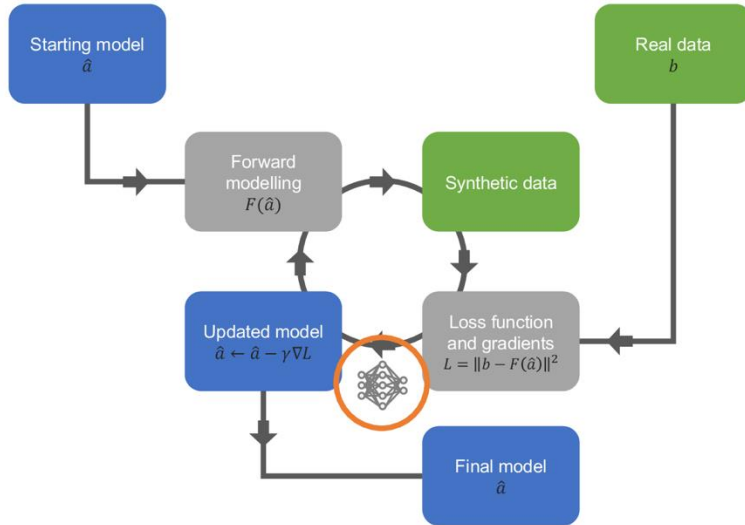


Pathak et al, FourCastNet: A Global Data-driven High-resolution Weather Model using Adaptive Fourier Neural Operators, ArXiv (2022)

Li et al, Fourier Neural Operator for Parametric Partial Differential Equations, ICLR (2021)

Example - inverse problems

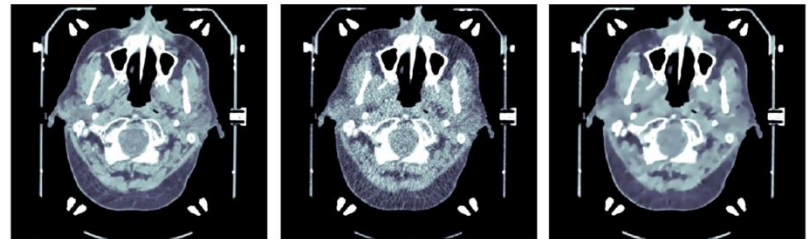
Learned gradient descent



Idea: use a neural network to **update** gradient before descent step

1. Initialise \hat{a}
2. Loop:
 1. Compute gradient, $g = \frac{\partial L(\hat{a})}{\partial \hat{a}}$
 2. Update gradient, $g \leftarrow NN(g, \hat{a}, b, \theta)$
 3. Take step, $\hat{a} \leftarrow \hat{a} - \gamma g$

Network is trained using many example inverse problems, and differentiating through entire algorithm **end-to-end**



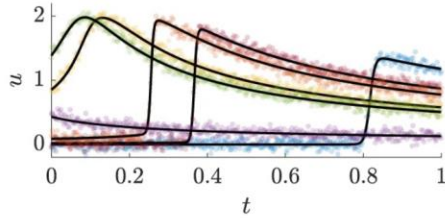
Ground truth

Traditional inversion

Learned gradient descent

Adler et al, Solving ill-posed inverse problems using iterative deep neural networks, Inverse Problems (2017)

Example - equation discovery



Ground truth: $u_t + uu_x - 0.0032u_{xx} = 0$

Discovered: $u_t + 1.002uu_x - 0.0032u_{xx} = 0$

$$u_t = \Lambda \phi$$

$$\Lambda = (\lambda_1 \lambda_2 \lambda_3 \lambda_4 \lambda_5 \dots) \quad \phi = \begin{pmatrix} \hat{u}_x \\ \hat{u}_{xx} \\ \hat{u}_t \\ \hat{u}_{tt} \\ \hat{u}_{xt} \\ \dots \end{pmatrix}$$

Idea: fit a neural network to observed data, then regress over a **library** of gradients to “discover” underlying equations

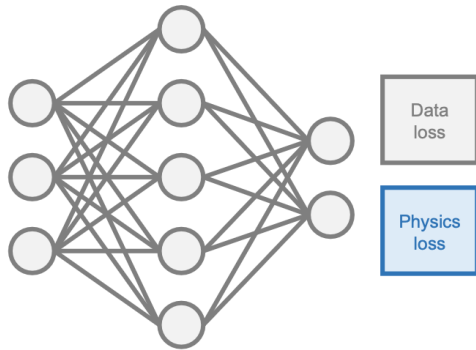
Input: $\tilde{u}(x, t) =$ (noisy) observational data

1. Fit a neural network to data, $\hat{u} = NN(x, t; \theta) \approx \tilde{u}$, using supervised learning
2. Compute various gradients of network, for example $\hat{u}_x, \hat{u}_{xx}, \hat{u}_t, \hat{u}_{tt}, \hat{u}_{xt} \dots$, at many random (x, t) locations
3. Carry out a (sparse) **linear regression** over combinations of these gradients to “discover” underlying equation

The state-of-the-art: scientific machine learning

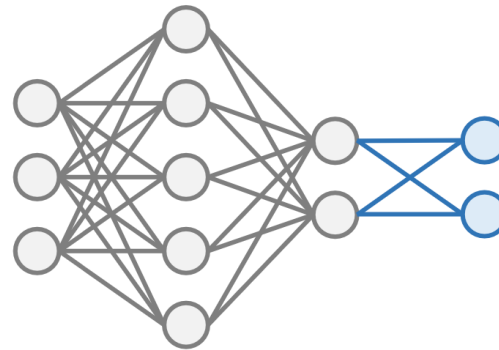
Ways to incorporate scientific principles into machine learning

Loss function



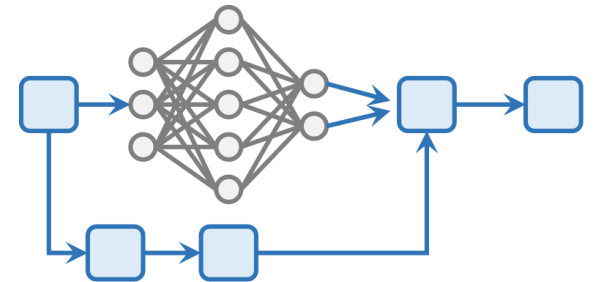
Example:
Physics-informed neural networks
(add governing equations to loss
function)

Architecture



Example:
Encoding symmetries / conservation laws
(e.g. energy conservation, rotational
invariance)

Hybrid approaches



Example:
Neural differential equations
(incorporating neural networks into PDE
models)

A plethora of SciML techniques

Scientific task

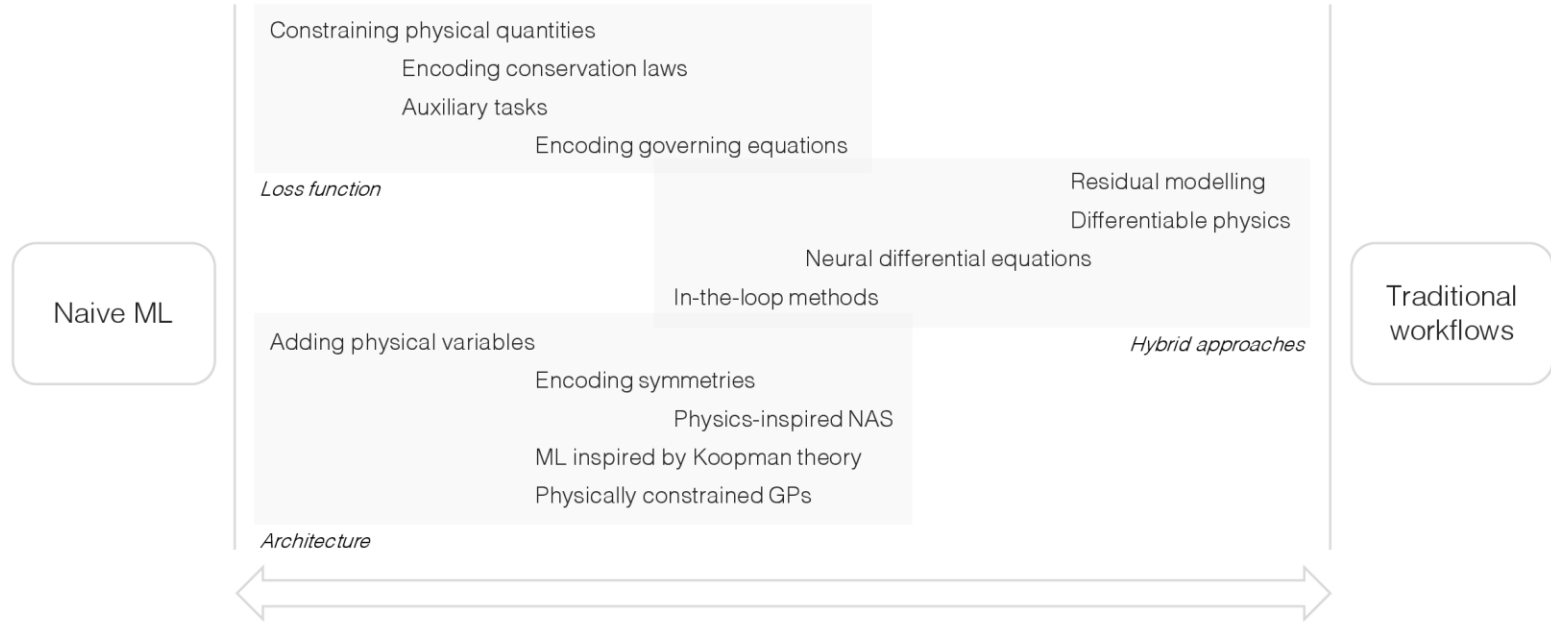
	Forward simulation	Inversion	Equation discovery
Architecture			
Adding physical variables	Daw et al.		
Encoding symmetries	Ling et al., Wang et al., Anderson et al., Schütt et al.		Udrescu et al.
Physics-inspired NAS	Ba et al., Panju and Ghodsi		
ML inspired by Koopman theory	Geneva and Zabarab, Lusch et al.		
Physically constrained GPs		Raissi et al., Raissi and Karniadakis	
Other approaches	Jumper et al., Mohan et al.		
Loss function			
Constraining physical quantities	Karpatne et al., Zhang et al., Benjamin et al., Erichson et al., Xie et al., Brehmer et al.		
Encoding conservation laws	Beucler et al., Zeng et al.		Greydanus et al., Toth et al., Cranmer et al.
Auxiliary tasks	de Oliveira et al.		
Encoding governing equations	Raissi et al., Jin et al., Jin et al., Chen et al., Kharazmi et al., Yang et al., Wang et al., Wang et al., Li et al., Zhu et al., Geneva and Zabarab, Gao et al.		Chen et al., Champion et al.
Hybrid approaches			
Residual modelling	Pawar et al.	Jiang et al.	
Differentiable physics		Ren et al., Minkov et al., Würfl et al., Zhang et al.	
Neural differential equations			Chen et al., Rackauckas et al., Long et al.
In-the-loop methods	Um et al., Rasp et al.	Adler and Öktem, Morningstar et al., Hammernik et al., Li et al., Lunz et al., Bora et al., Mosser et al.	

SciML technique

Hamiltonian neural networks
 Learned sub-grid processes
 Hidden physics models
 DeepONets
 Differentiable simulation
 Fourier neural operators
 Encoding conservation laws
 Encoding physical symmetries
 Solver-in-the-loop
 Physics-constrained Gaussian processes
 Physics-informed neural networks
 AI Feynman
 AlphaFold
 Learned regularisation
 Physics-informed neural operators
 Neural ODEs
 PDE-Net
 Algorithm unrolling

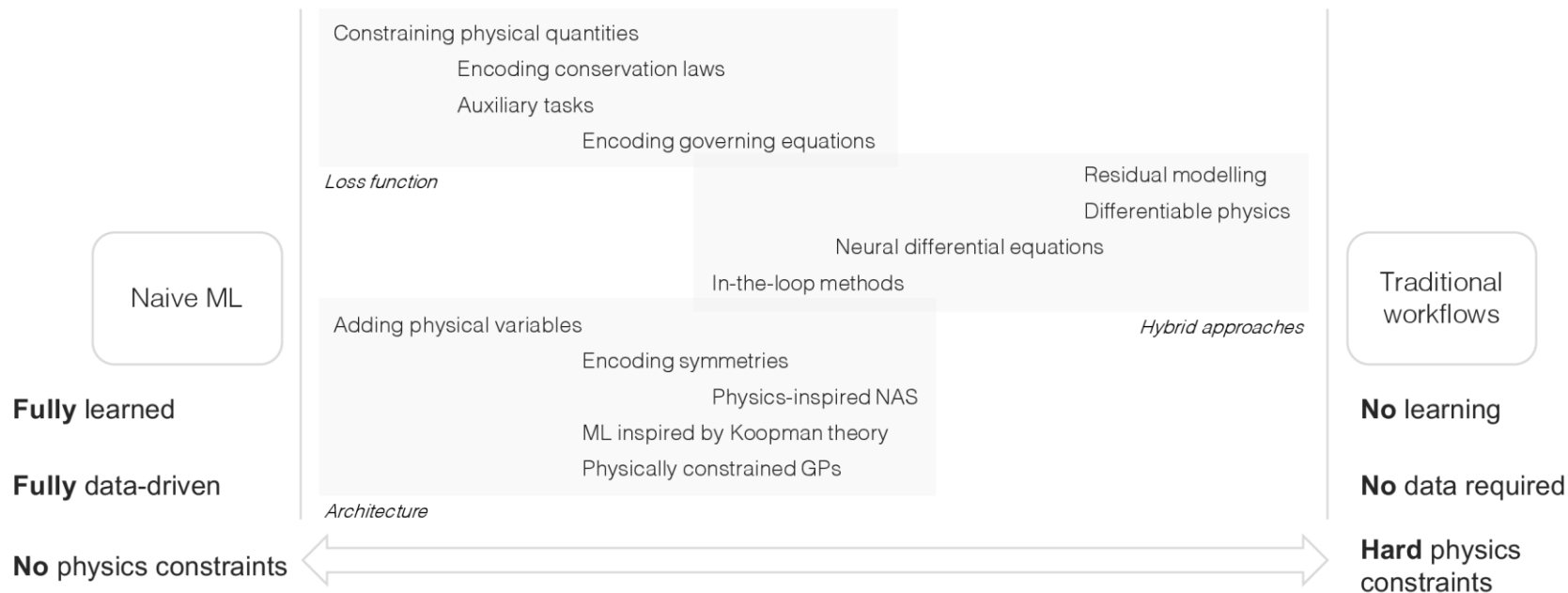
Source: B Moseley, Physics-informed machine learning: from concepts to real-world applications, PhD thesis, 2022

A plethora of SciML techniques



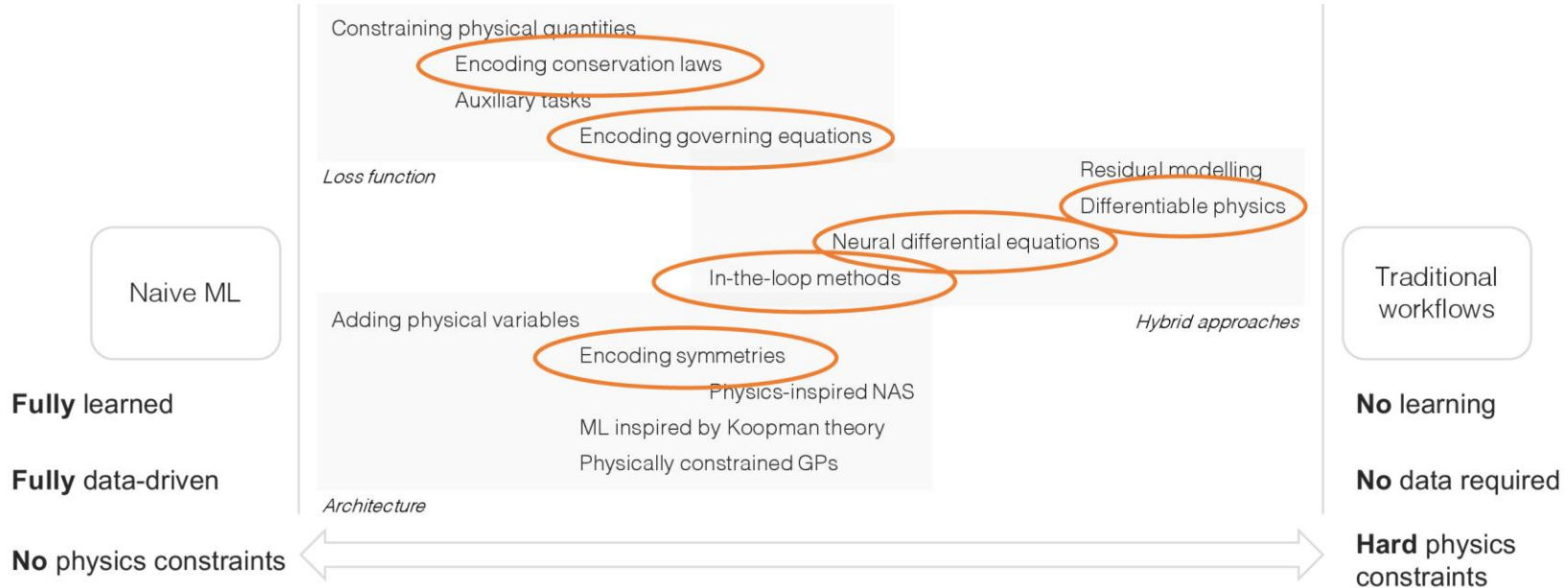
Source: B Moseley, Physics-informed machine learning: from concepts to real-world applications, PhD thesis, 2022

A plethora of SciML techniques

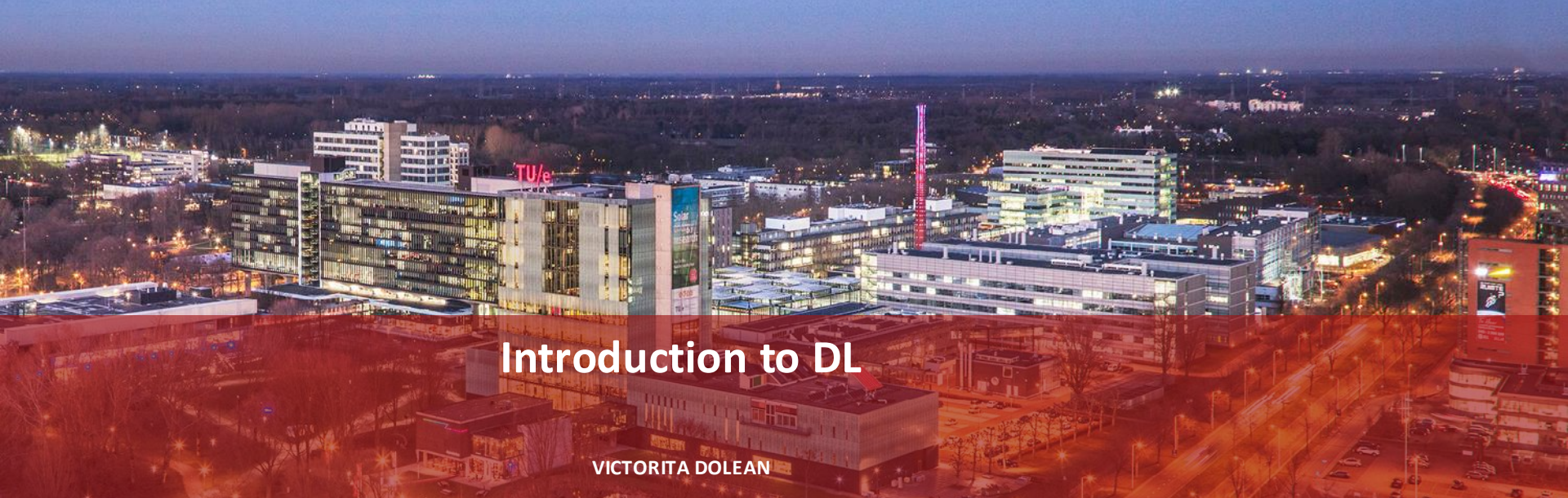


Source: B Moseley, Physics-informed machine learning: from concepts to real-world applications, PhD thesis, 2022

A plethora of SciML techniques



Source: B Moseley, Physics-informed machine learning: from concepts to real-world applications, PhD thesis, 2022

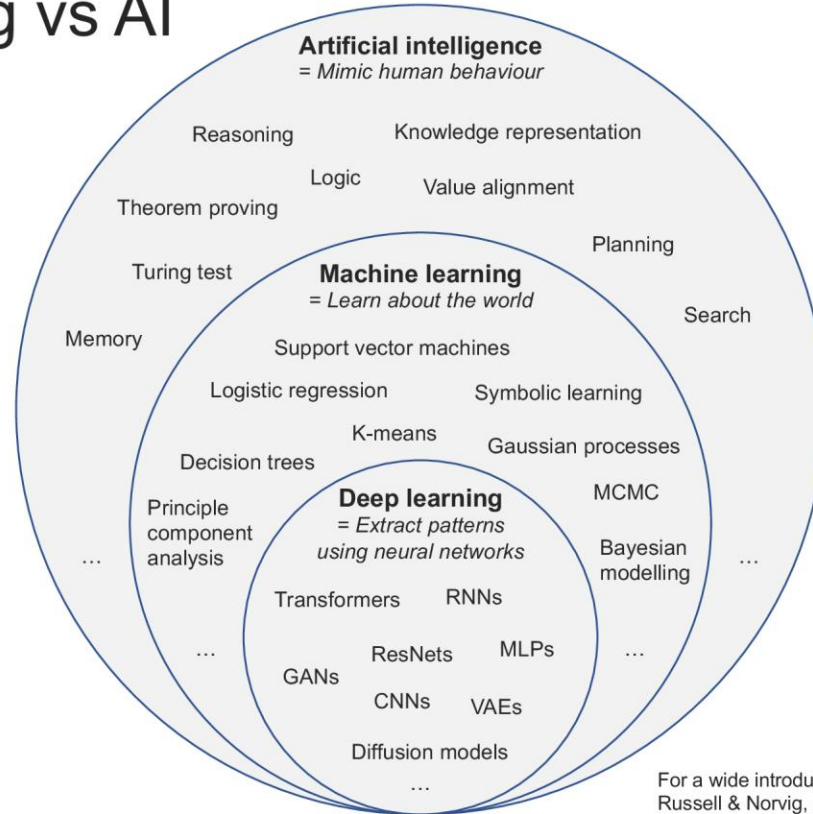


Introduction to DL

VICTORITA DOLEAN



Deep learning vs AI

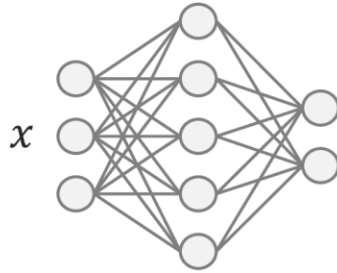


For a wide introduction to AI, see for example:
Russell & Norvig, Artificial Intelligence: A Modern Approach

What is a neural network?

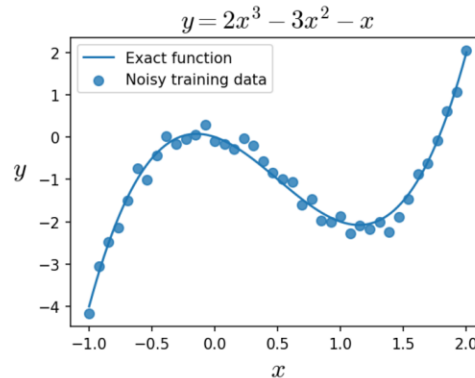


Neural networks are simply **flexible functions** fit to data



$$\hat{y} = NN(x; \theta)$$

Example dataset:



Goal: given training data, find a function (with flexible parameters θ) which approximates the true function,

$$\hat{y} = NN(x; \theta) \approx y(x)$$

Function fitting

Simple polynomial regression

$$\hat{y}(x; \theta) = \theta_4 x^3 + \theta_3 x^2 + \theta_2 x + \theta_1$$

To fit, use least-squares:

$$\theta^* = \min_{\theta} \sum_i^N (\hat{y}(x_i; \theta) - y_i)^2 \quad (1)$$

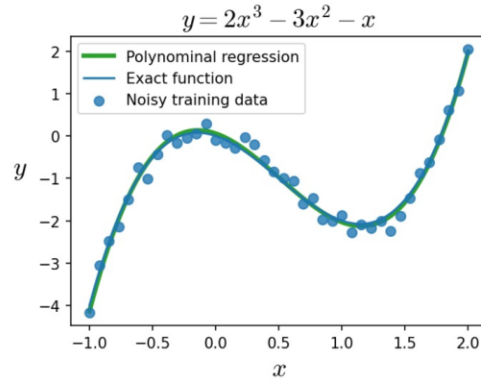
Re-write using linear algebra:

$$\begin{pmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \dots \end{pmatrix} = \begin{pmatrix} 1 & x_1 & x_1^2 & x_1^3 \\ 1 & x_2 & x_2^2 & x_2^3 \\ \dots & \dots & \dots & \dots \end{pmatrix} \begin{pmatrix} \theta_1 \\ \theta_2 \\ \theta_3 \\ \theta_4 \end{pmatrix} \text{ or } \hat{Y} = \Phi^T \theta$$

$$\theta^* = \min_{\theta} \|\Phi^T \theta - Y\|^2$$

In this case, it can be shown (1) has an analytical solution:

$$\theta^* = (\Phi^T \Phi)^{-1} \Phi^T Y$$



Neural network regression

$$\hat{y}(x; \theta) = NN(x; \theta)$$

To fit, use least-squares:

$$\theta^* = \min_{\theta} \sum_i^N (NN(x_i; \theta) - y_i)^2 \quad (2)$$

In general, no analytical solution to (2) exists, so we must use **optimisation**

For example, gradient descent:

$$\theta_j \leftarrow \theta_j - \gamma \frac{\partial \sum_i^N (NN(x_i; \theta) - y_i)^2}{\partial \theta_j}$$

or equally

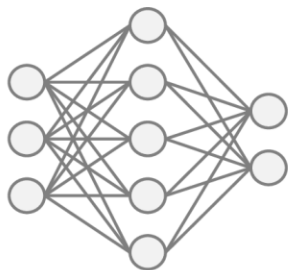
$$\theta_j \leftarrow \theta_j - \gamma \frac{\partial L(\theta)}{\partial \theta_j}$$

Where γ is the learning rate and $L(\theta)$ is the **loss function**

Neural network architecture

So, what exactly is $\hat{y} = NN(x; \theta)$?

This depends on the network **architecture** you choose (CNN, ResNet, Transformer, ... etc)



2-layer MLP

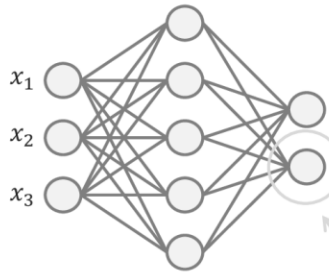
The most basic architecture is the **multilayer perceptron** (MLP) (aka **fully connected network**)

For example, a 2-layer MLP is defined as:

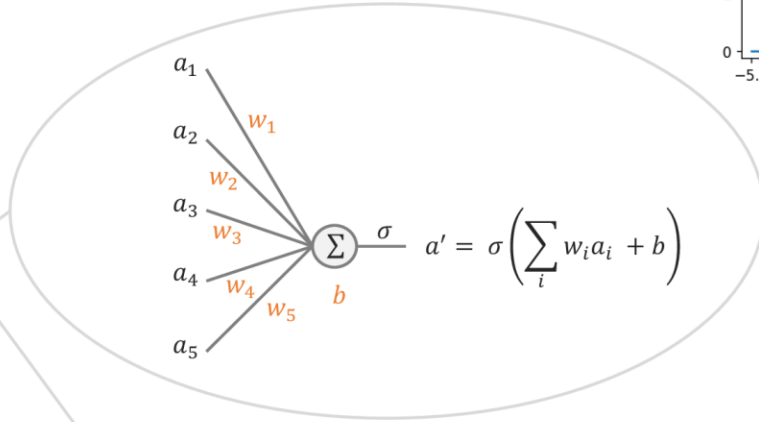
$$NN(x; \theta) = W_2 \sigma(W_1 x + b_1) + b_2$$

Where x is an input vector, W_1 and W_2 are learnable weight matrices, b_1 and b_2 are learnable bias vectors, and σ is an activation function, for example, $\sigma = \tanh(\cdot)$

Biological inspiration



2-layer MLP

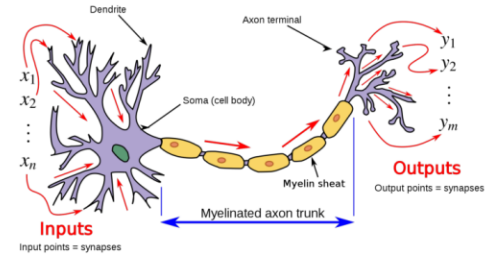
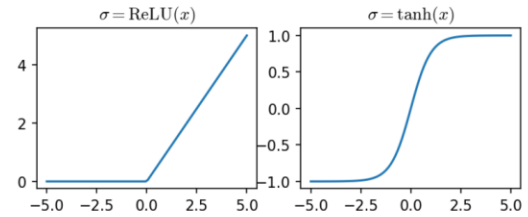


For last layer:

$$\begin{pmatrix} a'_1 \\ a'_2 \end{pmatrix} = \sigma \left(\begin{pmatrix} w_{11} & w_{12} & w_{13} & w_{14} & w_{15} \\ w_{21} & w_{22} & w_{23} & w_{24} & w_{25} \end{pmatrix} \begin{pmatrix} a_1 \\ a_2 \\ a_3 \\ a_4 \\ a_5 \end{pmatrix} + \begin{pmatrix} b_1 \\ b_2 \end{pmatrix} \right)$$

Entire network:

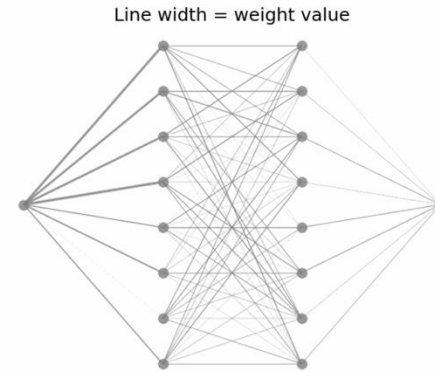
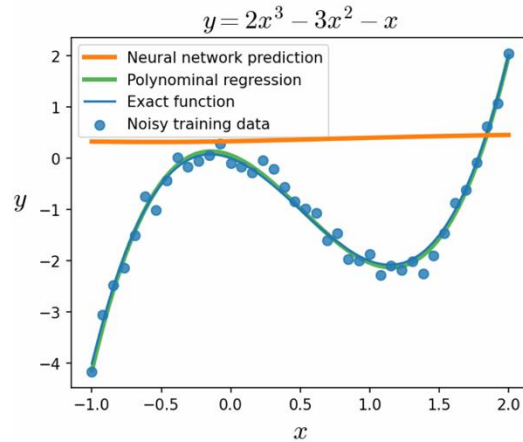
$$NN(\mathbf{x}; \theta) = \sigma(W_2 \sigma(W_1 \mathbf{x} + \mathbf{b}_1) + \mathbf{b}_2) = \mathbf{f} \circ \mathbf{g}(\mathbf{x}; \theta)$$



Biological neuron
(Source: Wikipedia)

Polynomial regression example

Training step 0



$$NN(x; \theta) = W_3(\sigma(W_2\sigma(W_1x + b_1) + b_2) + b_3)$$

Trained using gradient descent

$$\theta_j \leftarrow \theta_j - \gamma \frac{\partial \sum_i^N (NN(x_i; \theta) - y_i)^2}{\partial \theta_j}$$

Universal approximation

So why not just use linear regression?



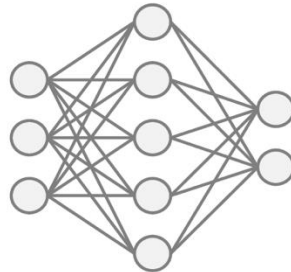
Neural networks are simply **flexible functions** fit to data



With enough parameters, neural networks can approximate any* arbitrarily complex function
= **universal approximation**

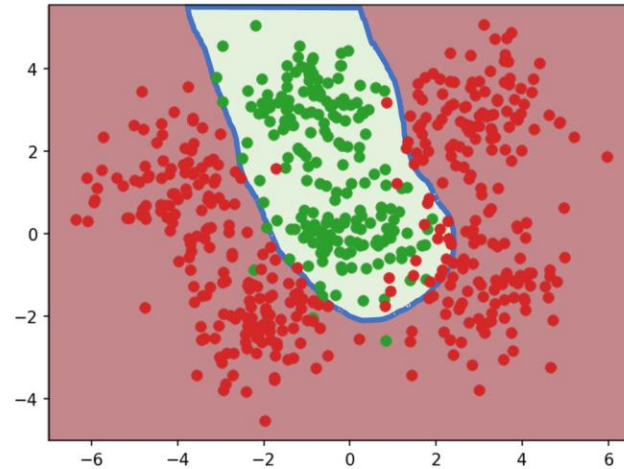
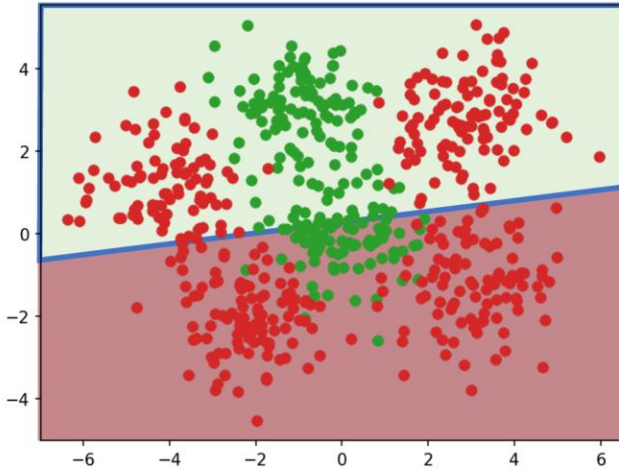


x = array of RGB values



$$\hat{y} = P(\text{dog} | x) = 1$$

Importance of activation functions

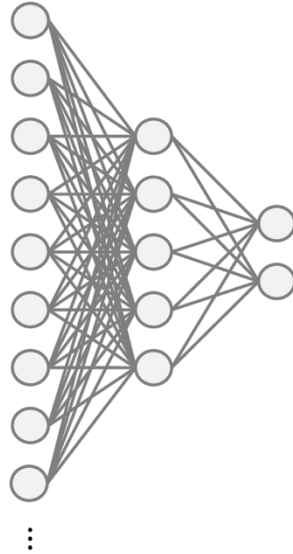


Non-linearities allow us to approximate arbitrary **non-linear** functions

MLPs use lots of parameters



=> Flatten to 1D =>



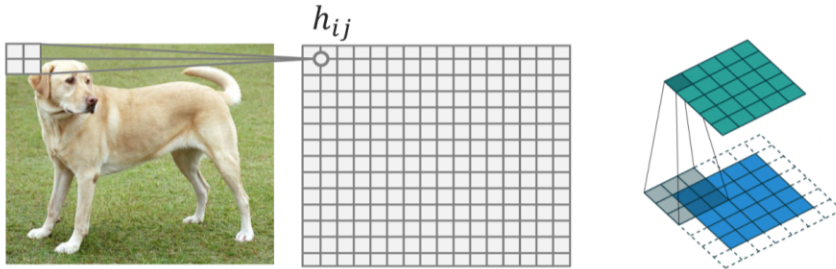
$$NN(x; \theta) = W_3(\sigma(W_2\sigma(W_1x + b_1) + b_2) + b_3)$$

Assume the image has shape 128×128 , and we have 100 hidden units in the first layer, then W_1 has shape $(100 \times (128 \times 128)) = (100 \times 16,384)$

= 1.6M parameters!

=> A simple MLP image classifier can have millions of parameters

Convolutional neural network (CNN)



Convolutional neural networks honor the **spatial correlations** in their inputs

Each neuron;

- Has a **limited** field of view
- **Shares** the same weights as the other neurons in the layer
- Mathematically, CNNs use cross-correlation

CNNs have translation equivariance (an inductive bias)

$$NN(x; \theta) = W_3 * (\sigma(W_2 * \sigma(\underbrace{W_1 * x + b_1}_h) + b_2) + b_3)$$

$$h_{ij} = \sum_{i'}^l \sum_{j'}^m W_{i'j'} x_{i+i', j+j'} + b$$

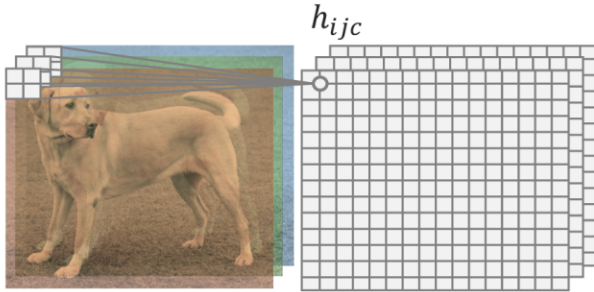
Let the size of the convolutional filter be 3 x 3

Then W_1 has shape (3 x 3)

= 9 parameters! (much, much smaller than a MLP)

Image source:
github.com/vdumoulin/conv_arithmetic

Convolutional neural network (CNN)



In practice, CNNs are usually extended so they can have multiple **channels** in the inputs and outputs of each layer

e.g. (R,G,B) image as input, where each channel is a color

Also:

- 1D and 3D CNNs follow analogously
- And we can add dilations and strides too

Then the convolutional layer is defined by:

$$h_{ijc} = \sum_{i'}^l \sum_{j'}^m \sum_{c'}^{C_{\text{in}}} W_{i'j'c'c} x_{i+i', j+j', c'} + b_c$$

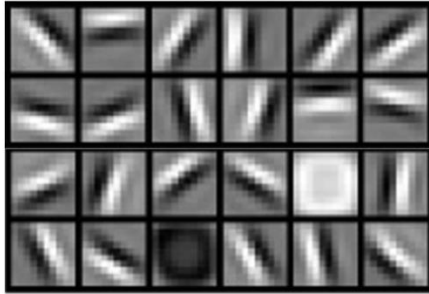
Let the size of the convolutional filter be 3×3

Then W has shape $(3 \times 3 \times C_{\text{in}} \times C_{\text{out}})$

= 81 parameters for 3 input and 3 output channels

Deep CNNs

First layer



Second layer



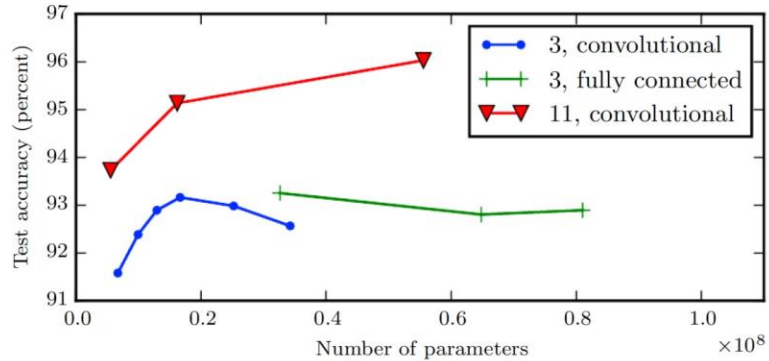
Third layer



Deep CNNs learn **hierarchical** features

Lee et al, Unsupervised Learning of Hierarchical Representations with Convolutional Deep Belief Networks, Communications of the ACM (2011)

Depth is key



Goodfellow et al, Multi-digit number recognition from street view imagery using deep convolutional neural networks, ICLR (2014)

Empirically, deep neural networks perform better than shallow neural networks

=> encode a very general belief that the true function is **composed** of simpler functions

How do we train neural networks?

Gradient descent

$$\hat{y}(x; \theta) = NN(x; \theta)$$

To fit, use least-squares:

$$\theta^* = \min_{\theta} \sum_i^N (NN(x_i; \theta) - y_i)^2 \quad (2)$$

In general, no analytical solution to (2) exists, so we must use **optimisation**

For example, gradient descent:

$$\theta_j \leftarrow \theta_j - \gamma \frac{\partial \sum_i^N (NN(x_i; \theta) - y_i)^2}{\partial \theta_j}$$

or equally

$$\theta_j \leftarrow \theta_j - \gamma \frac{\partial L(\theta)}{\partial \theta_j}$$

Where γ is the learning rate and $L(\theta)$ is the **loss function**

Note that

$$\frac{\partial L(\theta)}{\partial \theta_j} = \sum_i^N 2(NN(x_i; \theta) - y_i) \frac{\partial NN(x_i; \theta)}{\partial \theta_j}$$

Let's consider a fully connected network

$$NN(x; \theta) = W_3(\sigma(W_2 \sigma(\underbrace{W_1 x + b_1}_{\mathbf{g}}) + b_2) + b_3) = \mathbf{f} \circ \mathbf{g} \circ \mathbf{h}(x; \theta)$$

\mathbf{h}

How do we calculate $\frac{\partial NN(x_i; \theta)}{\partial W_1}$?

Note \mathbf{f} , \mathbf{g} , and \mathbf{h} are vector functions =>

Use the **multivariate chain rule** (= matrix multiplication of **Jacobians**)

$$\frac{\partial NN}{\partial W_1} = \frac{\partial \mathbf{f}}{\partial \mathbf{g}} \frac{\partial \mathbf{g}}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial W_1}$$

$$J = \frac{\partial \mathbf{f}}{\partial \mathbf{g}} = \begin{pmatrix} \frac{\partial f_1}{\partial g_1} & \dots & \frac{\partial f_1}{\partial g_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial g_1} & \dots & \frac{\partial f_m}{\partial g_n} \end{pmatrix}$$

Backpropagation

Forward pass:

$$\mathbf{x}_i \rightarrow \mathbf{h}_i = W_1 \mathbf{x}_i + \mathbf{b}_1 \rightarrow \mathbf{g}_i = W_2 \sigma(\mathbf{h}_i) + \mathbf{b}_2 \rightarrow f_i = W_3 \sigma(\mathbf{g}_i) + \mathbf{b}_3$$

Backward pass:

In practice:

 PyTorch



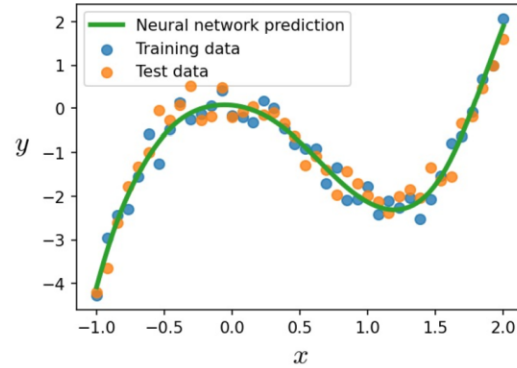
Autodifferentiation tracks all your forward computations and their gradients and applies the chain rule automatically for you, so you don't have to worry!

Train loss vs Train error vs Test error

Training loss:

$$L(\theta) = \frac{1}{N} \sum_i^N (NN(x_i; \theta) - y_i)^2$$

$$D_{\text{train}} = \{(x_1, y_1), \dots, (x_N, y_N)\}$$
$$x, y \sim p(x, y)$$



Train error:

$$L(\theta) = \frac{1}{N} \sum_i^N |NN(x_i; \theta) - y_i|$$

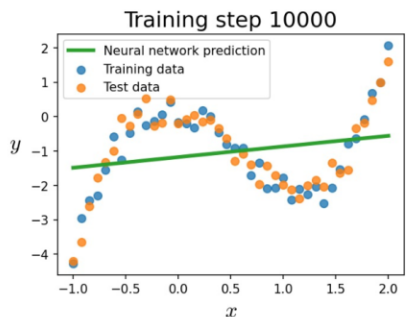
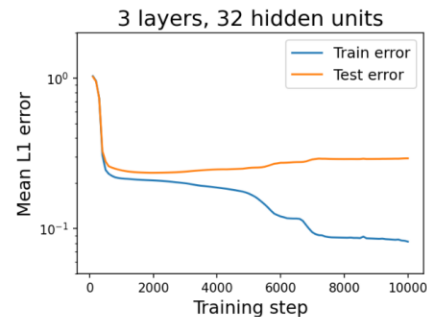
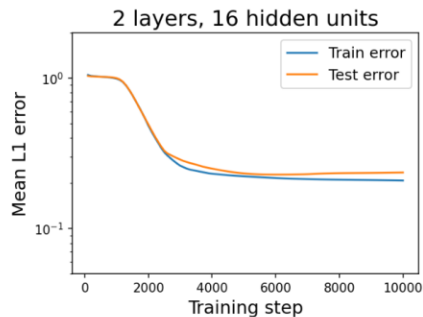
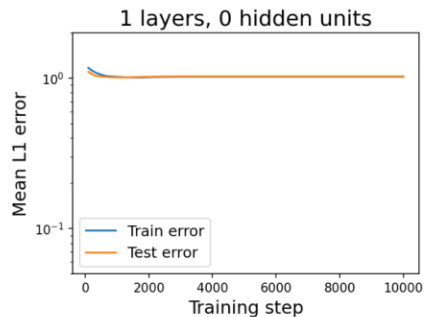
$$D_{\text{train}} = \{(x_1, y_1), \dots, (x_N, y_N)\}$$
$$x, y \sim p(x, y)$$

Test error:

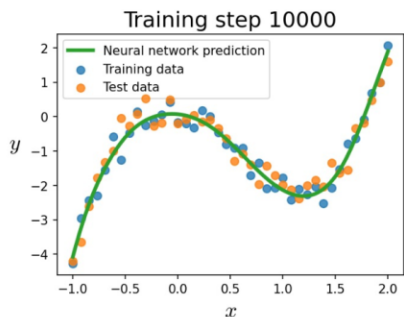
$$L(\theta) = \frac{1}{M} \sum_i^M |NN(x_i; \theta) - y_i|$$

$$D_{\text{test}} = \{(x_1, y_1), \dots, (x_M, y_M)\}$$
$$x, y \sim p(x, y)$$

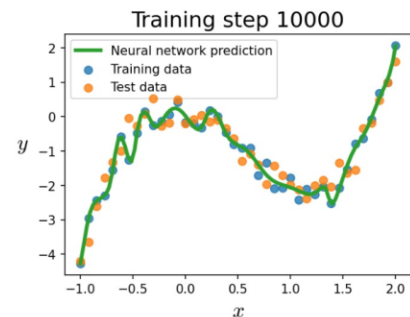
Overfitting vs underfitting



Underfit



Ideal model



Overfit

Understanding error sources

Training loss:

$$L(\theta) = \frac{1}{N} \sum_i^N (NN(x_i; \theta) - y_i)^2$$
$$D = \{(x_1, y_1), \dots, (x_N, y_N)\}$$

This is also known as the **empirical loss**, and can be more generally written as:

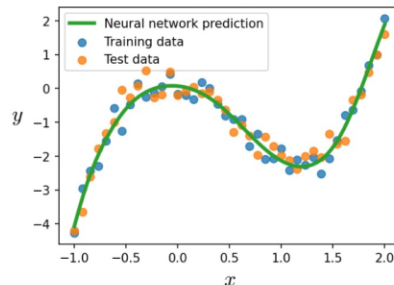
$$L(\theta) = \frac{1}{N} \sum_i^N l(NN(x_i; \theta), y_i), \quad x, y \sim p(x, y)$$

But what we really want to minimise is the **expected loss**:

$$\mathcal{L}(\theta) = \iint l(NN(x; \theta), y) p(x, y) dx dy$$
$$= E_{(x,y) \sim p} [l(x, y; \theta)]$$

Why is this impossible in practice?

Because, usually, we only have access to a finite number of samples from $p(x, y)$ (our training dataset, D)



This leads to

1) **Estimation error** (finite amount of training data) (aka generalisation error)

$$\mathcal{E}_{\text{est}} = \mathcal{L}(\underline{NN}) - \mathcal{L}(\underline{NN^*})$$

NN which minimises
empirical loss

NN which minimises
expected loss

We also have

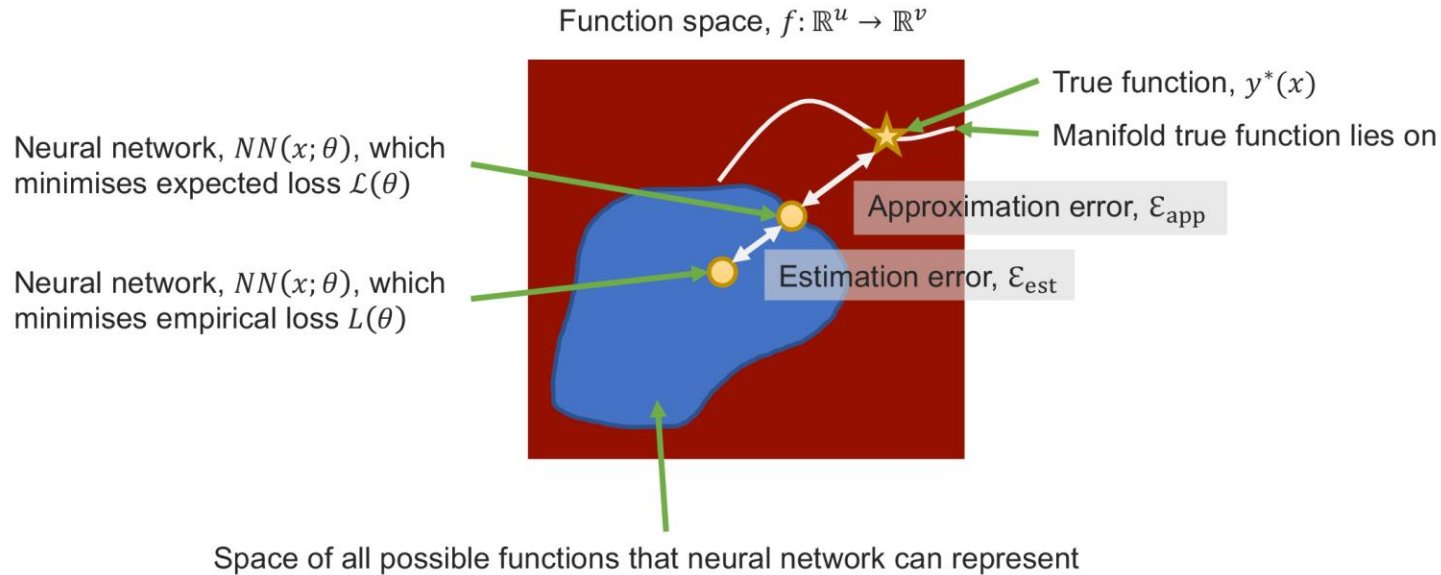
2) **Approximation error** (limited expressivity of neural network)

$$\mathcal{E}_{\text{app}} = \mathcal{L}(\underline{NN^*}) - \mathcal{L}(\underline{y^*})$$

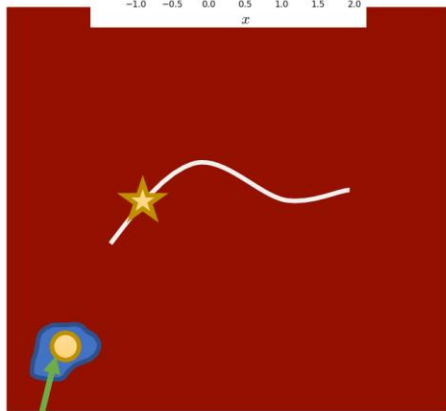
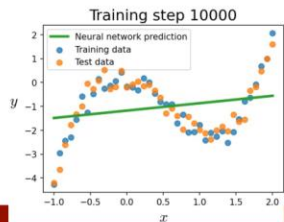
NN which minimises
expected loss

True function which minimises
expected loss

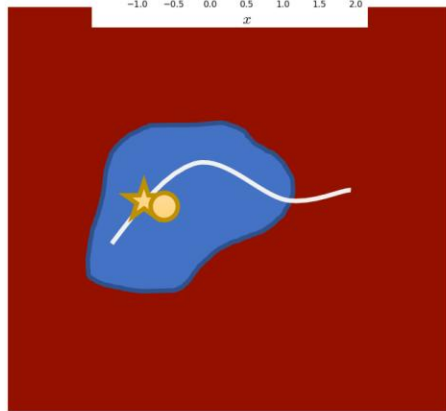
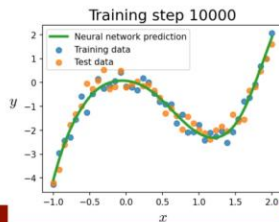
Function space



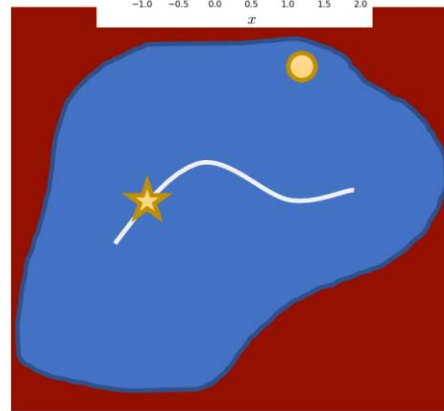
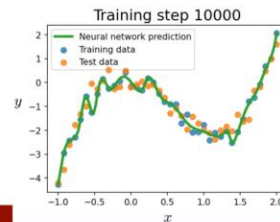
Function space



Underfit



Ideal model



Overfit

Neural network, $NN(x; \theta)$, which minimises empirical loss $L(\theta)$

Decomposing expected loss

Let $h(x, D) = NN(x; \theta)$ be the function learned by the neural network using a **particular** training dataset $D = \{(x_1, y_1), \dots, (x_N, y_N)\}$

Then we would like to understand what the expected loss is averaged over **all possible training datasets**, i.e.

$$E_{(x,y) \sim p, D \sim p^N} [l(x, y, D)]$$

Assume a simple least-squares loss function, e.g.

$$E_{x,y,D} [(h(x, D) - y)^2]$$

Then it can be shown*:

$$E_{x,y,D} [(h(x, D) - y)^2] = \underbrace{E_x [(\bar{h}(x) - \bar{y}(x))^2]}_{\text{Bias}^2} + \underbrace{E_{x,D} [(h(x, D) - \bar{h}(x))^2]}_{\text{Variance}} + \underbrace{E_{x,y} [(\bar{y}(x) - y)^2]}_{\text{Irreducible noise}}$$

Where

$$\bar{y}(x) = E_{y|x}[y] = \int y p(y|x) dy = \text{Average observation}$$

$$\bar{h}(x) = E_D[h(x, D)] = \int h(x, D) p(D) dD = \text{Average model}$$

*for full derivation, see e.g. <https://www.cs.cornell.edu/courses/cs4780/2018fa/lectures/lecturenote12.html>

Decomposing expected loss

From previous slide:

$$E_{x,y,D}[(h(x,D) - y)^2] = E_x \left[\underbrace{(\bar{h}(x) - \bar{y}(x))^2}_{\text{Bias}^2} \right] + E_{x,D} \left[\underbrace{(h(x,D) - \bar{h}(x))^2}_{\text{Variance}} \right] + E_{x,y} \left[\underbrace{(\bar{y}(x) - y)^2}_{\text{Irreducible noise}} \right]$$

Where

$$\bar{y}(x) = E_{y|x}[y] = \int y p(y|x) dy = \text{Average observation}$$

$$\bar{h}(x) = E_D[h(x,D)] = \int h(x,D) p(D) dD = \text{Average model}$$

Bias = “average” model error

Variance = sensitivity of model prediction to training data

Irreducible noise = inherent noise in the training observations; you can never beat this ($\bar{y}(x)$ is the best possible predictor of y)

Decomposing expected loss

From previous slide:

$$E_{x,y,D}[(h(x,D) - y)^2] = \underbrace{E_x \left[(\bar{h}(x) - \bar{y}(x))^2 \right]}_{\text{Bias}^2} + \underbrace{E_{x,D} \left[(h(x,D) - \bar{h}(x))^2 \right]}_{\text{Variance}} + \underbrace{E_{x,y} [(\bar{y}(x) - y)^2]}_{\text{Irreducible noise}}$$

Where

$$\bar{y}(x) = E_{y|x}[y] = \int y p(y|x) dy = \text{Average observation}$$

$$\bar{h}(x) = E_D[h(x,D)] = \int h(x,D) p(D) dD = \text{Average model}$$

Bias = “average” model error

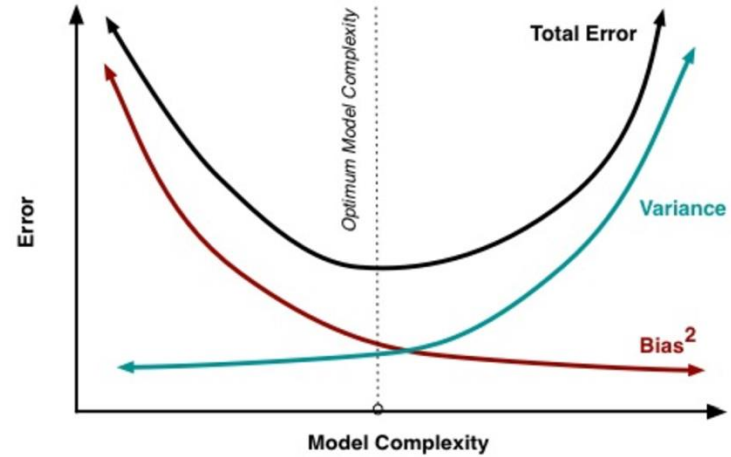
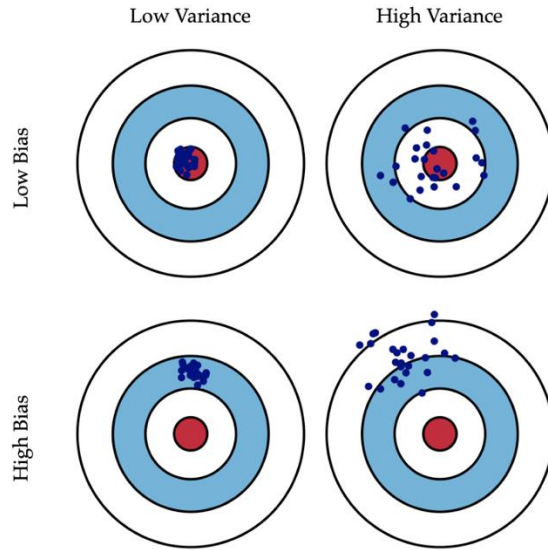
Variance = sensitivity of model prediction to training data

Irreducible noise = inherent noise in the training observations; you can never beat this ($\bar{y}(x)$ is the best possible predictor of y)

We want $E_{x,y,D}[(h(x,D) - y)^2]$ to be as low as possible.

But, typically, there is a bias-variance **tradeoff!**

Bias-variance tradeoff

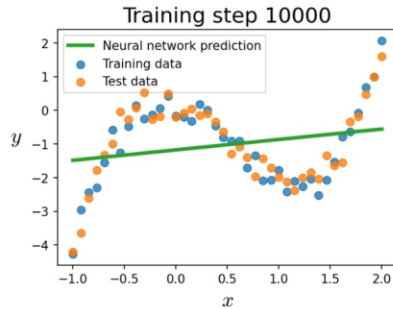
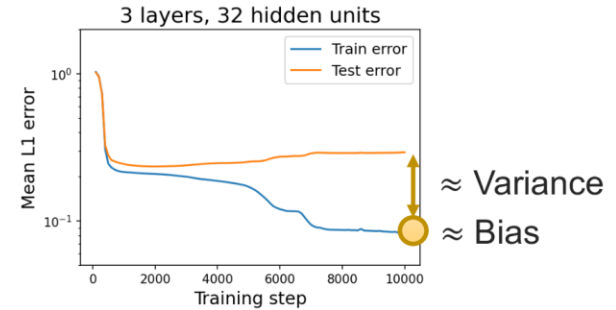
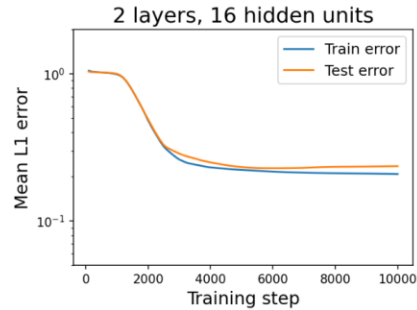
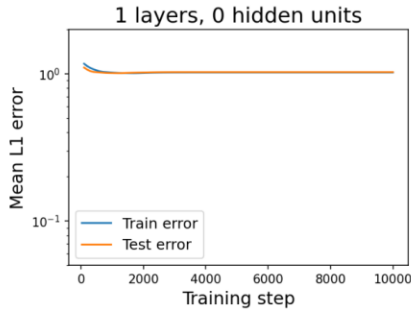


Source: <http://scott.fortmann-roe.com/docs/BiasVariance.html>

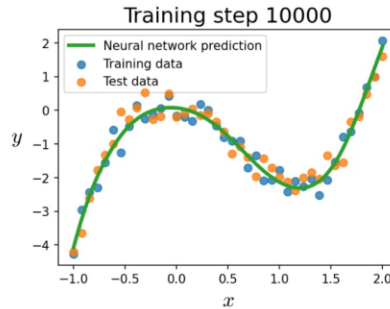
Bias = “average” model error

Variance = sensitivity of model prediction to training data

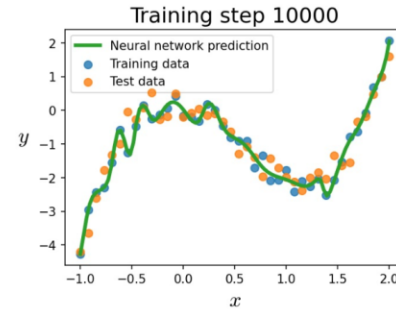
Overfitting vs underfitting



High bias
Low variance



Low bias
Low variance



Low bias
High variance

Improving model performance

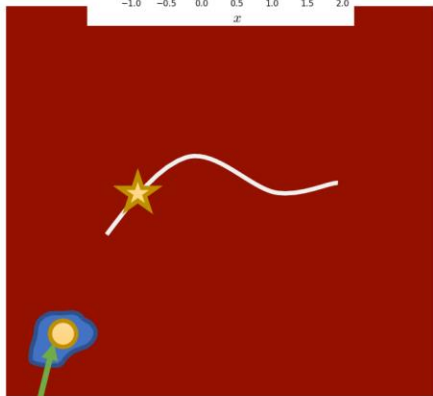
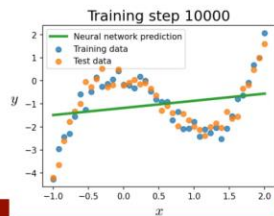
If you have high bias (or underfitting)

- Increase model **complexity** (e.g. number of free parameters)
- Or modify model architecture (shift location of hypothesis space)

If you have high variance (or overfitting)

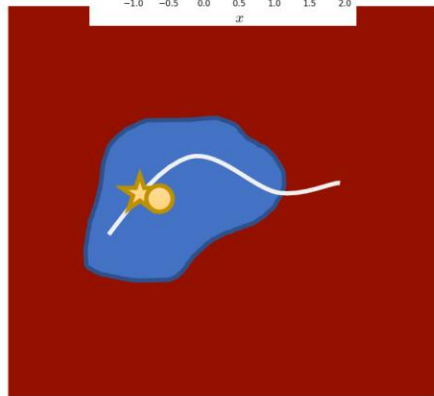
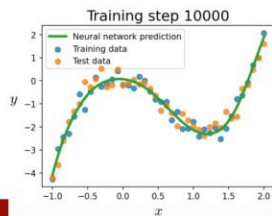
- Increase the amount of training **data**
- Or **regularise** the neural network in some other way

Function space

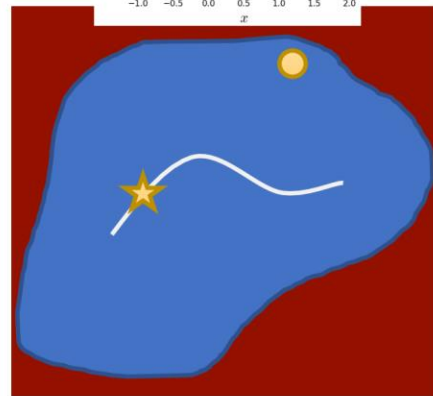
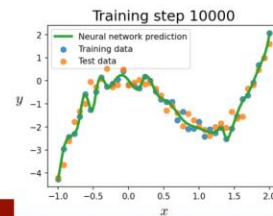


Underfit

Neural network, $NN(x; \theta)$, which minimises empirical loss $L(\theta)$



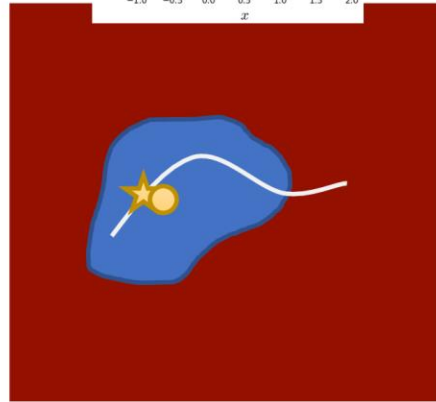
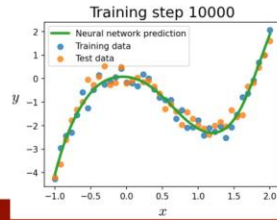
Ideal model



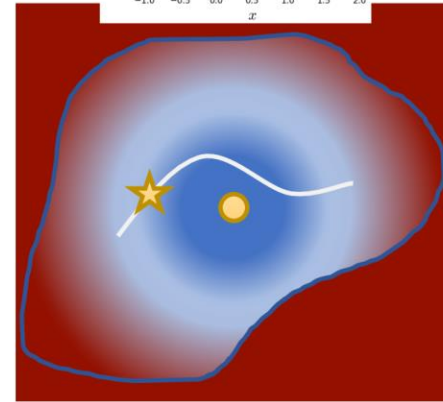
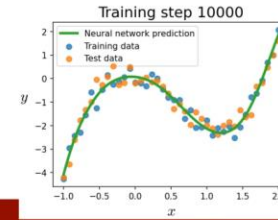
Overfit

Regularisation

- As well as changing the boundaries of the function space (hard constraint)
- We can **prefer** certain regions of the function space over others (soft constraint)



Ideal model



Regularised

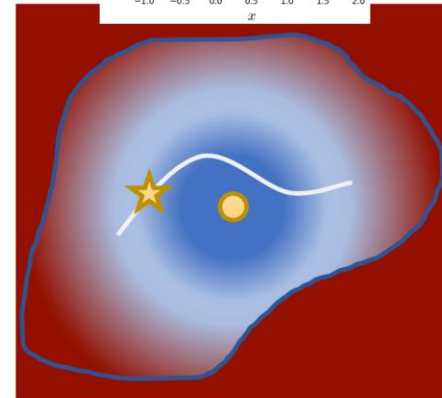
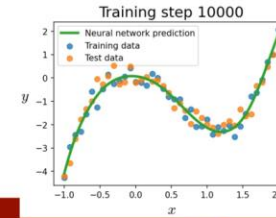
Regularisation



Regularisation = **restrict** the space of possible functions a neural network can learn / **prefer** certain regions of the function space / impart a **prior** on the learning algorithm

Ways to regularise neural networks:

- Reduce model complexity
- Modify model architecture
- Increase amount of (or augment) training data
- Weight regularization / pruning
- Additional loss terms
- Dropout
- Early stopping



Regularised

Additional loss terms

Weight regularisation:

$$L(\theta) = \frac{1}{N} \sum_i^N (NN(x_i; \theta) - y_i)^2 + \lambda \|\theta\|^2$$

Alternatively, we can use any **additional** arbitrary loss term as a regulariser:

$$L(\theta) = \frac{1}{N} \sum_i^N (NN(x_i; \theta) - y_i)^2 + R(D, \theta)$$

Examples include:

Structural similarity index measure (SSIM) loss (perceptual image quality measure)

$$\text{SSIM}(y, \hat{y}) = \frac{(2\mu_y\mu_{\hat{y}} + c_1)(2\sigma_{y\hat{y}} + c_2)}{(\mu_y^2 + \mu_{\hat{y}}^2 + c_1)(\sigma_y^2 + \sigma_{\hat{y}}^2 + c_2)}$$

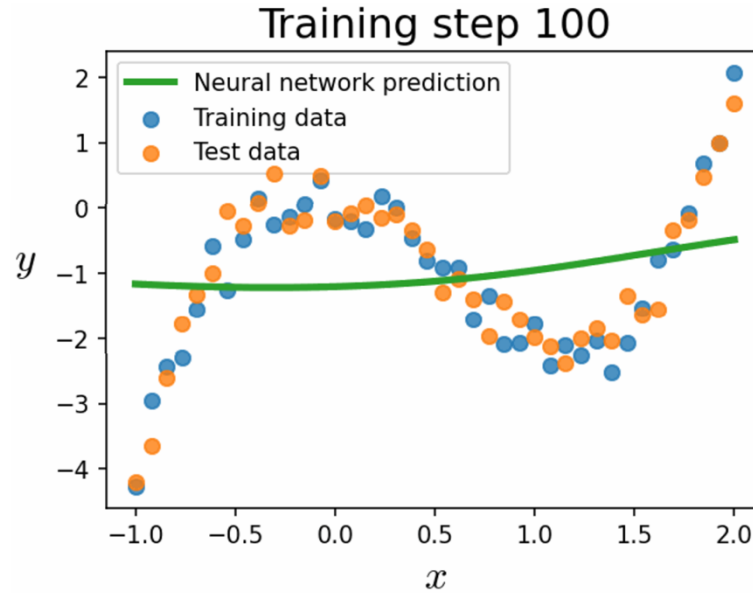
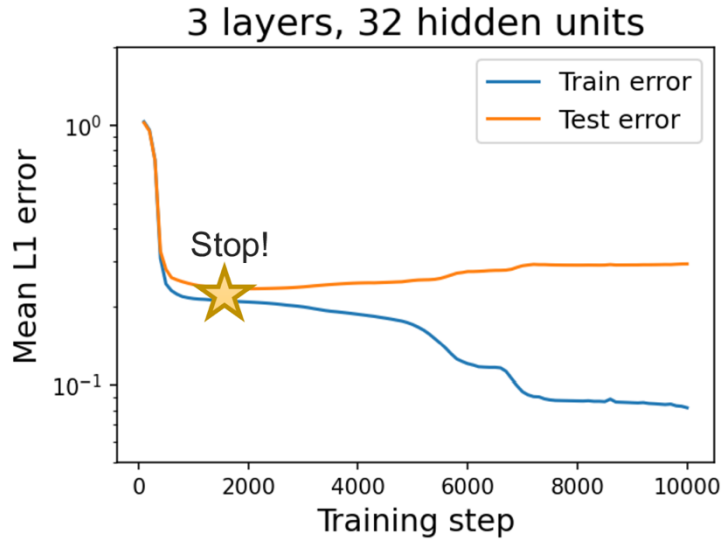
Smoothness loss (e.g. in frequency domain)

$$L(\hat{y}) = w(\omega) \text{FFT}(\hat{y})$$

Feature space loss (e.g. using a pretrained CNN)

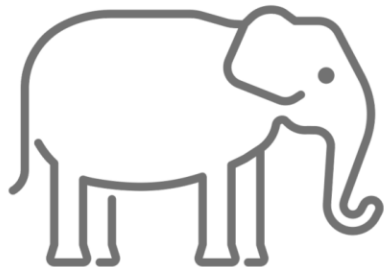
$$F(y, \hat{y}) = \|\text{CNN}(y)_{\text{hidden}} - \text{CNN}(\hat{y})_{\text{hidden}}\|^2$$

Early stopping



Related to the **spectral bias** of neural networks

Another source of error...



Up until now, we have assumed that the learned neural network parameters **globally** minimise the empirical loss, i.e.

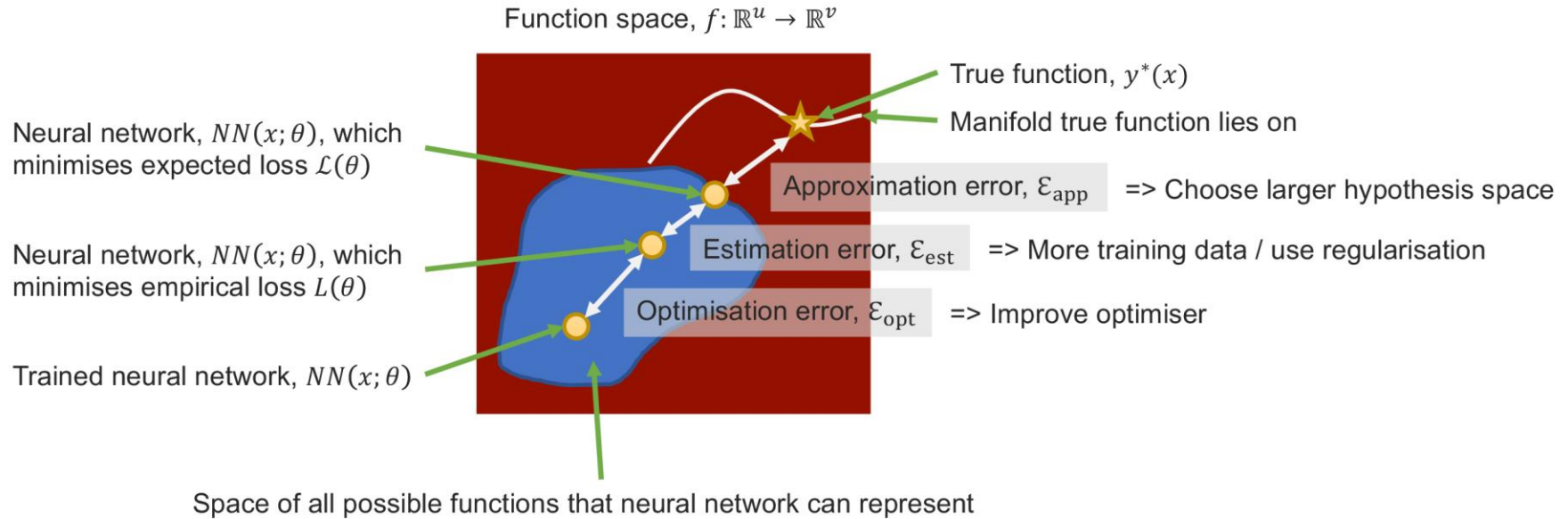
$$\theta^* = \min_{\theta} L(\theta)$$

⚠ But, in practice, finding the global minima is **extremely** hard

So, we must also consider the effect of **optimisation error** on the performance of neural networks

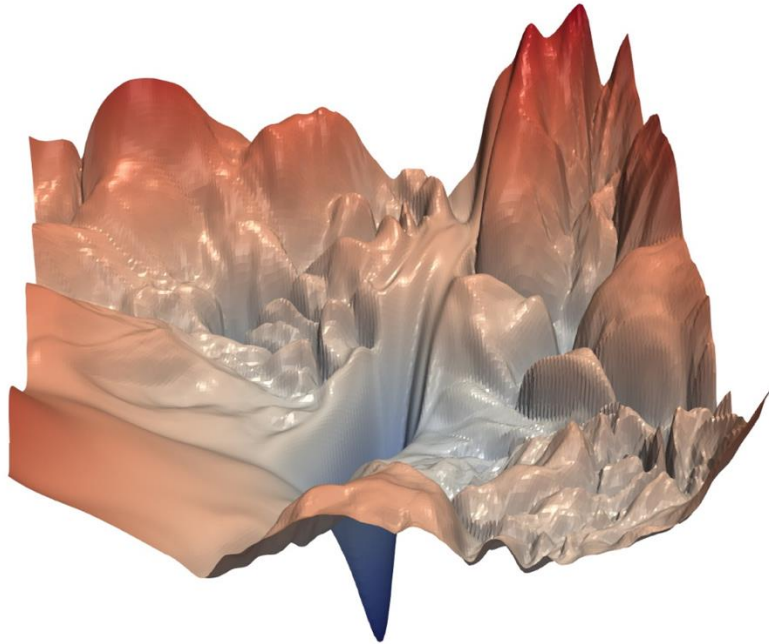
⚠ => Even though a large enough neural network can represent any arbitrary function (universal approximation theorem), there is **no guarantee** finding this network is tractable!

All error sources



$$\mathcal{E} = \epsilon_{\text{app}} + \epsilon_{\text{est}} + \epsilon_{\text{opt}}$$

Loss landscape



2D visualization of the loss function

$$l(\alpha, \beta) = L(\theta^* + \alpha u + \beta v)$$

Where u and v are 2 randomly sampled directional vectors in the parameter space

Loss surface for ResNet-56 without skip connections trained on CIFAR-10 (natural images)

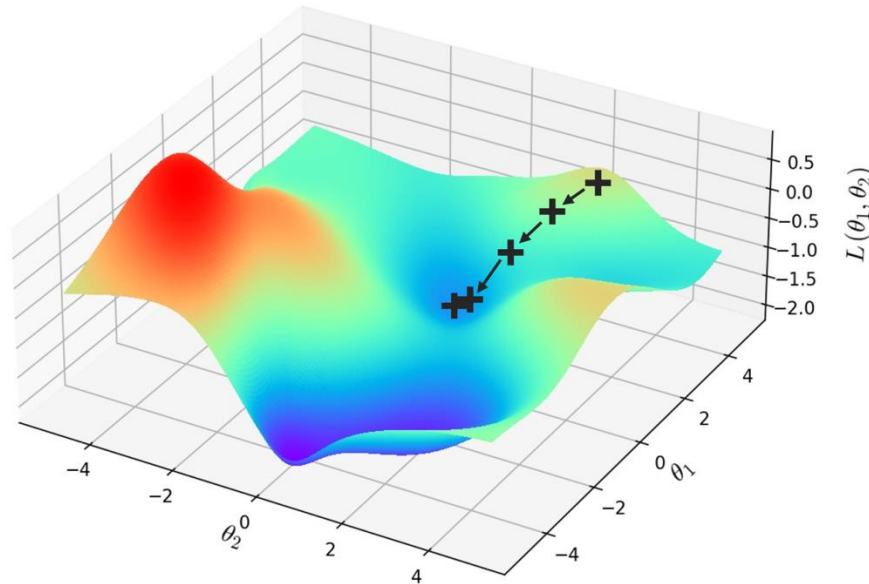


Loss function is very high dimensional and contains many **local minima**

Li et al, Visualizing the Loss Landscape of Neural Nets, NeurIPS (2018)

Gradient descent

- Problem 2): Loss function has lots of **local** minima



$$L(\theta) = \frac{1}{N} \sum_i^N (NN(x_i; \theta) - y_i)^2$$

1. Initialise weights randomly

2. Loop:

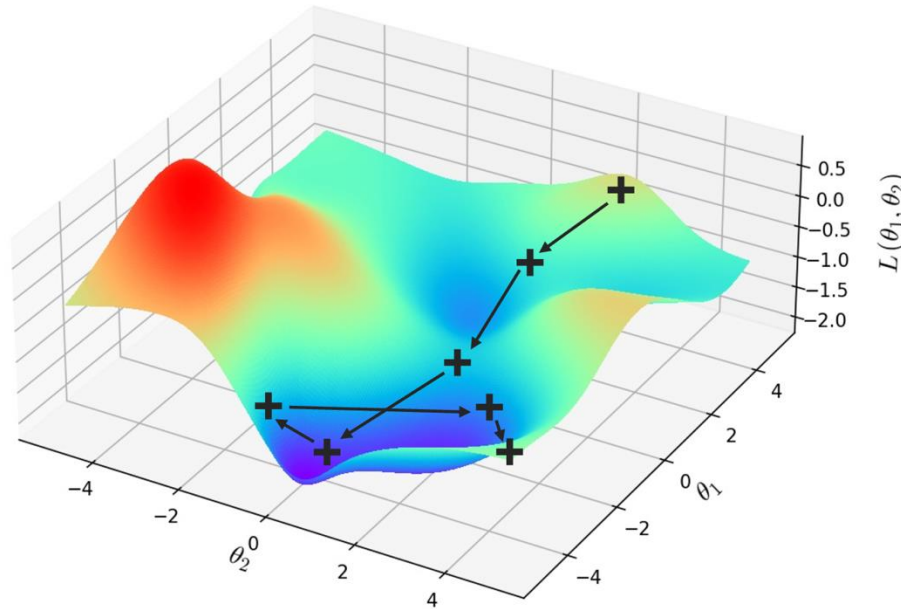
1. Compute gradient, $\frac{\partial L(\theta)}{\partial \theta_j}$

2. Update weights,

$$\theta_j \leftarrow \theta_j - \gamma \frac{\partial L(\theta)}{\partial \theta_j}$$

Gradient descent with large learning rate

- Can escape local minima (and global minima too!)
- Typically reduce learning rate with training step (e.g. exponentially decaying)



$$L(\theta) = \frac{1}{N} \sum_i^N (NN(x_i; \theta) - y_i)^2$$

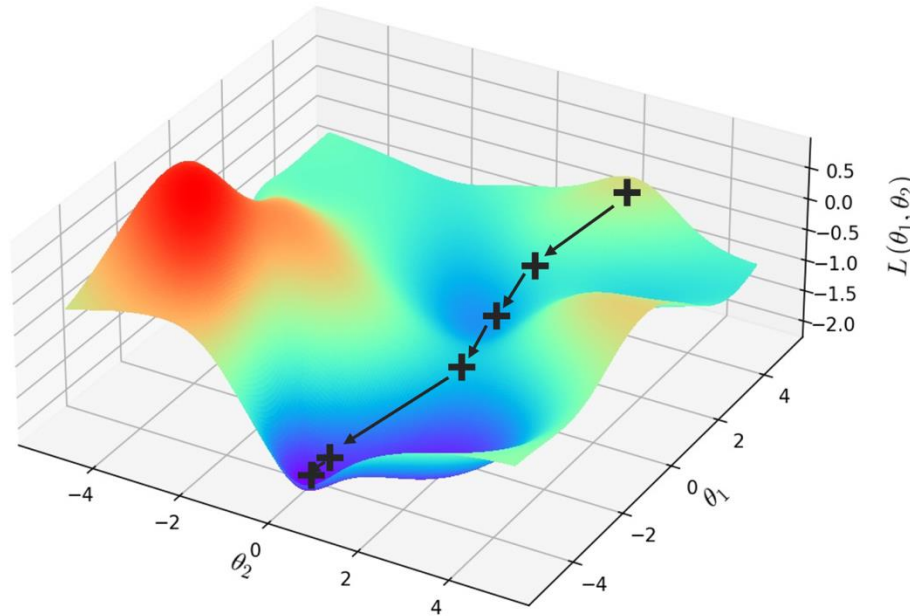
1. Initialise weights randomly
2. Loop:

1. Compute gradient, $\frac{\partial L(\theta)}{\partial \theta_j}$

2. Update weights,

$$\theta_j \leftarrow \theta_j - \gamma \frac{\partial L(\theta)}{\partial \theta_j}$$

Gradient descent with adaptive learning rate



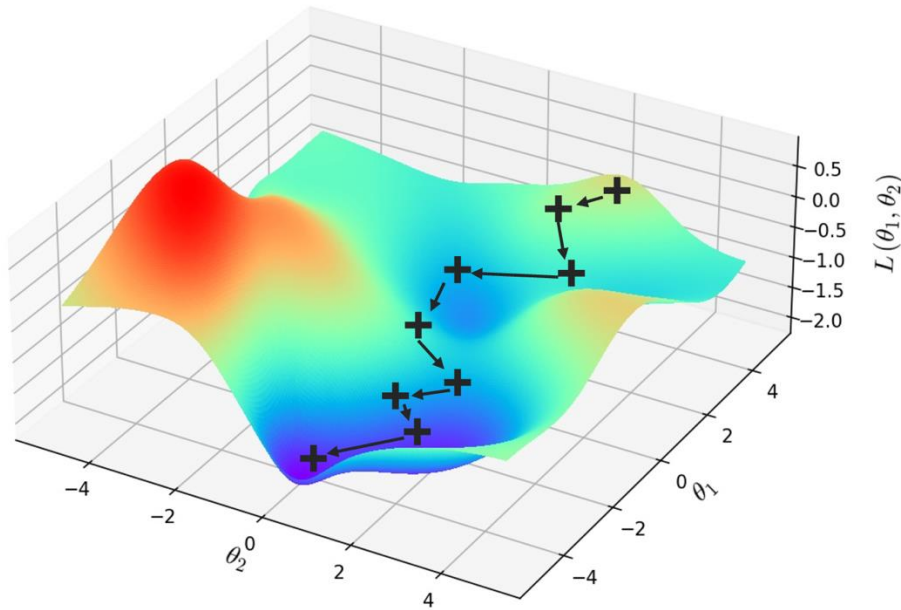
“Intelligent” learning rate

Can adapt learning rate based on

- How large gradient is
- Local curvature
- How fast learning is happening
- Values of weights
- .. many other ideas

(Mini-batch) stochastic gradient descent

$$L(\theta) = \frac{1}{N} \sum_i^N (NN(x_i; \theta) - y_i)^2$$



1. Initialise weights randomly

2. Loop:

1. **Randomly** sample mini-batch (subset) of B training points

2. Compute **mini-batch** gradient, $\frac{\partial L_B(\theta)}{\partial \theta_j}$

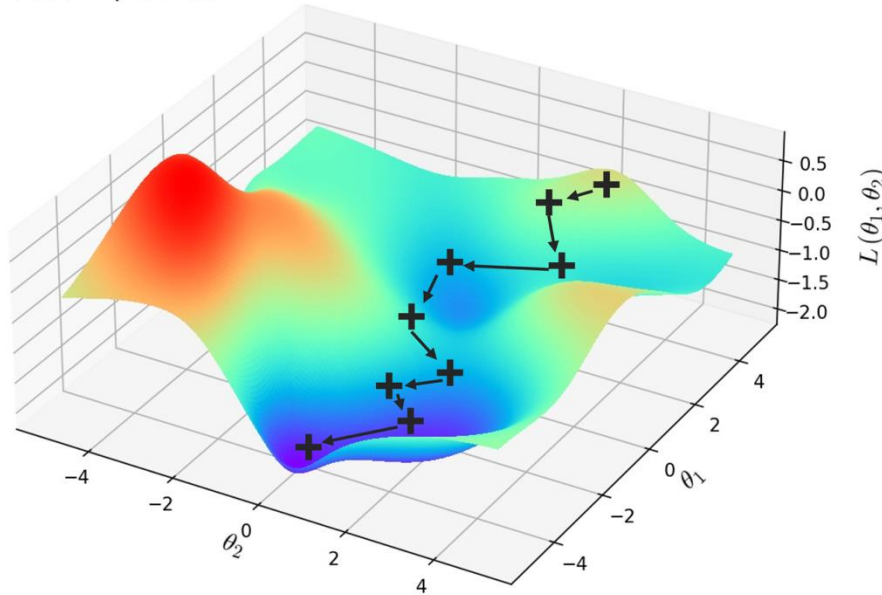
3. Update weights,

$$\theta_j \leftarrow \theta_j - \gamma \frac{\partial L_B(\theta)}{\partial \theta_j}$$

(Mini-batch) stochastic gradient descent



- Produces noisy (approximate) gradient estimates ✓
- Can escape local minima ✓
- Much **more** efficient gradient computation (batch size of 100s vs millions) ✓
- Widely used in practice



$$L(\theta) = \frac{1}{N} \sum_i^N (NN(x_i; \theta) - y_i)^2$$

1. Initialise weights randomly
2. Loop:
 1. **Randomly** sample mini-batch (subset) of B training points
 2. Compute **mini-batch** gradient, $\frac{\partial L_B(\theta)}{\partial \theta_j}$
 3. Update weights,

$$\theta_j \leftarrow \theta_j - \gamma \frac{\partial L_B(\theta)}{\partial \theta_j}$$

Adam optimiser

Adaptive moment estimation

Kingma et al, Adam: A Method for Stochastic Optimization, ICLR, (2015)

$$L(\theta) = \frac{1}{N} \sum_i^N (NN(x_i; \theta) - y_i)^2$$

1. Initialise θ randomly, and $m_0 \leftarrow 0$ (first moment) $v_0 \leftarrow 0$ (second moment)
2. For $t = 1$ to ... :

1. Compute gradient, $\frac{\partial L}{\partial \theta_{t-1}}$

2. Update moments, $m_t \leftarrow \beta_1 m_{t-1} + (1 - \beta_1) \frac{\partial L}{\partial \theta_{t-1}}$, $v_t \leftarrow \beta_2 v_{t-1} + (1 - \beta_2) \left(\frac{\partial L}{\partial \theta_{t-1}} \right)^2$

3. Apply bias correction, $\hat{m}_t \leftarrow \frac{m_t}{1 - \beta_1^t}$, $\hat{v}_t \leftarrow \frac{v_t}{1 - \beta_2^t}$

4. Update weights,

$$\theta_t \leftarrow \theta_{t-1} - \gamma \frac{\hat{m}_t}{\sqrt{\hat{v}_t + \epsilon}}$$

Typically, $\beta_1 = 0.9$, $\beta_2 = 0.999$

Adam optimiser

Idea 1: keep some **momentum** when moving in parameter space

$$m_t \leftarrow \beta_1 m_{t-1} + (1 - \beta_1) \frac{\partial L}{\partial \theta_{t-1}}$$

β_1 = “forgetting factor”

Essentially computes (exponentially-weighted) moving average of gradient

Idea 2: **adapt** the step size (per parameter) based on the **magnitude** of the gradient, again using some momentum

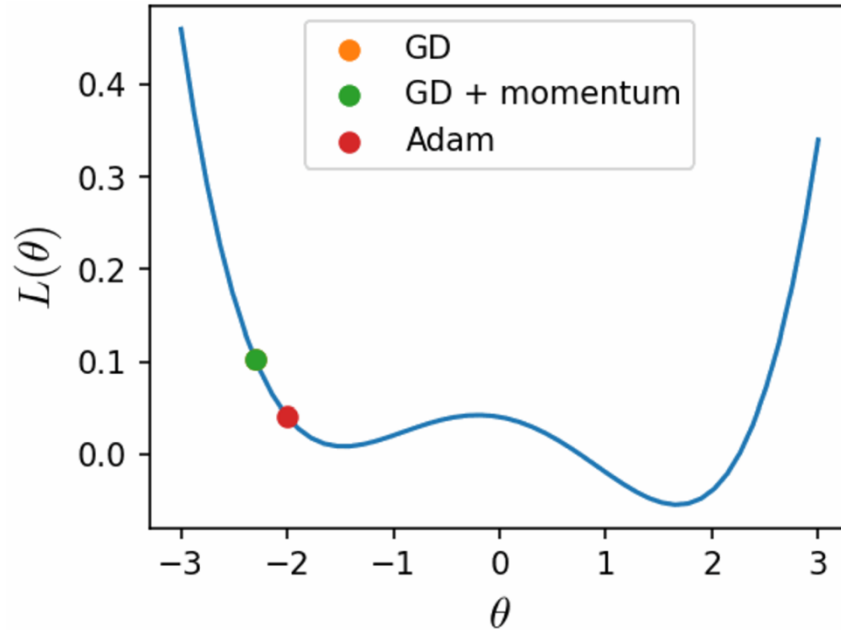
$$v_t \leftarrow \beta_2 v_{t-1} + (1 - \beta_2) \left(\frac{\partial L}{\partial \theta_{t-1}} \right)^2$$

“Normalises” step size

$$\theta_t \leftarrow \theta_{t-1} - \gamma \frac{\hat{m}_t}{\sqrt{\hat{v}_t + \epsilon}}$$

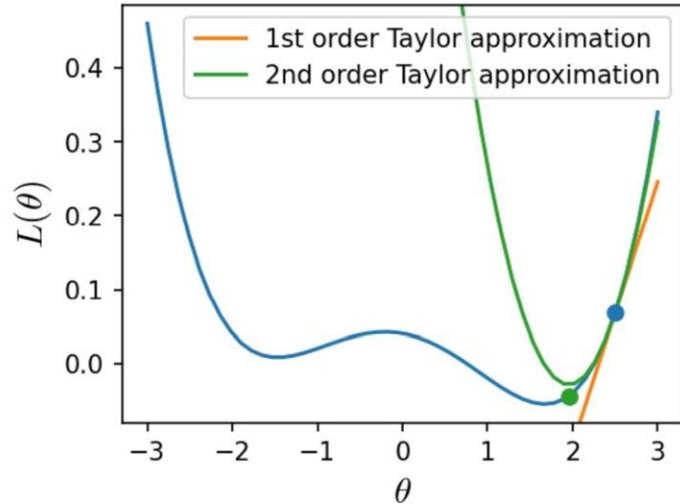
Note - the bias correction gives a more accurate estimate of the moving average (removes effect of initial value)

Adam optimiser Widely used in practice!



Higher-order optimisation

Make use of **higher-order** derivatives (e.g. curvature) of loss function to make bigger steps



Taylor approximation:

$$L(\theta) = L(a) + L'(a)(\theta - a) + \frac{1}{2}L''(a)(\theta - a)^2 + \dots$$

Newton's method:

Instead of following gradient, step towards the local **minima** of the 2nd order approximation

$$\theta \leftarrow \theta - \gamma \frac{L'(\theta)}{L''(\theta)}$$

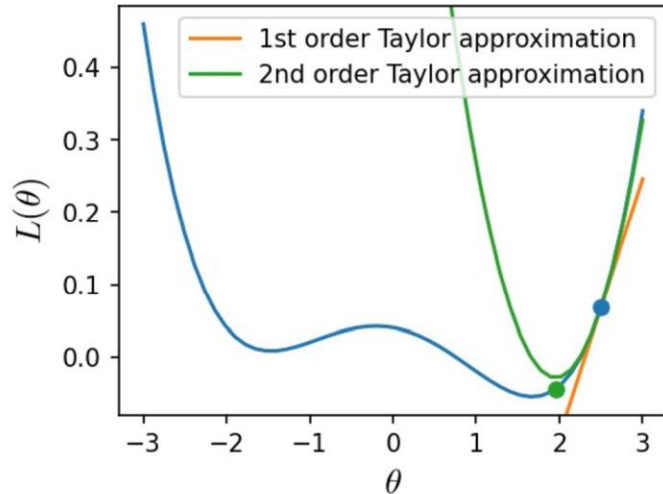
For a multivariate function:

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \gamma (H(L(\boldsymbol{\theta})))^{-1} \nabla L(\boldsymbol{\theta})$$

Where $H(L(\boldsymbol{\theta}))$ is the **Hessian** matrix

Higher-order optimisation

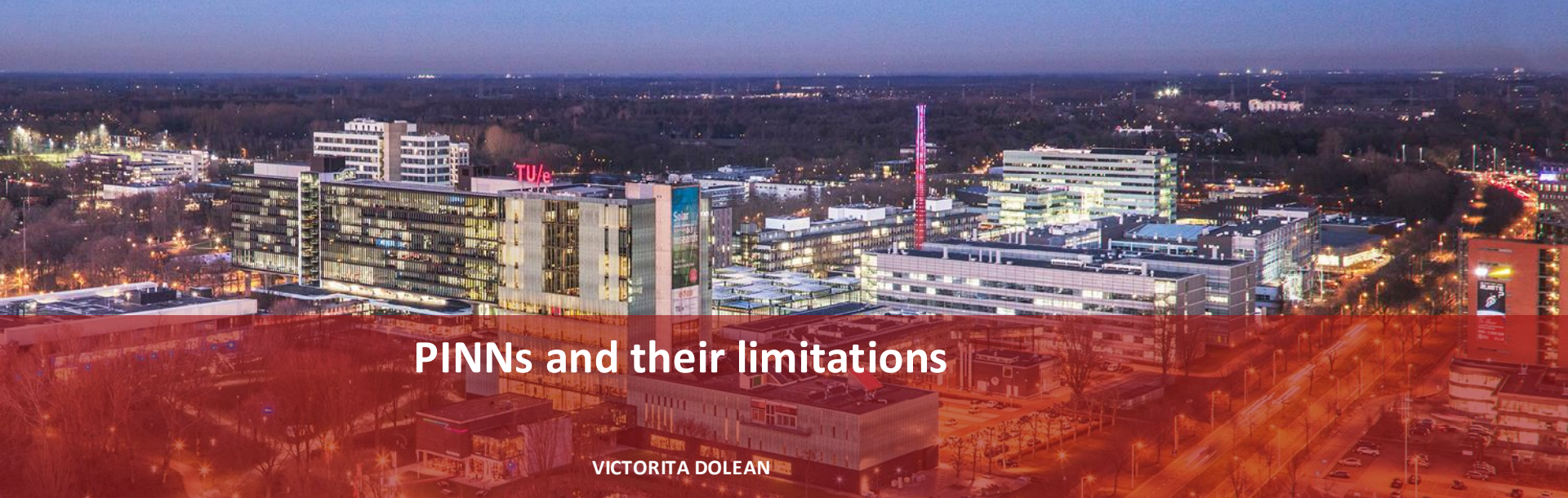
Make use of **higher-order** derivatives (e.g. curvature) of loss function to make bigger steps



$$\theta \leftarrow \theta - \gamma(H(L(\theta)))^{-1}\nabla L(\theta)$$

Note

- It is usually prohibitively expensive to compute the Hessian ($p \times p$) and its inverse
- So most modern methods (e.g. L-BFGS) use approximations
- These are called **quasi-Newton** methods



PINNs and their limitations

VICTORITA DOLEAN



What is a PINN?

Damped harmonic oscillator:

$$m \frac{d^2 u}{dt^2} + \mu \frac{du}{dt} + ku = 0$$

Initial conditions:

$$u(t = 0) = 1$$

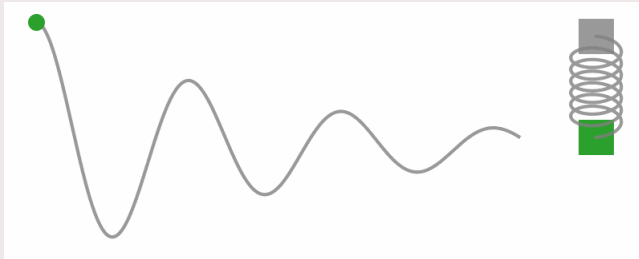
$$u_t(t = 0) = 0$$

u = displacement

m = mass of oscillator

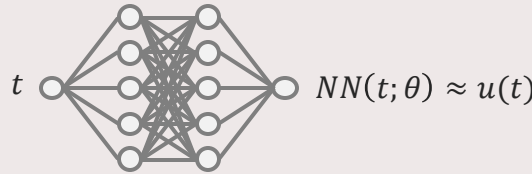
μ = coefficient of friction

k = spring constant



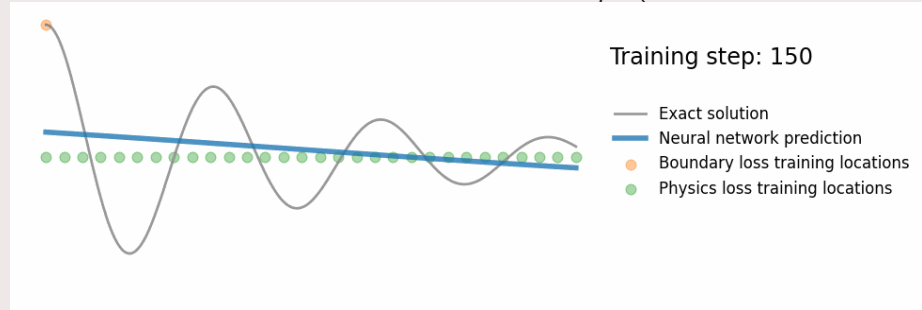
Key idea: use a neural network to directly approximate the solution

$$NN(t; \theta) \approx u(t)$$



Train the network using the loss function:

$$L(\theta) = \lambda_1 (NN(t = 0; \theta) - 1)^2 + \lambda_2 \left(\frac{dNN}{dt}(t = 0; \theta) - 0 \right)^2 + \frac{1}{N_p} \sum_i^{N_p} \left(\left[m \frac{d^2}{dt^2} + \mu \frac{d}{dt} + k \right] NN(t_i; \theta) \right)^2$$



Raissi et al, Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations, JCP (2018)

Lagaris et al, Artificial neural networks for solving ordinary and partial differential equations, IEEE (1998)

PINNs are a general framework for solving PDEs

Given a PDE and its boundary/initial conditions

$$\begin{aligned}\mathcal{D}[u(x)] &= f(x), & x \in \Omega \subset \mathbb{R}^d \\ \mathcal{B}_k[u(x)] &= g_k(x), & x \in \Gamma_k \subset \partial\Omega\end{aligned}$$

Where \mathcal{D} is some differential operator, \mathcal{B}_k are a set of boundary operators, and $u(x)$ is the solution to the PDE

PINNs train a neural network to **approximate** the solution to the PDE $NN(x; \theta) \approx u(x)$ using the following loss function:

$$L(\theta) = L_b(\theta) + L_p(\theta)$$

$$L_b(\theta) = \sum_k \frac{\lambda_k}{N_{bk}} \sum_j \|\mathcal{B}_k[NN(x_{kj}; \theta)] - g_k(x_{kj})\|^2$$

$$L_p(\theta) = \frac{1}{N_p} \sum_i \|\mathcal{D}[NN(x_i; \theta)] - f(x_i)\|^2$$

For example, the 1+1D viscous Burgers' equation:

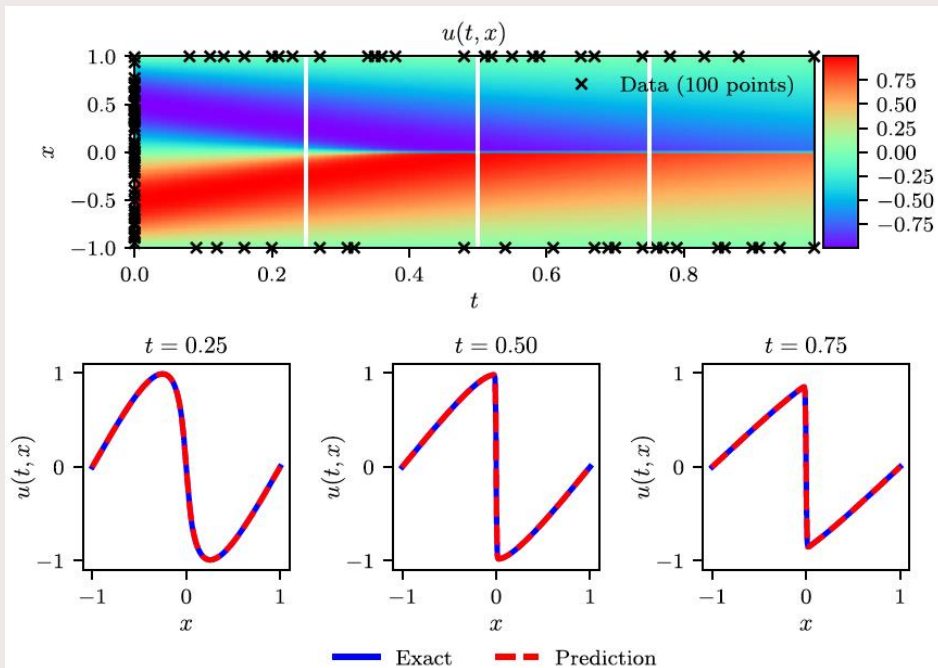
$$\frac{\partial u}{\partial t} + u \frac{\partial u}{\partial x} - \nu \frac{\partial^2 u}{\partial x^2} = 0$$

$$\begin{aligned}u(x, 0) &= -\sin(\pi x) \\ u(-1, t) &= u(+1, t) = 0\end{aligned}$$

$$u(x, t) \approx NN(x, t; \theta)$$

$$\begin{aligned}L_b(\theta) &= \frac{\lambda_1}{N_{b1}} \sum_j^{N_{b1}} (NN(x_j, 0; \theta) + \sin(\pi x_j))^2 \\ &+ \frac{\lambda_2}{N_{b2}} \sum_k^{N_{b2}} (NN(-1, t_k; \theta) - 0)^2 \\ &+ \frac{\lambda_3}{N_{b3}} \sum_l^{N_{b3}} (NN(+1, t_l; \theta) - 0)^2 \\ L_p(\theta) &= \frac{1}{N_p} \sum_i^{N_p} \left(\left(\frac{\partial NN}{\partial t} + NN \frac{\partial NN}{\partial x} - \nu \frac{\partial^2 NN}{\partial x^2} \right) (x_i, t_i; \theta) \right)^2\end{aligned}$$

PINNs for solving viscous Burgers' equation



$$L_b(\theta) = \frac{\lambda_1}{N_{b1}} \sum_j^{N_{b1}} (NN(x_j, 0; \theta) + \sin(\pi x_j))^2$$

$$+ \frac{\lambda_2}{N_{b2}} \sum_k^{N_{b2}} (NN(-1, t_k; \theta) - 0)^2$$

$$+ \frac{\lambda_3}{N_{b3}} \sum_l^{N_{b3}} (NN(+1, t_l; \theta) - 0)^2$$

$$L_p(\theta) = \frac{1}{N_p} \sum_i^{N_p} \left(\left(\frac{\partial NN}{\partial t} + NN \frac{\partial NN}{\partial x} - v \frac{\partial^2 NN}{\partial x^2} \right) (x_i, t_i; \theta) \right)^2$$

$$v = 0.01/\pi$$

$N_p = 10,000$ (Latin hypercube sampling)

$N_{b1} + N_{b2} + N_{b3} = 100$

Fully connected network with 9 layers, 20 hidden units
(3021 free parameters)

Tanh activation function

L-BFGS optimiser

PINNs – an entire research field

SPRINGER LINK

Find a journal | Publish with us | Track your research | Search

Home > Journal of Scientific Computing > Article

Scientific Machine Learning Through Physics–Informed Neural Networks: Where we are and What’s Next

Open access | Published: 26 July 2022
Volume 92, article number 88, (2022) | [Cite this article](#)

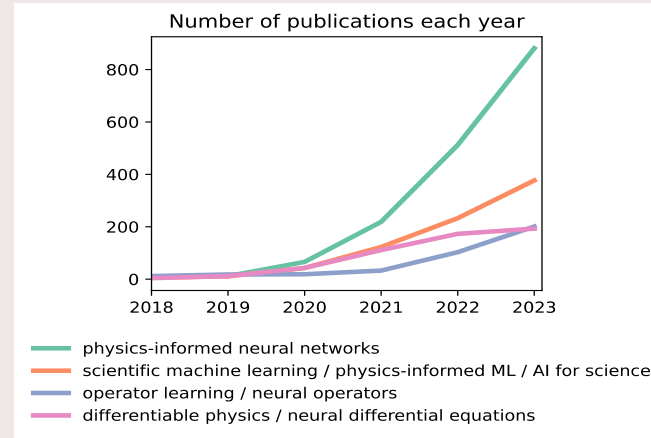
[Download PDF](#) | You have full access to this [open access](#) article

Salvatore Cuomo, Vincenzo Schiano Di Cola, Fabio Giampaolo, Gianluigi Rozza, Maziar Raissi & Francesco Piccialli

67k Accesses | 274 Citations | 7 Altmetric | [Explore all metrics](#)

Abstract

Physics-Informed Neural Networks (PINN) are neural networks (NNs) that encode model equations, like Partial Differential Equations (PDE), as a component of the neural network itself. PINNs are nowadays used to solve PDEs, fractional equations, integral-differential equations, and stochastic PDEs. This novel methodology has arisen as a multi-task learning framework in which a NN must fit observed data while reducing a PDE residual. This article provides a comprehensive review of the literature on PINNs: while the primary goal of the students to characterize these networks and their related advantages and



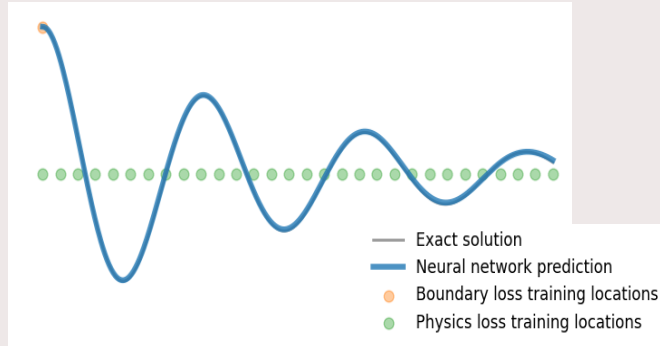
Source: Scopus keyword search (Feb 2024)

- The basic concepts behind PINNs were introduced in the 1990s (Lagaris et al, IEEE (1998) and others)
- Raissi et al, JCP (2019) is the seminal paper which reimplemented and extended PINNs using a modern deep learning framework

Raissi et al, Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations, JCP (2018)

Lagaris et al, Artificial neural networks for solving ordinary and partial differential equations, IEEE (1998)

PINN training loop



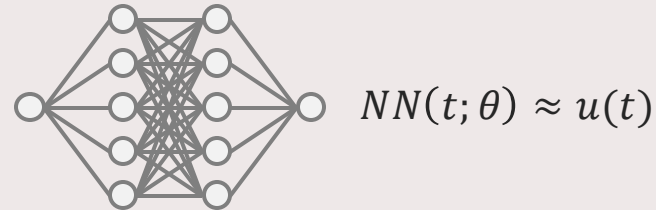
Training loop:

1. Sample boundary/ physics training points
2. Compute network outputs
3. Compute 1st and 2nd order gradient of network output **with respect to network input**
4. Compute loss
5. Compute gradient of loss function **with respect to network parameters**
6. Take gradient descent step

$$\text{Boundary loss } L_b(\theta) \left\{ \begin{array}{l} L(\theta) = \lambda_1 (NN(t=0; \theta) - 1)^2 \\ + \lambda_2 \left(\frac{dNN}{dt}(t=0; \theta) - 0 \right)^2 \end{array} \right.$$

$$\text{Physics loss } L_p(\theta) \left\{ + \frac{1}{N_p} \sum_i^{N_p} \left(\left[m \frac{d^2}{dt^2} + \mu \frac{d}{dt} + k \right] NN(t_i; \theta) \right)^2 \right.$$

How can we compute the gradients (e.g. $\frac{dNN}{dt}$ and $\frac{dL}{d\theta}$) required in steps 3 and 5?



Gradient computation for PINNs

Assume network is an MLP, e.g.:

$$NN(t; \theta) = W_3(\sigma(W_2\sigma(W_1t + \mathbf{b}_1) + \mathbf{b}_2) + b_3 \approx u(t)$$

Analogous to backpropagation, we can use the multivariate **chain rule** to compute gradients with respect to network **inputs**:

$$\begin{aligned} NN(t; \theta) &= f \circ \mathbf{g} \circ \mathbf{h}(t; \theta) \\ f &= W_3\sigma(\mathbf{g}) + b_3 \\ \mathbf{g} &= W_2\sigma(\mathbf{h}) + \mathbf{b}_2 \\ \mathbf{h} &= W_1t + \mathbf{b}_1 \end{aligned}$$

$$\frac{\partial NN}{\partial t} = \frac{\partial f}{\partial \mathbf{g}} \frac{\partial \mathbf{g}}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial t}$$

$$\frac{\partial NN}{\partial t} = W_3 \text{diag}(\sigma'(\mathbf{g}))W_2 \text{diag}(\sigma'(\mathbf{h}))W_1$$

Extending computational graph



Computing gradients can simply be thought of as **extending** the network's computational graph:

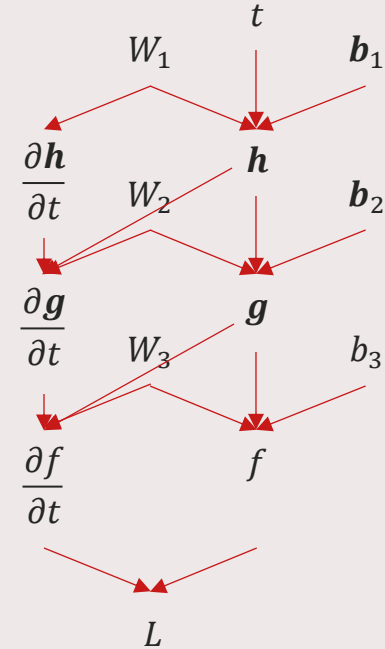
$$NN(t; \theta) = W_3(\sigma(W_2\sigma(W_1t + \mathbf{b}_1) + \mathbf{b}_2) + \mathbf{b}_3)$$

$$\begin{aligned} f &= W_3\sigma(\mathbf{g}) + \mathbf{b}_3 \\ \mathbf{g} &= W_2\sigma(\mathbf{h}) + \mathbf{b}_2 \\ \mathbf{h} &= W_1t + \mathbf{b}_1 \end{aligned}$$

$$\frac{\partial NN}{\partial t} = W_3 \text{diag}(\sigma'(\mathbf{g}))W_2 \text{diag}(\sigma'(\mathbf{h}))W_1$$



We can recursively apply **autodifferentiation** to compute gradients and extend the graph



PINN training loop

Training loop:

1. Sample boundary/ physics training points
2. Compute network outputs
3. Compute 1st and 2nd order gradient of network output **with respect to network input** **<= (recursively) apply autodiff, extending graph**
4. Compute loss
5. Compute gradient of loss function **with respect to network parameters** **<= apply autodiff on extended graph**
6. Take gradient descent step



We can recursively apply **autodifferentiation** to compute gradients and extend the graph

```
# PINN training psuedocode

#2.
t.requires_grad_(True)# tells PyTorch to start tracking graph
theta.requires_grad_(True)
u = NN(t, theta)

#3.
dudt = torch.autograd.grad(u, t, torch.ones_like(u),
                            create_graph=True)[0]
d2udt2 = torch.autograd.grad(dudt, t, torch.ones_like(u),
                              create_graph=True)[0]

#4.
physics_loss = torch.mean((m*d2udt2 + mu*dudt + k*u)**2)
loss = physics_loss + lambda_*boundary_loss

#5.
dtheta = torch.autograd.grad(loss, theta)[0]
```

Computational cost of higher order derivatives



Note, gradient computation roughly **doubles** the size of the computational graph*

⇒ The cost of evaluating $\frac{\partial^n NN}{\partial t^n}$ grows **exponentially** with n (!)

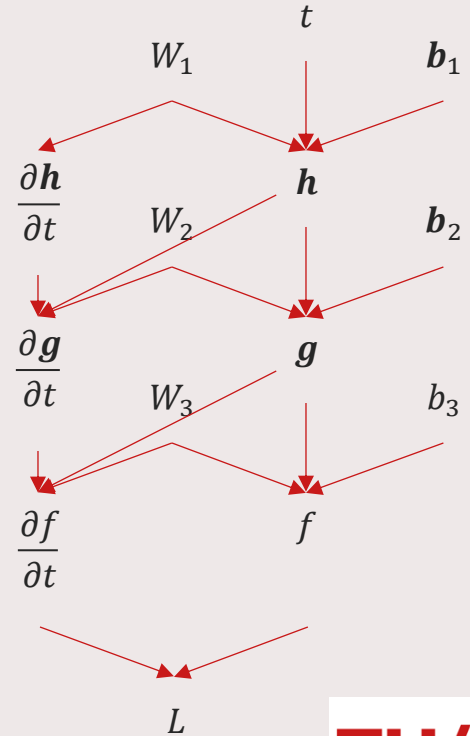
⇒ Most time is spent computing gradients, not the forward pass, when training PINNs

*More precisely, given some $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$, its Jacobian $J \in \mathbb{R}^{m \times n}$ and some vector $v \in \mathbb{R}^n$

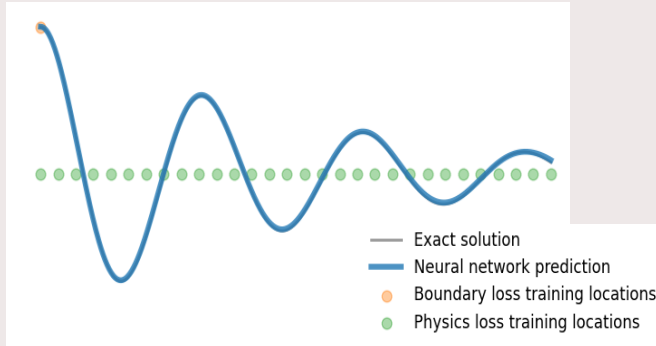
$$TIME\{f, Jv\} \leq \omega TIME\{f\}$$

With a constant $\omega \in [2, 2.5]$ using autodifferentiation

For detailed derivation, see eg: Griewank and Walther, Evaluating Derivatives, Ch 3.1, SIAM (2008))



Other important considerations

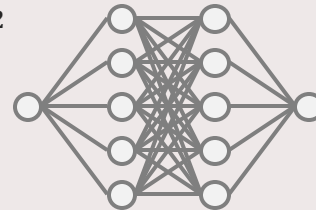


Boundary loss
 $L_b(\theta)$

$$\left\{ \begin{aligned} L(\theta) &= \lambda_1 (NN(t=0; \theta) - 1)^2 \\ &+ \lambda_2 \left(\frac{dNN}{dt}(t=0; \theta) - 0 \right)^2 \end{aligned} \right.$$

Physics loss
 $L_p(\theta)$

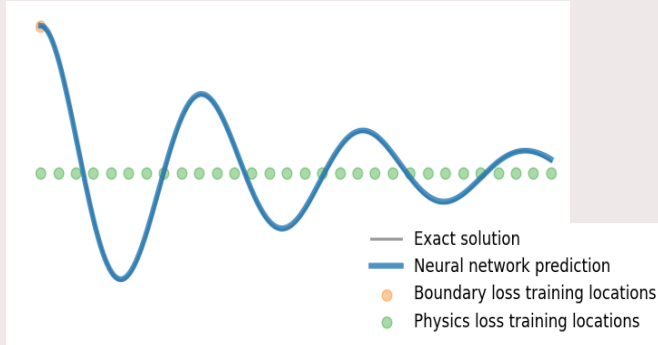
$$\left\{ + \frac{1}{N_p} \sum_i^{N_p} \left(\left[m \frac{d^2}{dt^2} + \mu \frac{d}{dt} + k \right] NN(t_i; \theta) \right)^2 \right.$$



$$NN(t; \theta) \approx u(t)$$

- λ values often significantly affect convergence
- “Enough” collocation points $\{t_i\}_{i=1}^{N_p}$ must be chosen so that the learned solution is accurate across the full domain
- Similarly, “enough” boundary points must be chosen such that the learned solution is unique
- Training points are usually sampled uniformly, randomly or quasi randomly
- PINNs still suffer from approximation, estimation and optimisation error (just like normal neural networks).

Other important considerations



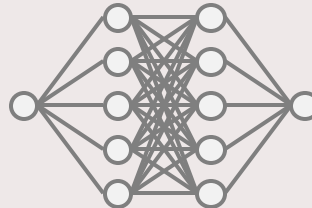
From a mathematical perspective:

- PINNs provide a way to solve PDEs:
 - Neural network is a mesh-free, **functional** approximation of PDE solution
 - Physics loss is used to assert solution is **consistent** with PDE
 - Boundary loss is used to assert boundary/initial conditions, to ensure solution is **unique**

$$\text{Boundary loss } L_b(\theta) \left\{ \begin{aligned} &L(\theta) = \lambda_1 (NN(t=0; \theta) - 1)^2 \\ &+ \lambda_2 \left(\frac{dNN}{dt}(t=0; \theta) - 0 \right)^2 \end{aligned} \right.$$

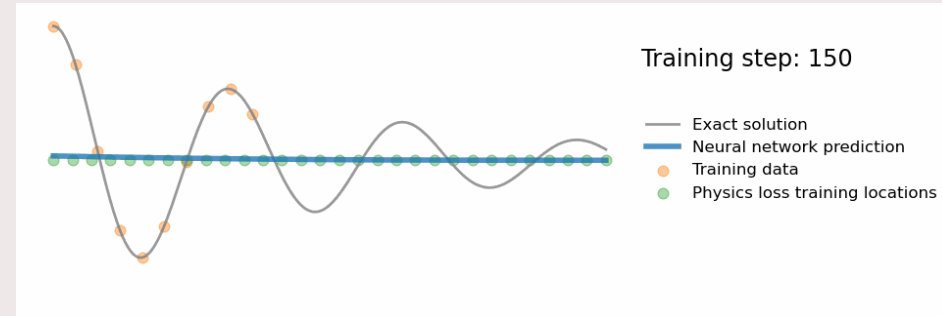
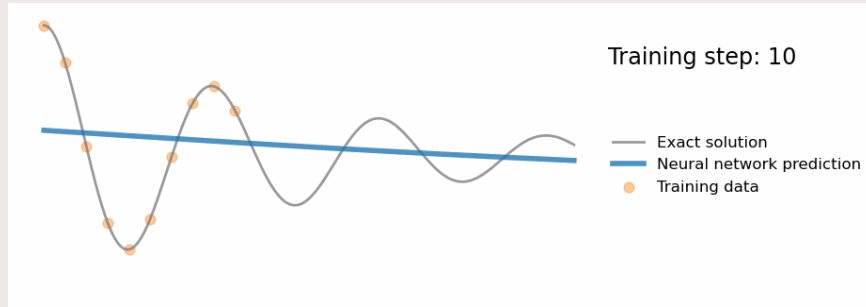
What can we say about PINNs from a ML perspective?
 Physics loss is an **unsupervised** regulariser, which adds prior knowledge

$$\text{Physics loss } L_p(\theta) \left\{ + \frac{1}{N_p} \sum_i^{N_p} \left(\left[m \frac{d^2}{dt^2} + \mu \frac{d}{dt} + k \right] NN(t_i; \theta) \right)^2 \right.$$



$$NN(t; \theta) \approx u(t)$$

PINNs from a ML perspective



$$L(\theta) = \frac{1}{N} \sum_i^N (NN(t_i; \theta) - \underline{u_i})^2$$

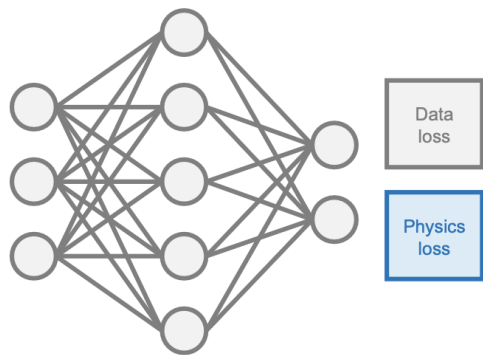
$$L(\theta) = \frac{1}{N} \sum_i^N (NN(t_i; \theta) - u_i)^2 \quad \text{Supervised loss}$$
$$+ \frac{\lambda}{M} \sum_j^M \left(\left[m \frac{d^2}{dt^2} + \mu \frac{d}{dt} + k \right] NN(t_j; \theta) \right)^2 \quad \text{Physics loss}$$

From a ML perspective:

- Physics loss is an **unsupervised** regulariser, which adds prior knowledge

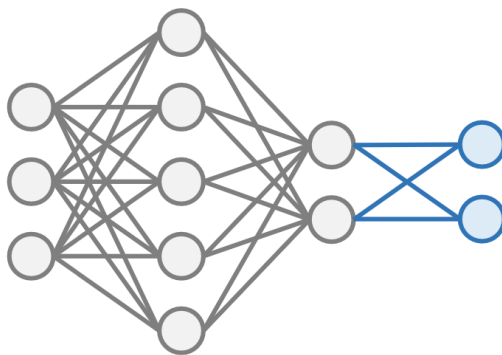
Ways to incorporate scientific principles into machine learning

Loss function



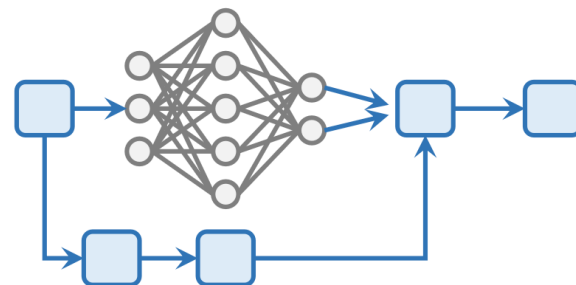
Example:
Physics-informed neural networks
(add governing equations to loss
function)

Architecture



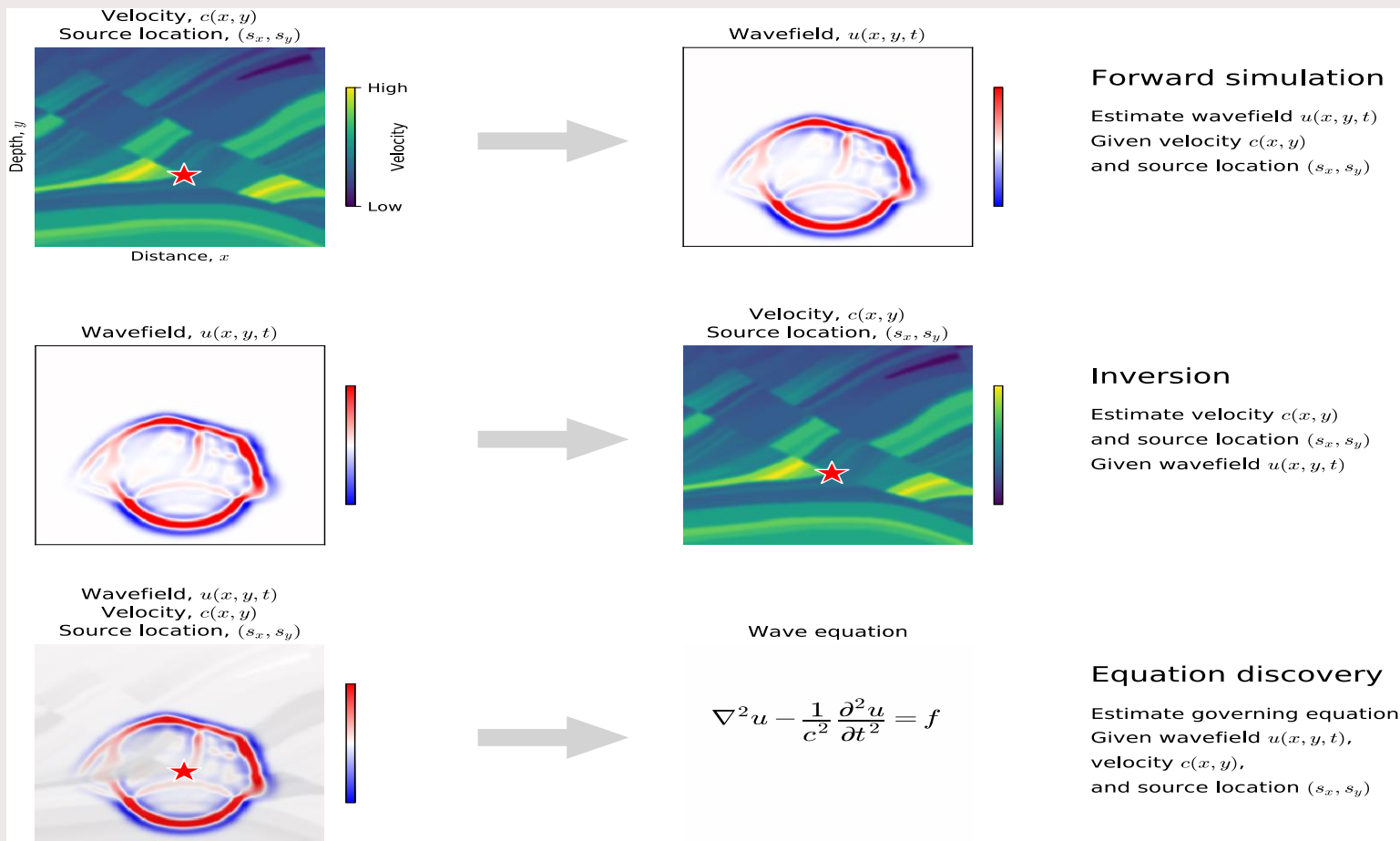
Example:
Encoding symmetries / conservation laws
(e.g. energy conservation, rotational
invariance)

Hybrid approaches



Example:
Neural differential equations
(incorporating neural networks into PDE
models)

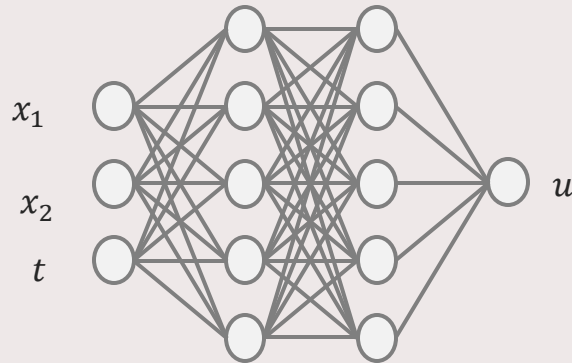
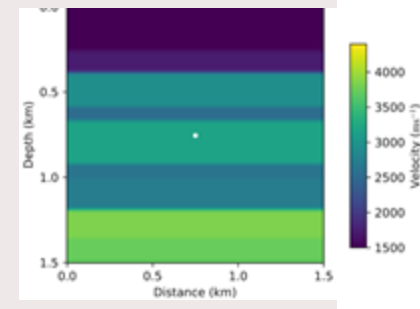
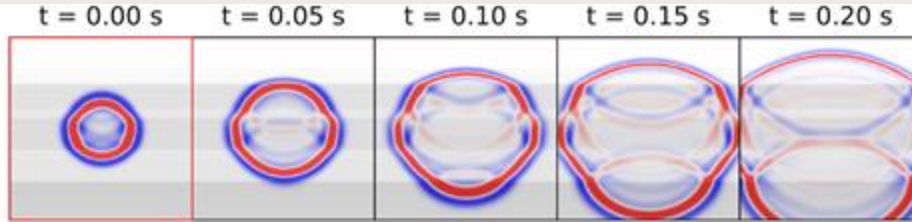
Key scientific tasks



1. PINNs for solving wave equation – forward problem

Velocity model, $c(x)$

Ground truth FD simulation



Moseley et al, Solving the wave equation with physics-informed deep learning, ArXiv (2020)

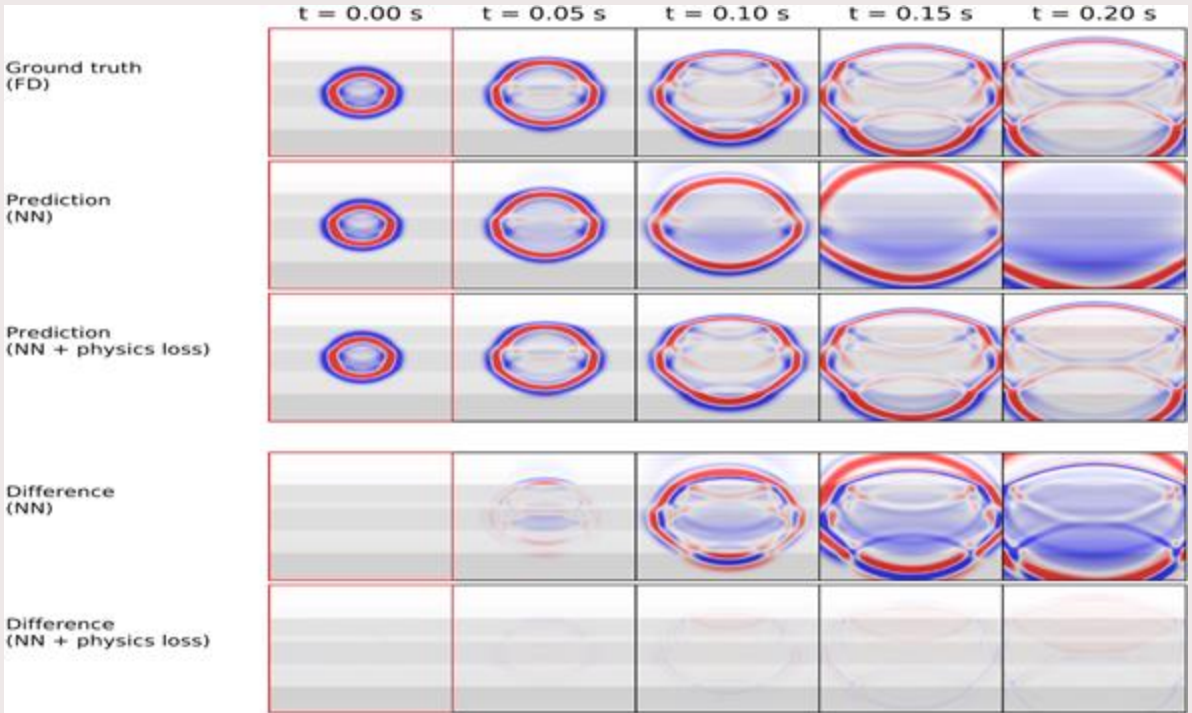
Boundary data from FD simulation (first 0.02 seconds)

Collocation points randomly sampled over entire domain (up to 0.2 seconds)

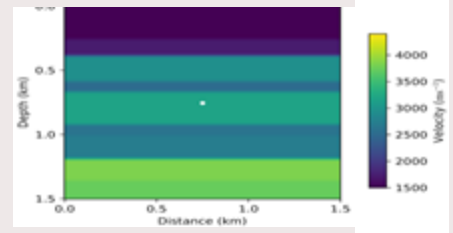
$$L_b(\theta) = \frac{\lambda}{N_b} \sum_j^{N_b} \left(NN(x_j, t_j; \theta) - u_{FD}(x_j, t_j) \right)^2$$

$$L_p(\theta) = \frac{1}{N_p} \sum_i^{N_p} \left(\left[\nabla^2 - \frac{1}{c(x_i)^2} \frac{\partial^2}{\partial t^2} \right] NN(x_i, t_i; \theta) \right)^2$$

PINNs for solving wave equation



Velocity model, $c(x)$



Moseley et al, Solving the wave equation with physics-informed deep learning, ArXiv (2020)

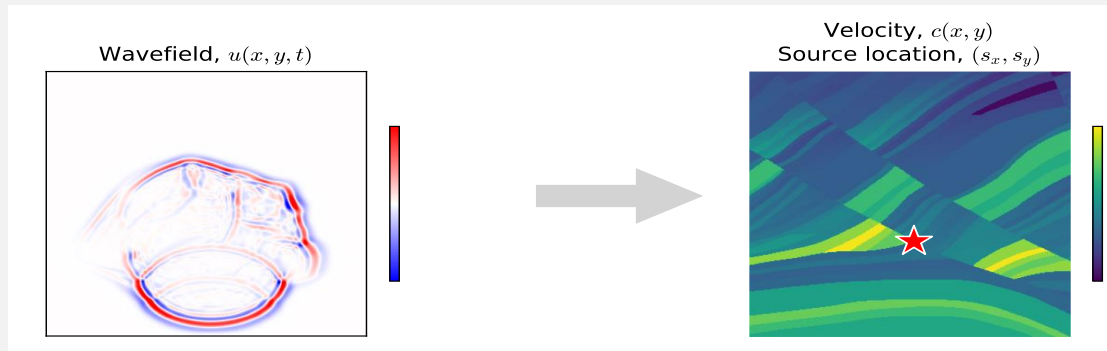
- Mini-batch size $N_b = N_p = 500$ (random sampling)
- Fully connected network with 10 layers, 1024 hidden units
- Softplus activation
- Adam optimiser

Training time: **~1 hour**

What is an inverse problem?

Wave equation:

$$\nabla^2 u - \frac{1}{c(x)^2} \frac{\partial^2 u}{\partial t^2} = 0$$



- Fundamentally, inverse problems are **search** problems
- It is often useful to frame them as an optimisation problem, for example:

$$\min_{\hat{a}} \|b - F(\hat{a})\|^2$$

b = observed
wavefield $u(x, t)$

$$b = F(a)$$

a = (velocity model
 $c(x)$, source
location)

a = set of input conditions

F = physical model of the system

b = resulting properties given F and a

2. PINNs for inversion

PINNs for solving **forward** simulation:

$$L(\theta) = L_b(\theta) + L_p(\theta)$$

$$L_b(\theta) = \sum_k \frac{\lambda_k}{N_{bk}} \sum_j^{N_{bk}} \left\| \mathcal{B}_k[NN(x_{kj}; \theta)] - g_k(x_{kj}) \right\|^2$$

Boundary loss

$$L_p(\theta) = \frac{1}{N_p} \sum_i^{N_p} \left\| \mathcal{D}[NN(x_i; \theta)] - f(x_i) \right\|^2$$

Physics loss

For example:

$$D = \left[\nabla^2 - \frac{1}{c(x)^2} \frac{\partial^2}{\partial t^2} \right]$$
$$f = 0$$

PINNs for solving **inverse** problems:

$$L(\theta, \phi) = L_p(\theta, \phi) + L_d(\theta)$$

$$L_p(\theta, \phi) = \frac{1}{N_p} \sum_i^{N_p} \left\| \mathcal{D}[NN(x_i; \theta); \phi] - f(x_i) \right\|^2$$

Physics loss

$$L_d(\theta) = \frac{\lambda}{N_d} \sum_l^{N_d} \left\| NN(x_l; \theta) - u_l \right\|^2$$

Data loss

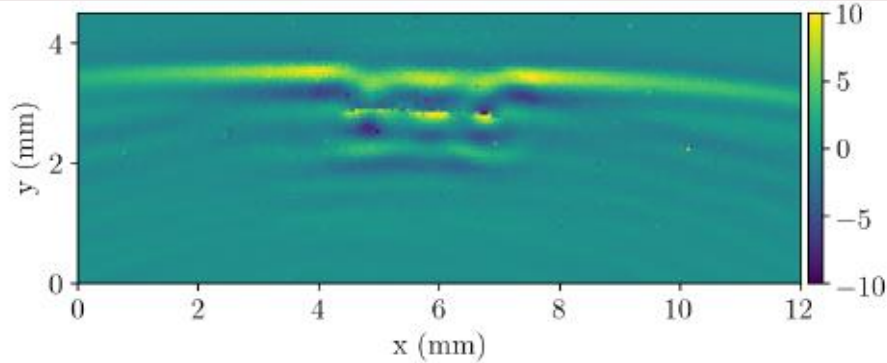
where ϕ are unknown, underlying PDE parameters we wish to invert for, and $\{x_l, u_l\}$ are a set of (potentially noisy) observational data



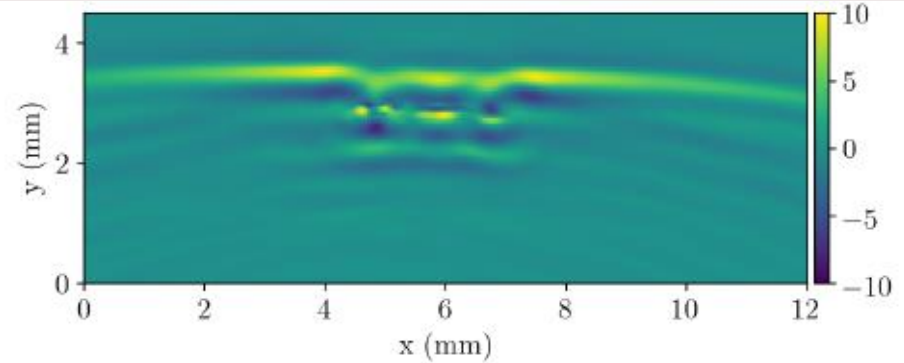
We **simultaneously** learn θ and ϕ when training the PINN

Wave equation

Shukla et al, Physics-Informed Neural Network for Ultrasound
 Nondestructive Quantification of Surface Breaking Cracks, Journal of
 Nondestructive Evaluation (2020)



(a) Actual data at $t = 12.38 \mu\text{s}$.

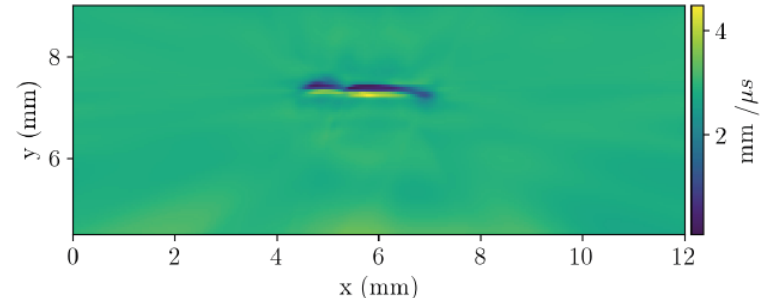


(b) Data recovered from PINN simulation at $t = 12.38 \mu\text{s}$.

$$L_p(\theta, \phi) = \frac{1}{N_p} \sum_i^{N_p} \left(\left[\nabla^2 - \frac{1}{c(x_i; \phi)^2} \frac{\partial^2}{\partial t^2} \right] NN(x_i, t_i; \theta) \right)^2$$

$$L_d(\theta) = \frac{\lambda}{N_d} \sum_l^{N_d} (NN(x_l, t_l; \theta) - u_{obs l})^2$$

Treat velocity model as **another** neural network, and simultaneously learn it



(d) Speed $v(x, y)$ recovered from PINN simulation.

3. PINNs for equation discovery

PINNs for solving **inverse** problems:

$$L(\theta, \phi) = L_p(\theta, \phi) + L_d(\theta)$$

$$L_p(\theta, \phi) = \frac{1}{N_p} \sum_i^{N_p} \|\mathcal{D}[NN(x_i; \theta); \phi] - f(x_i)\|^2 \quad \text{Physics loss}$$

$$L_d(\theta) = \frac{\lambda}{N_d} \sum_l^{N_d} \|NN(x_l; \theta) - u_l\|^2 \quad \text{Data loss}$$

Where $\{x_l, u_l\}$ are a set of (potentially noisy) observational data

But how do we learn the **entire** differential operator \mathcal{D} , rather than its parameters ϕ ?

How do we learn an **entire** differential operator \mathcal{D} ?

Build a **library** of n operators, such as:

$$\phi = (1, \partial_x, \partial_t, \partial_{xx}, \partial_{tt}, \partial_{xt})^T$$

Then assume the differential operator can be represented as

$$\mathcal{D} = \Lambda \phi$$

Where Λ is a (sparse) matrix of shape (d_u, n)

E.g. for 1D damped harmonic oscillator:

$$\mathcal{D} = (k \quad \mu \quad m \quad 0) \begin{pmatrix} 1 \\ d_t \\ d_{tt} \\ d_{ttt} \end{pmatrix} \\ = m \frac{d^2}{dt^2} + \mu \frac{d}{dt} + k$$

3. PINNs for equation discovery

How do we learn an **entire** differential operator \mathcal{D} ? PINNs for **equation discovery**:

Build a **library** of n operators, such as:

$$\phi = (1, \partial_x, \partial_t, \partial_{xx}, \partial_{tt}, \partial_{xt})^T$$

Then assume the differential operator can be represented as

$$\mathcal{D} = \Lambda \phi$$

Where Λ is a (sparse) matrix of shape (d_u, n)

E.g. for 1D damped harmonic oscillator:

$$\begin{aligned} \mathcal{D} &= (k \quad \mu \quad m \quad 0) \begin{pmatrix} 1 \\ d_t \\ d_{tt} \\ d_{ttt} \end{pmatrix} \\ &= m \frac{d^2}{dt^2} + \mu \frac{d}{dt} + k \end{aligned}$$

$$L(\theta, \Lambda) = L_p(\theta, \Lambda) + L_d(\theta)$$

$$L_p(\theta, \Lambda) = \frac{1}{N_p} \sum_i^{N_p} \|\Lambda \phi[NN(x_i; \theta)]\|^2 + \|\Lambda\|^2$$

$$L_d(\theta) = \frac{\lambda}{N_d} \sum_l^{N_d} \|NN(x_l; \theta) - u_l\|^2$$

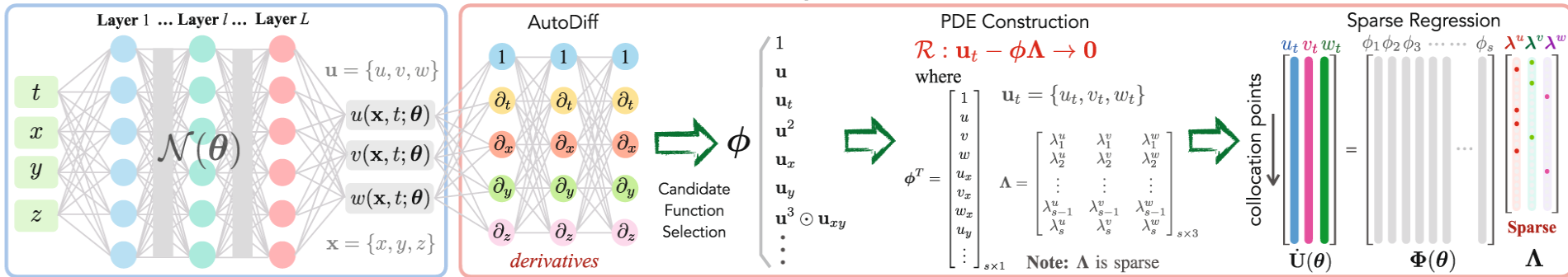
Where Λ are treated as **learnable** parameters and $\{x_l, u_l\}$ are a set of (potentially noisy) observational data

Typically, some regularization / prior on Λ (e.g. sparsity) is needed, as this optimisation problem can be very **ill-posed**

PINNs for equation discovery

DNN with Unknown Parameters θ

Physical Law with Unknown Parameters Λ



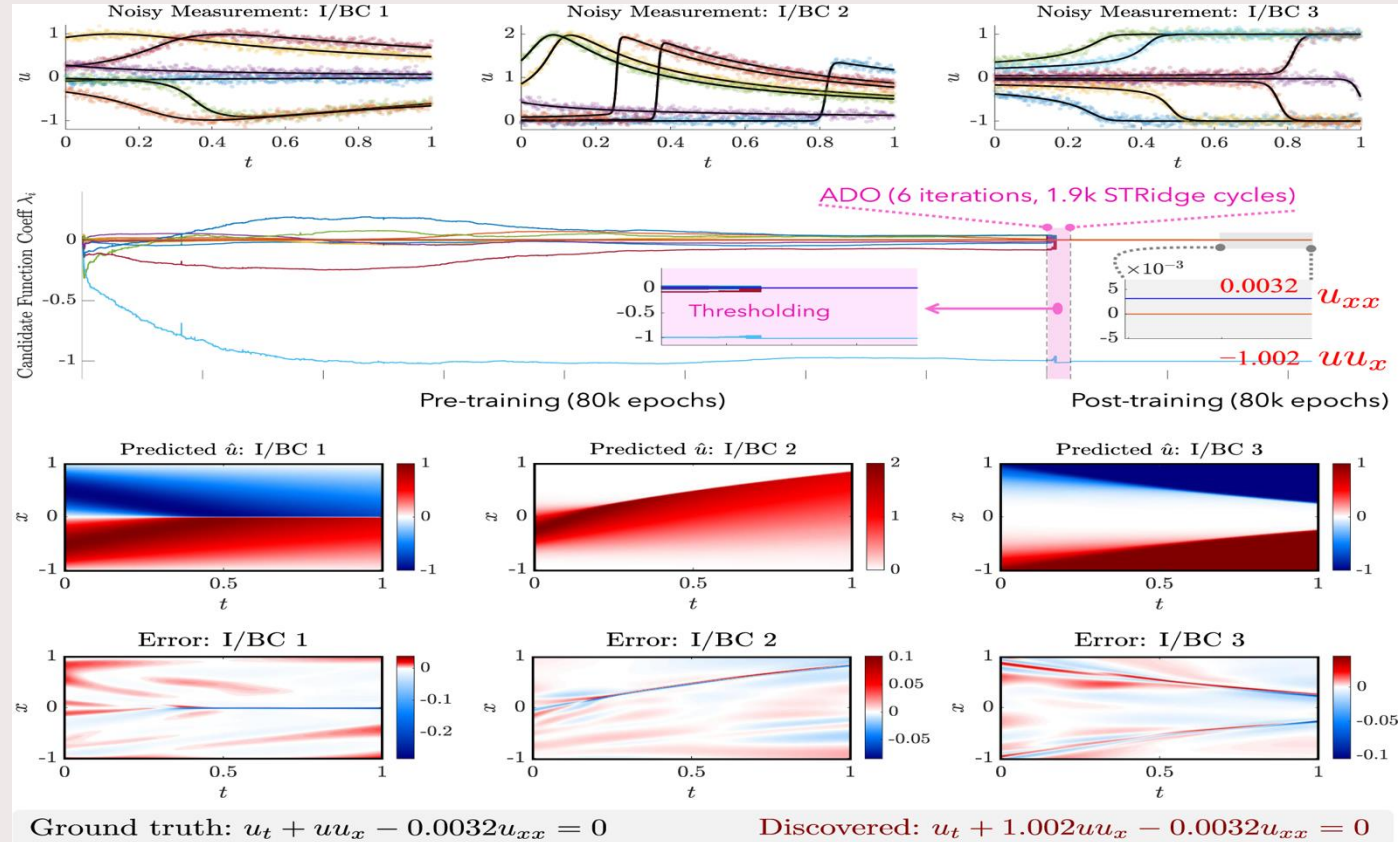
$$\underbrace{\mathcal{L}_d(\theta; \mathcal{D}_u)}_{\text{measurement}} = \frac{1}{N_m} \|\mathbf{u}^\theta - \mathbf{u}^m\|_2^2 \rightarrow \underbrace{\mathcal{L}(\theta, \Lambda; \mathcal{D}_u, \mathcal{D}_c)}_{\text{total loss}} = \underbrace{\mathcal{L}_d(\theta; \mathcal{D}_u)}_{\text{data loss}} + \underbrace{\alpha \mathcal{L}_p(\theta, \Lambda; \mathcal{D}_c)}_{\text{physics loss}} + \underbrace{\beta \|\Lambda\|_0}_{\text{regularization}} \leftarrow \underbrace{\text{Residual Loss}}_{\text{collocation points}} = \frac{1}{N_c} \|\dot{\mathbf{U}}(\theta) - \Phi(\theta) \Lambda\|$$

$$\text{Solution by ADO: } \hat{\Lambda}_{k+1} := \arg \min_{\Lambda} \left[\|\dot{\mathbf{U}}(\hat{\theta}_k) - \Phi(\hat{\theta}_k) \Lambda\|_2^2 + \beta \|\Lambda\|_0 \right] \text{ by STRidge} \quad \hat{\theta}_{k+1} := \arg \min_{\theta} \left[\mathcal{L}_d(\theta; \mathcal{D}_u) + \alpha \mathcal{L}_p(\theta, \hat{\Lambda}_{k+1}; \mathcal{D}_c) \right] \text{ by DNN training}$$

- Trains by alternating between updating Λ and θ

Chen et al, Physics-informed learning of governing equations from scarce data, Nature communications (2021)

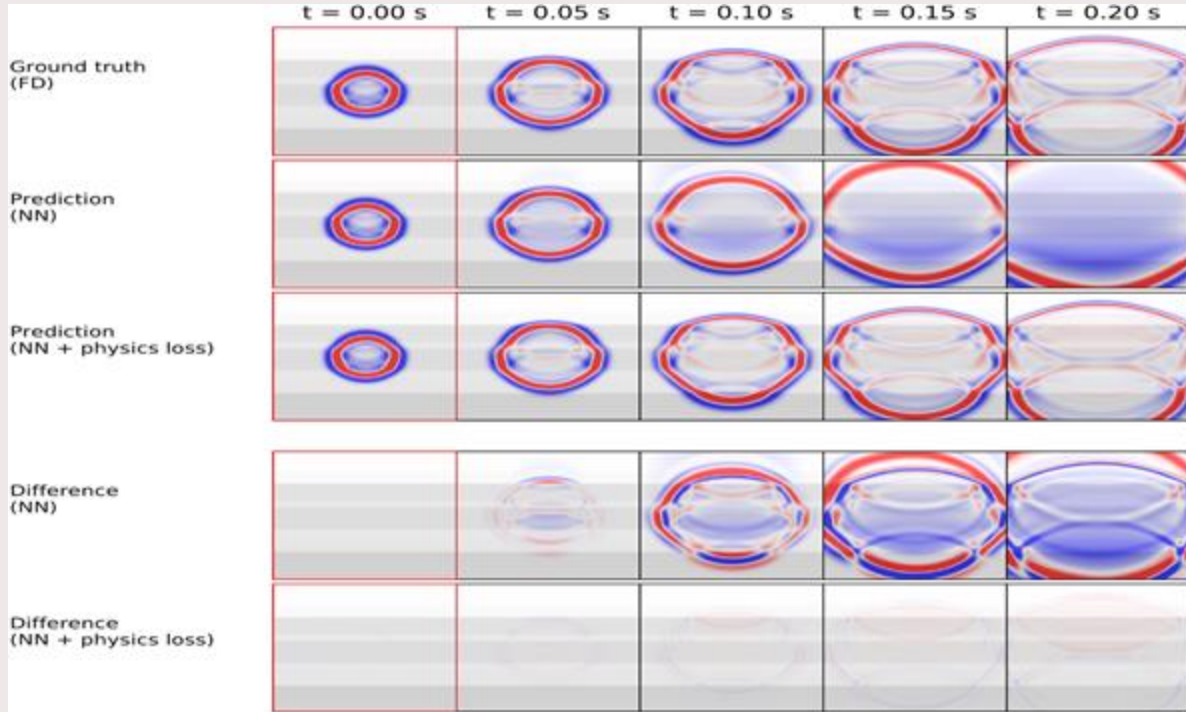
PINNs for equation discovery



- PINN “discovering” Burgers’ equation
- By combining datasets sampled under three different I/BCs with 10% noise

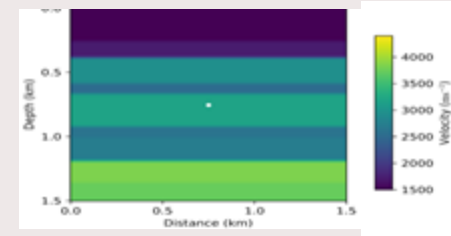
Limitations of PINNs?

Solution time



PINNs need to be retrained for each IC/BC!

Velocity model, $c(x)$



Moseley et al, Solving the wave equation with physics-informed deep learning, ArXiv (2020)

Mini-batch size $N_b = N_p = 500$ (random sampling)

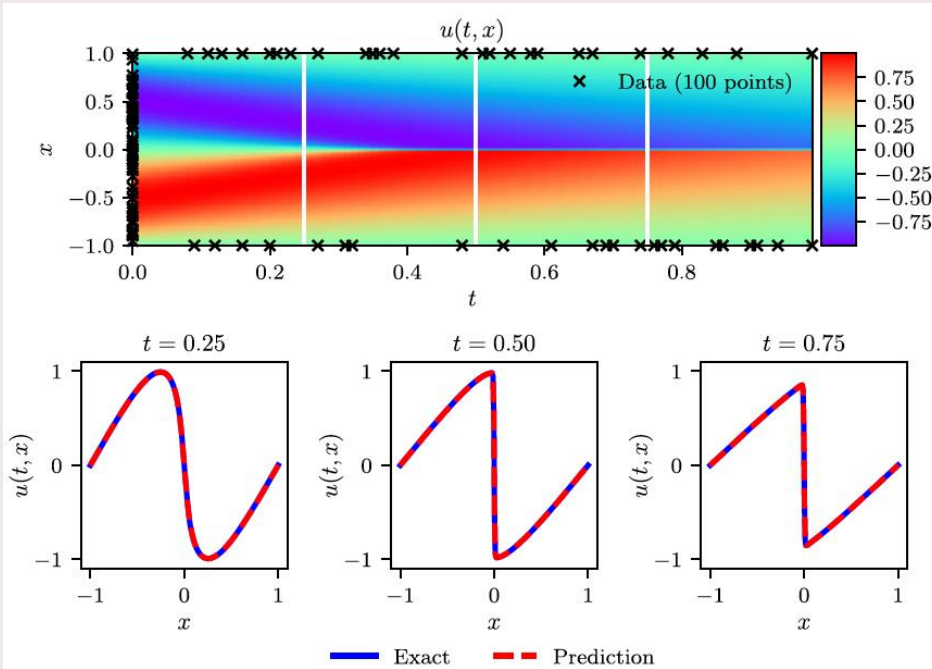
Fully connected network with 10 layers, 1024 hidden units

Softplus activation

Adam optimiser

Training time: **~1 hour**

Competing loss terms



$$L_b(\theta) = \frac{\lambda_1}{N_{b1}} \sum_j^{N_{b1}} (NN(x_j, 0; \theta) + \sin(\pi x_j))^2$$

$$+ \frac{\lambda_2}{N_{b2}} \sum_k^{N_{b2}} (NN(-1, t_k; \theta) - 0)^2$$

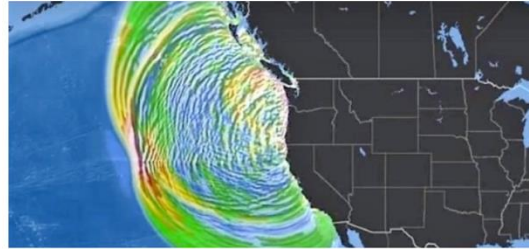
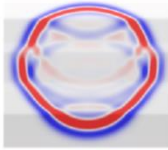
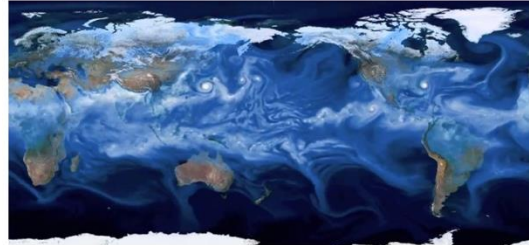
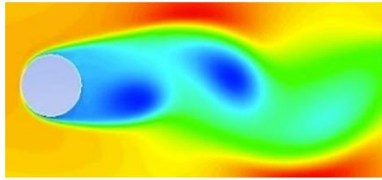
$$+ \frac{\lambda_3}{N_{b3}} \sum_l^{N_{b3}} (NN(+1, t_l; \theta) - 0)^2$$

$$L_p(\theta) = \frac{1}{N_p} \sum_i^{N_p} \left(\left(\frac{\partial NN}{\partial t} + NN \frac{\partial NN}{\partial x} - v \frac{\partial^2 NN}{\partial x^2} \right) (x_i, t_i; \theta) \right)^2$$

How do we choose $\lambda_1, \lambda_2, \lambda_3$? Their relative size can influence the quality of the solution and reinforcement of the BC.

There can be competing terms in the loss functions!

Scaling to more complex problems



Majority of PINN research focuses on **toy/simplified** problems,
as proof-of-principle studies

It is often challenging to **scale** traditional scientific algorithms to:

- More complex phenomena (**multi-scale, multi-physics**)
- Large domains / higher frequency solutions
- Incorporate **real, noisy** and **sparse** data

How do PINNs cope in this setting?

Why SciML models do not scale? ChatGPT says:

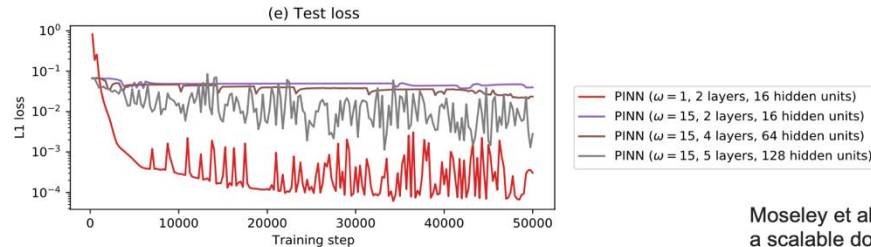
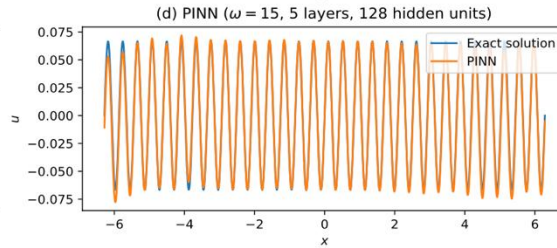
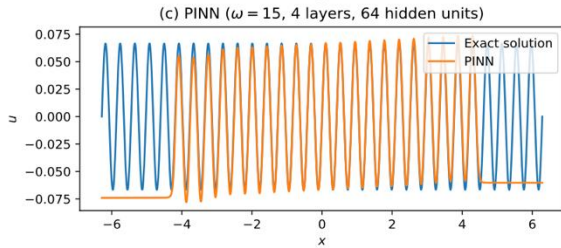
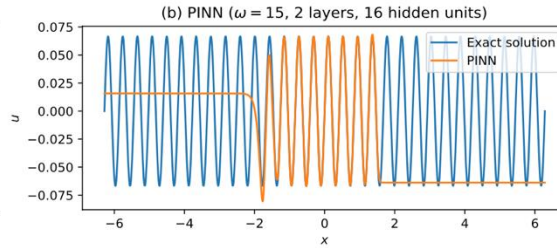
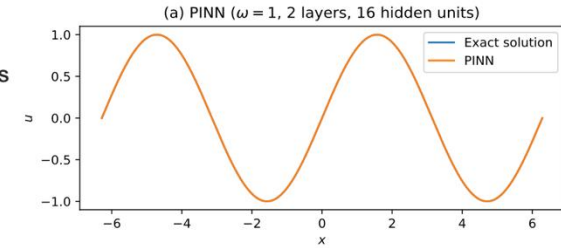
Scaling Scientific Machine Learning (SciML) models presents several challenges. These issues arise from the **unique characteristics of scientific data and models, which differ from those typically used in traditional machine learning**. Here are key reasons why SciML models may not scale well:

1. Complexity of Scientific Data and Models

- **High Dimensionality:** Scientific datasets are often high-dimensional, involving numerous variables and parameters. This can lead to the "curse of dimensionality," where the computational complexity grows exponentially with the number of variables.
- **Non-linear Dynamics:** Many scientific models, such as those in physics, chemistry, or biology, involve highly non-linear relationships, which are challenging for machine learning models to capture efficiently at larger scales.
- **Coupled Systems:** SciML models often represent coupled systems (e.g., climate models, fluid dynamics), where different subsystems interact with each other. This introduces additional layers of complexity that make it difficult to scale models across multiple systems.

Scaling PINNs to higher frequencies

321 free parameters



PINN solving:

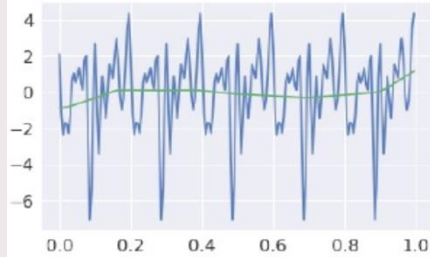
$$\frac{du}{dx} = \cos(\omega x)$$
$$u(0) = 0$$

66,433 free parameters

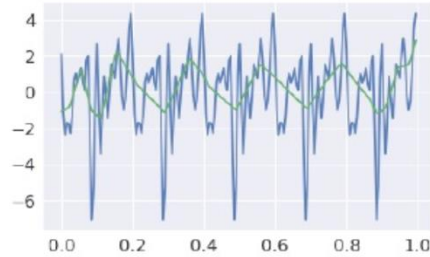
Problem: PINNs **struggle** to solve high-frequency / multiscale problems

Moseley et al, Finite Basis Physics-Informed Neural Networks (FBPINNs): a scalable domain decomposition approach for solving differential equations, ACM (2023)

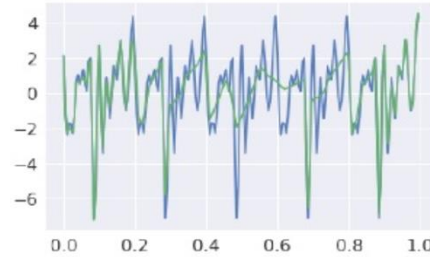
Spectral bias issue



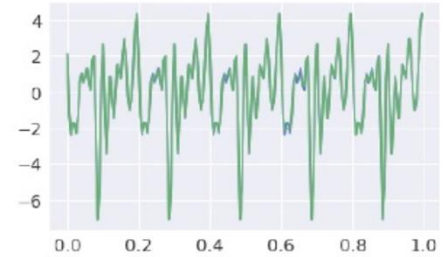
(a) Iteration 100



(b) Iteration 1000



(c) Iteration 10000

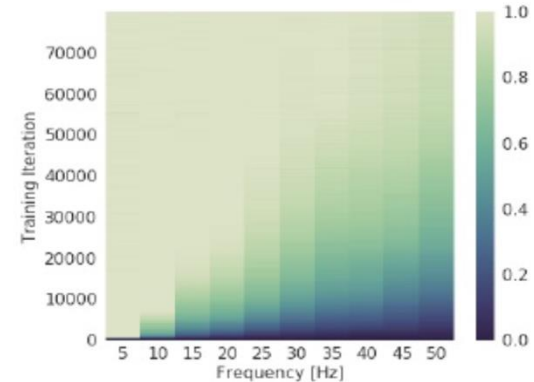


(d) Iteration 80000

NNs prioritise learning **lower** frequency functions first

Under certain assumptions can be proved via neural tangent kernel theory

Rahaman, N., et al, On the spectral bias of neural networks. 36th International Conference on Machine Learning, ICML (2019)



This behaviour can limit the model's ability to scale effectively, especially when the goal is to approximate complex scientific data with both large-scale and fine-scale dynamics.

PINNs – an entire research field

SPRINGER LINK

Find a journal Publish with us Track your research Search

Home > Journal of Scientific Computing > Article

Scientific Machine Learning Through Physics-Informed Neural Networks: Where we are and What's Next

Open access | Published: 26 July 2022

Volume 92, article number 88, (2022) Cite this article

Download PDF

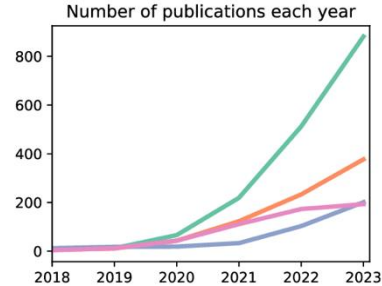
You have full access to this open access article

Salvatore Cuomo, Vincenzo Schiano Di Cola, Fabio Giampaolo, Gianluigi Rozza, Maziar Raissi & Francesco Piccialli

67k Accesses 274 Citations 7 Altmetric Explore all metrics

Abstract

Physics-Informed Neural Networks (PINN) are neural networks (NNs) that encode model equations, like Partial Differential Equations (PDE), as a component of the neural network itself. PINNs are nowadays used to solve PDEs, fractional equations, integral-differential equations, and stochastic PDEs. This novel methodology has arisen as a multi-task learning framework in which a NN must fit observed data while reducing a PDE residual. This article provides a comprehensive review of the literature on PINNs: while the primary goal of the study was to characterize these networks and their related advantages and



physics-informed neural networks
scientific machine learning / physics-informed ML / AI for science
operator learning / neural operators
differentiable physics / neural differential equations

Source: Scopus keyword search (Feb 2024)

[\[HTML\] Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations](#)

[M Raissi](#), [P Perdikaris](#), [GE Karniadakis](#) - Journal of Computational physics, 2019 - Elsevier

... We introduce **physics-informed neural networks – neural networks** that are trained to solve supervised learning tasks while respecting any given laws of physics described by general ...

☆ Save 📄 Cite Cited by 8017 Related articles All 7 versions

100-1000s of PINN applications / extensions!

Advantages / limitations of PINNs

Advantages

- **Mesh-free**
- Can jointly solve **forward** and **inverse** problems
- Often performs well on “messy” problems (where some observational data is available)
- Tractable, analytical solution gradients (e.g. for sensitivity analysis)
- Mostly **unsupervised**

Limitations

- **Computational cost** often high (especially for forward-only problems)
- Can be hard to **optimise** (and convergence properties less well understood)
- Challenging to **scale** to larger domains, multi-scale, multi-physics problems

But.. many PINN extensions exist!



PINNs and their extensions

VICTORITA DOLEAN

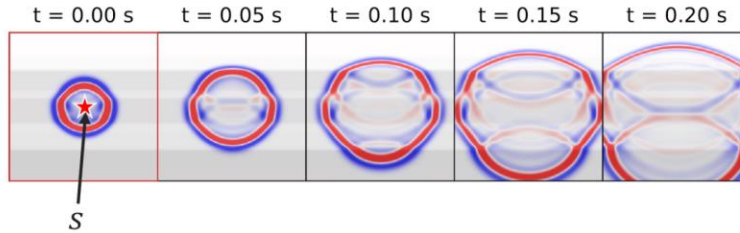


Limitation 1) – computational cost

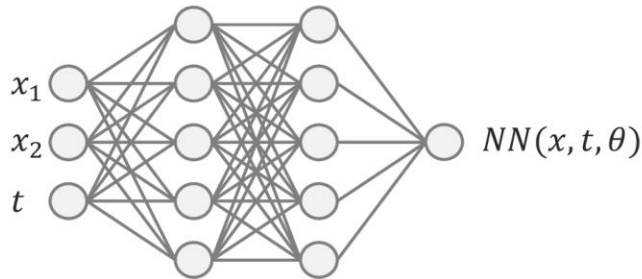
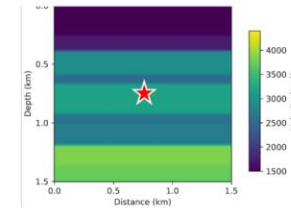
Conditioned PINNs

Idea: add I/BCs / other PDE parameters as **additional** network input parameters

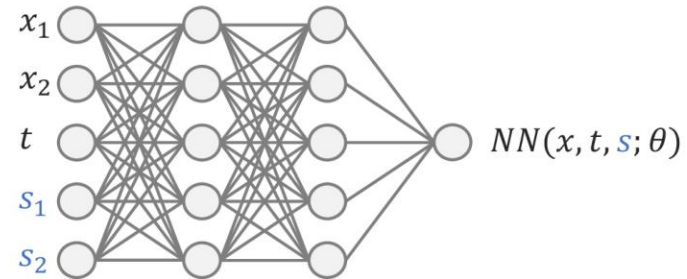
Ground truth FD simulation



Velocity model, $c(x)$



PINN

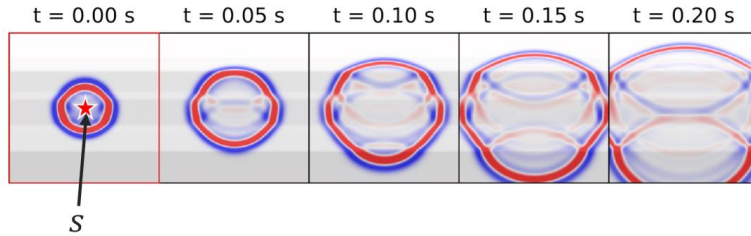


Conditioned PINN

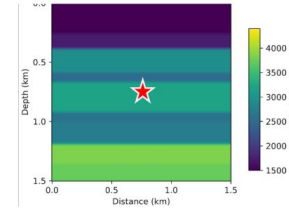
Conditioned PINNs

Idea: add I/BCs / other PDE parameters as **additional** network input parameters

Ground truth FD simulation



Velocity model, $c(x)$



$$L_b(\theta) = \frac{\lambda}{N_b} \sum_j^{N_b} \left(NN(x_j, t_j; \theta) - u_{FD}(x_j, t_j) \right)^2$$

$$L_p(\theta) = \frac{1}{N_p} \sum_i^{N_p} \left(\left[\nabla^2 - \frac{1}{c(x_i)^2} \frac{\partial^2}{\partial t^2} \right] NN(x_i, t_i; \theta) \right)^2$$

PINN

$$L_b(\theta) = \frac{\lambda}{N_b} \sum_j^{N_b} \left(NN(x_j, t_j, s_j; \theta) - \underline{u_{FD}(x_j, t_j, s_j)} \right)^2$$

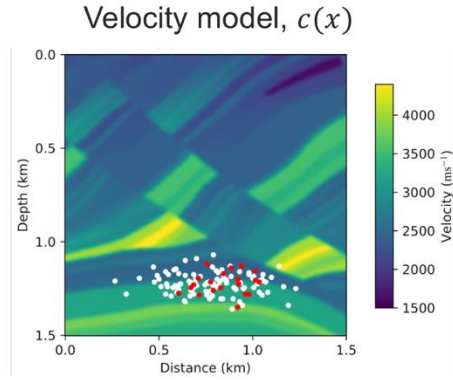
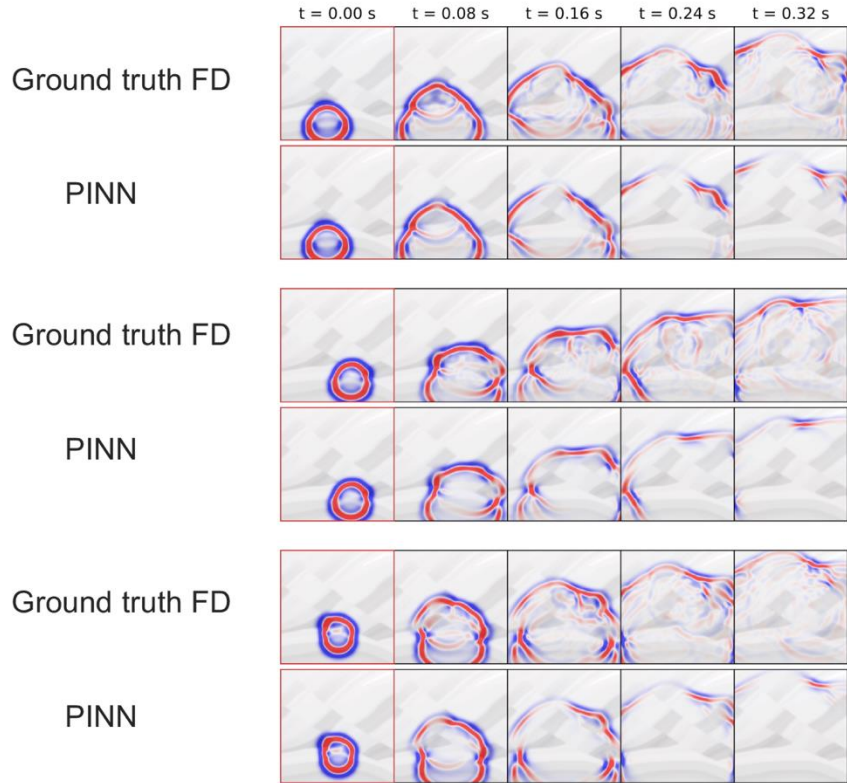
$$L_p(\theta) = \frac{1}{N_p} \sum_i^{N_p} \left(\left[\nabla^2 - \frac{1}{c(x_i)^2} \frac{\partial^2}{\partial t^2} \right] NN(x_i, t_i, \underline{s_i}; \theta) \right)^2$$

Conditioned PINN

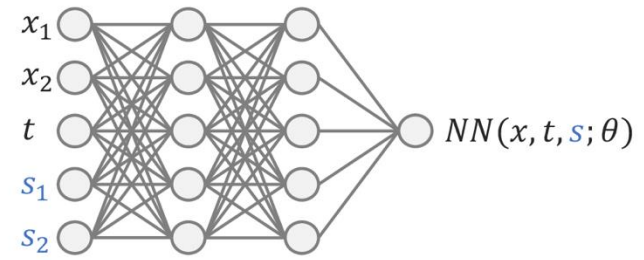
Boundary examples from many **different** source locations

Randomly sampled input coordinates **and** source locations

Conditioned PINNs



White = random source locations used for training
Red = source locations used for testing



Means the network **does not** need to be retrained for each simulation => much faster!

Aka a **surrogate** model

Moseley et al, Solving the wave

Physics-informed deep operator networks (DeepONets)

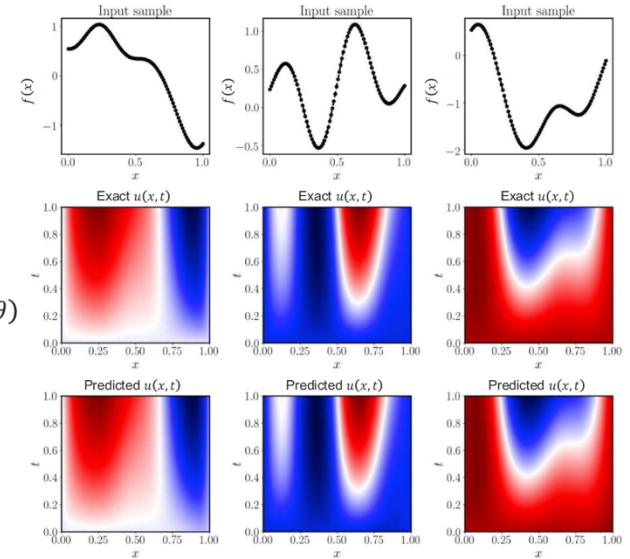
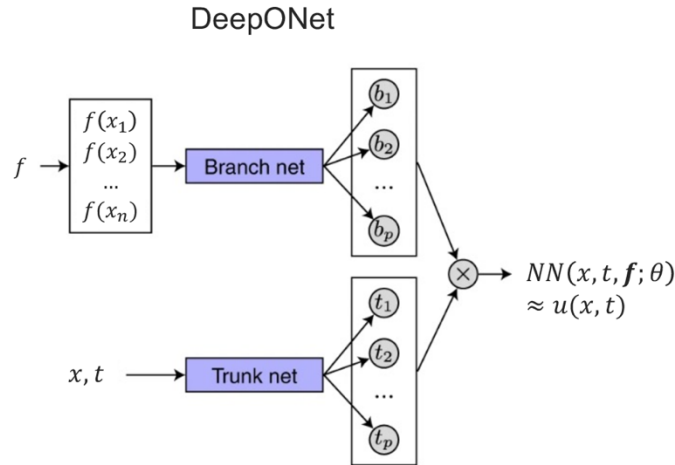
Conditioned PINN for solving reaction-diffusion equation:

$$\frac{\partial u}{\partial t} = D \frac{\partial^2 u}{\partial x^2} + ku^2 + f(x)$$

Input:
 n discretised values of $f(x)$

Output:
 $NN(x, t, f; \theta) \approx u(x, t)$

Trained using many examples of $f(x)$

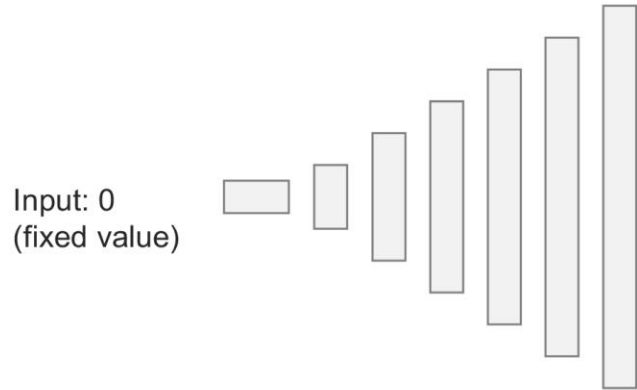


Wang et al, Learning the solution operator of parametric partial differential equations with physics-informed DeepONets, Science Advances (2021)

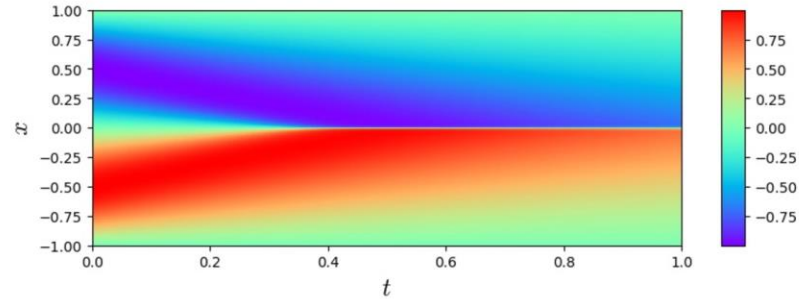
Lu et al, Learning nonlinear operators via DeepONet based on the universal approximation theorem of operators, Nature Machine Intelligence (2021)

Discretised PINNs

Idea: **discretise** solution and use e.g. convolutional network to learn spatial/temporal correlations



Architecture:
(Transposed) convolutional layers



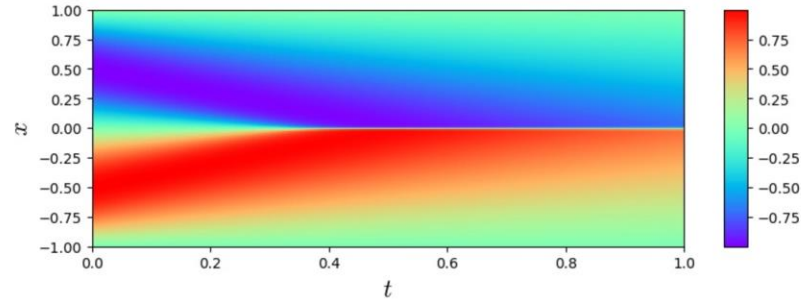
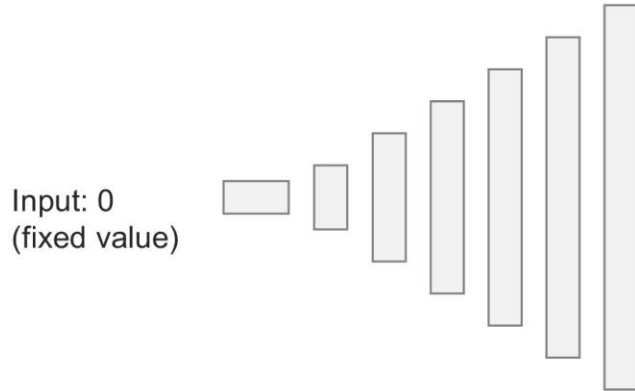
Output: **discretised** solution
(a matrix of values)
e.g. 128 x 256

$$NN(\theta)_{ij} \approx u(x = x_i, t = t_j)$$

Instead of: $NN(x, t; \theta) \approx u(x, t)$

Discretised PINNs

Idea: **discretise** solution and use e.g. convolutional network to learn spatial/temporal correlations



$$NN(x, t; \theta) \approx u(x, t)$$

$$L_p(\theta) = \frac{1}{N_p} \sum_i^{N_p} \left(\left(\frac{\partial NN}{\partial t} + NN \frac{\partial NN}{\partial x} - v \frac{\partial^2 NN}{\partial x^2} \right) (x_i, t_i; \theta) \right)^2$$

PINN

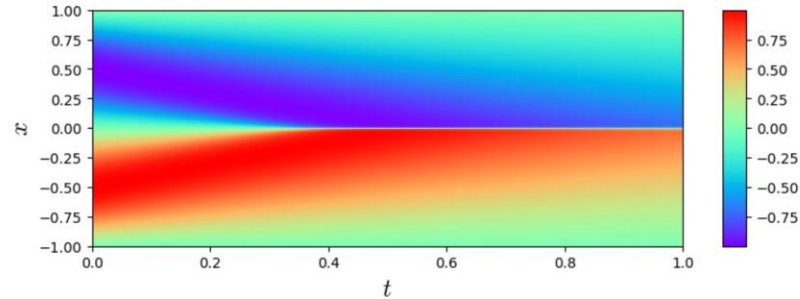
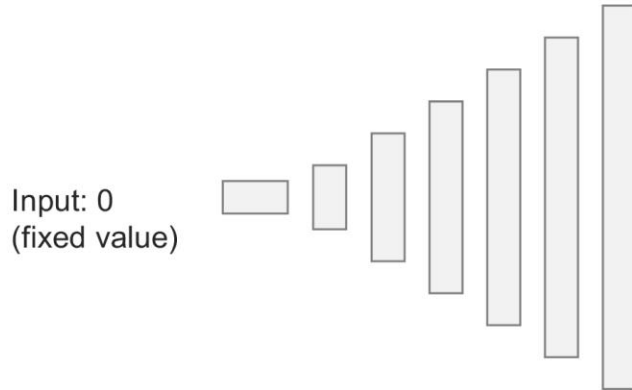
$$NN(\theta)_{ij} \approx u(x = x_i, t = t_j)$$

$$L_p(\theta) = \frac{1}{N_x N_t} \sum_i^{N_x} \sum_j^{N_t} \left(\frac{\delta NN(\theta)_{ij}}{\delta t} + NN(\theta)_{ij} \frac{\delta NN(\theta)_{ij}}{\delta x} - v \frac{\delta^2 NN(\theta)_{ij}}{\delta x^2} \right)^2$$

Discretised PINN

Discretised PINNs

Idea: **discretise** solution and use e.g. convolutional network to learn spatial/temporal correlations



Derivatives in loss function are **approximated** using **finite difference** filters, e.g.

$$\frac{\delta NN(\theta)_{ij}}{\delta t} = \frac{NN(\theta)_{ij} - NN(\theta)_{i,j-1}}{t_j - t_{j-1}}$$

$$NN(\theta)_{ij} \approx u(x = x_i, t = t_j)$$

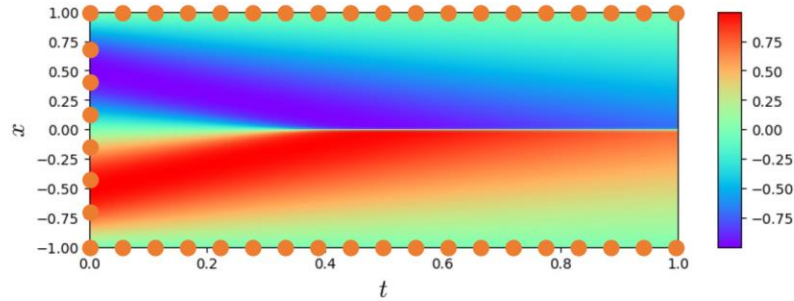
$$L_p(\theta) = \frac{1}{N_x N_t} \sum_i \sum_j \left(\frac{\delta NN(\theta)_{ij}}{\delta t} - NN(\theta)_{ij} \frac{\delta NN(\theta)_{ij}}{\delta x} - v \frac{\delta^2 NN(\theta)_{ij}}{\delta x^2} \right)^2$$

Discretised PINN

Autodiff is still used to update θ

Discretised PINNs

Idea: **discretise** solution and use e.g. convolutional network to learn spatial/temporal correlations



Initial/boundary conditions are asserted by **padding** the edges of the output solution with **appropriate** values (=hard constraint), e.g.

$$NN(\theta)_{i,j=0} = -\sin(\pi x_i)$$

$$NN(\theta)_{i=0,j} = NN(\theta)_{i=128,j} = \underline{0}$$

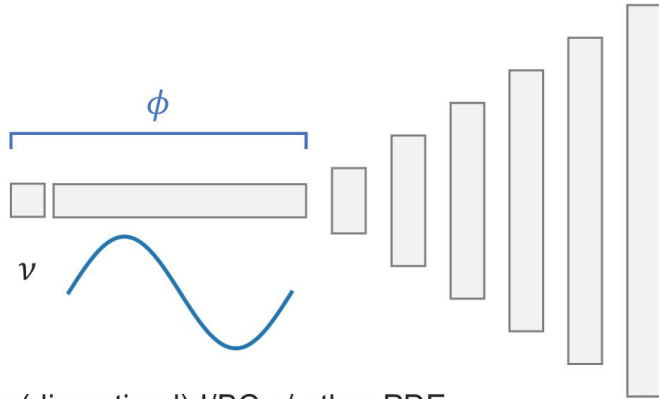
$$NN(\theta)_{ij} \approx u(x = x_i, t = t_j)$$

$$L_p(\theta) = \frac{1}{N_x N_t} \sum_i^{N_x} \sum_j^{N_t} \left(\frac{\delta NN(\theta)_{ij}}{\delta t} + NN(\theta)_{ij} \frac{\delta NN(\theta)_{ij}}{\delta x} - v \frac{\delta^2 NN(\theta)_{ij}}{\delta x^2} \right)^2$$

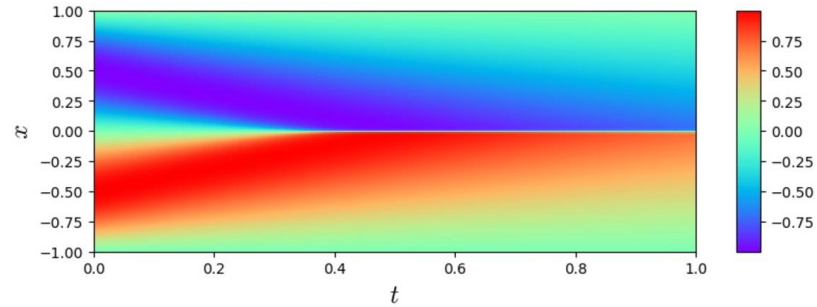
Discretised PINN

Conditioned discretised PINNs

Idea: **discretise** solution and use e.g. convolutional network to learn spatial/temporal correlations
 And condition on I/BCs / other PDE parameters



Input: (discretised) I/BCs / other PDE parameters



$$NN(\phi; \theta)_{ij} \approx u(x = x_i, t = t_j)$$

$$L_p(\theta) = \frac{1}{N_x N_t N_\phi} \sum_k^{N_\phi} \sum_i^{N_x} \sum_j^{N_t} \left(\frac{\delta NN(\phi_k; \theta)_{ij}}{\delta t} + NN(\phi_k; \theta)_{ij} \frac{\delta NN(\phi_k; \theta)_{ij}}{\delta x} - \phi_{k0} \frac{\delta^2 NN(\phi_k; \theta)_{ij}}{\delta x^2} \right)^2$$

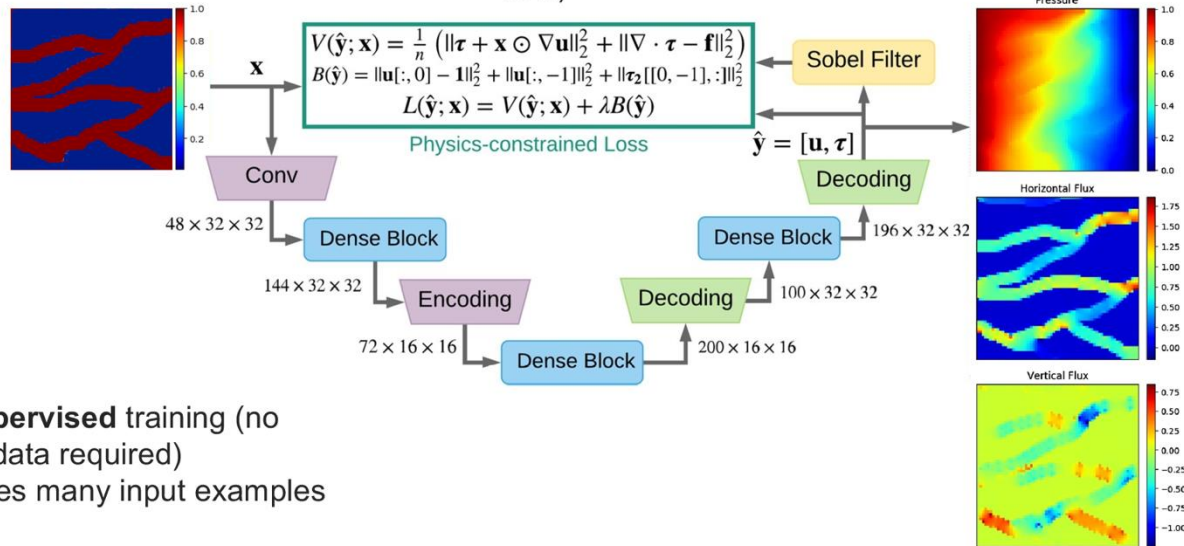
Conditioned discretised PINN

Conditioned discretised PINNs

Input:
permeability field

PDE: steady-state flow in
heterogeneous media (Darcy
flow)

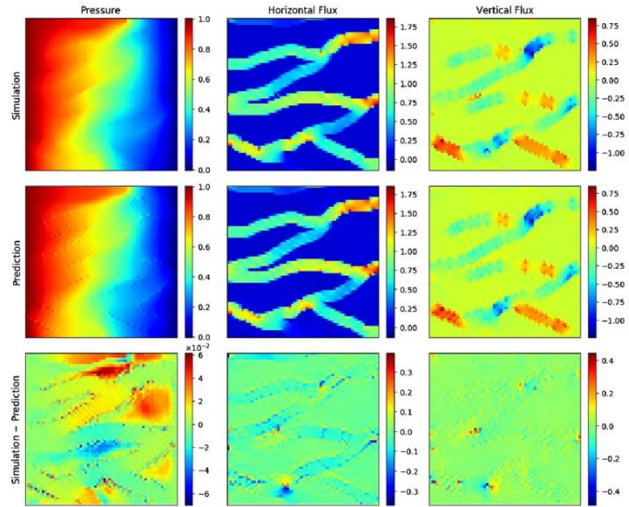
Output: fluid pressure
and flux field



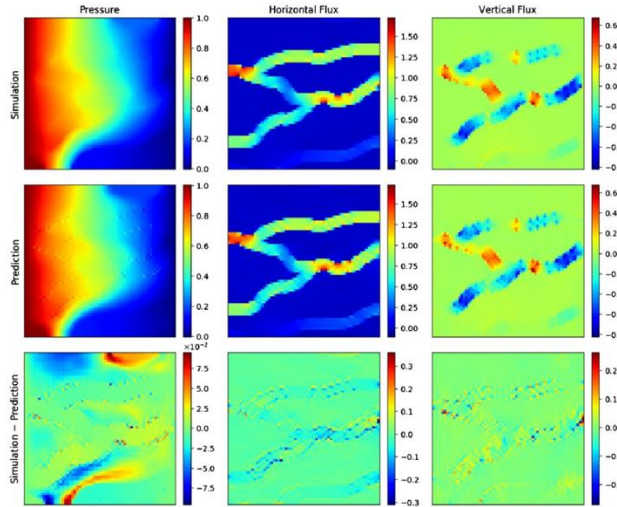
- Fully **unsupervised** training (no simulation data required)
- Only requires many input examples to train

Zhu, Y et al, Physics-constrained deep learning for high-dimensional surrogate modeling and uncertainty quantification without labeled data. Journal of Computational Physics (2019)

Conditioned discretised PINNs

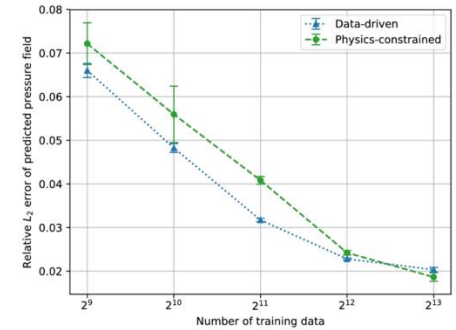


(c) Channelized, test 1.

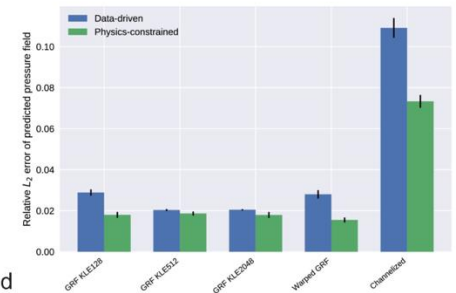


(d) Channelized, test 2.

Physics-informed vs fully data-driven CNN



Generalisation to different input distributions



Zhu, Y et al, Physics-constrained deep learning for high-dimensional surrogate modeling and uncertainty quantification without labeled data. Journal of Computational Physics (2019)

Advantages / limitations of discretised PINNs

Advantages

- Allows the use of CNNs to exploit spatial correlations between inputs/outputs of PDE
- Can be extended to:
 - Irregular geometries (Graph NNs)
 - Explicit time dependence (RNNs)
 - Mixed continuous/discrete input coordinates

Limitations

- Relies on approximate FD gradients
- Only outputs discretised solution; needs to be retrained to output on larger domains / finer grids

Geneva et al, Modeling the dynamics of PDE systems with physics-constrained deep auto-regressive networks, JCP (2020)

Gao et al, PhyGeoNet: Physics-informed geometry-adaptive convolutional neural networks for solving parameterized steady-state PDEs on irregular domain, JCP (2021)

Training PINNs with finite differences

Most time is spent **computing gradients**, not the forward pass, when training PINNs

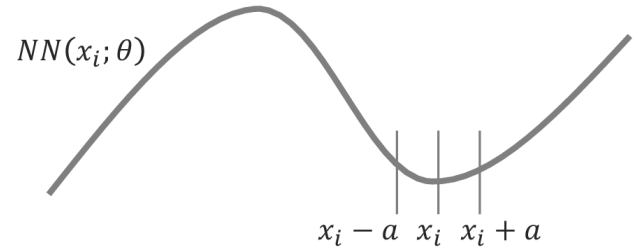
$$L_p(\theta) = \frac{1}{N_p} \sum_i^{N_p} \|\mathcal{D}[NN(x_i; \theta)] - f(x_i)\|^2$$

Idea: instead of using exact gradients from autodifferentiation, use **approximate** gradients from **finite differences**

For each collocation point x_i :

1. Sample a stencil of input points around x_i
2. Run forward pass of network with all these points
3. Approximate derivatives in \mathcal{D} using finite differences
4. Compute loss function using these derivatives

Autodiff is still used to update θ



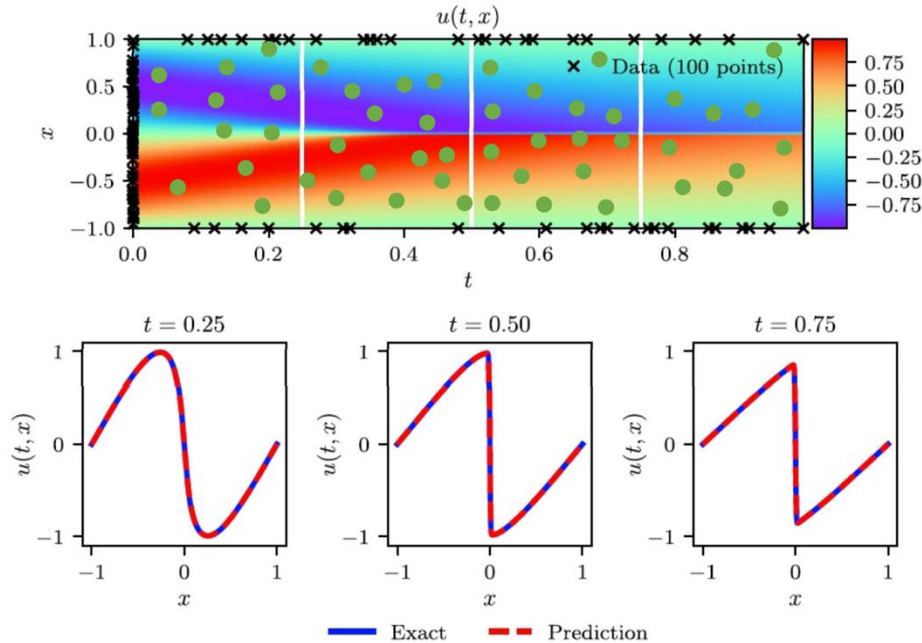
$$\frac{\partial NN(x_i; \theta)}{\partial x} \approx \frac{NN(x_i + a; \theta) - NN(x_i - a; \theta)}{2a}$$

Can offer 2-4x speedups depending on PDE

But choosing a suitable value of a is critical

Limitation 2) – poor convergence

Competing loss terms



Raissi et al, Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations, JCP (2018)

$$L_b(\theta) = \frac{\lambda_1}{N_{b1}} \sum_j^{N_{b1}} (NN(\underline{x}_j, 0; \theta) + \underline{\sin(\pi x_j)})^2$$

$$+ \frac{\lambda_2}{N_{b2}} \sum_k^{N_{b2}} (NN(-1, \underline{t}_k; \theta) - \underline{0})^2$$

$$+ \frac{\lambda_3}{N_{b3}} \sum_l^{N_{b3}} (NN(+1, \underline{t}_l; \theta) - \underline{0})^2$$


$$L_p(\theta) = \frac{1}{N_p} \sum_i^{N_p} \left(\left(\frac{\partial NN}{\partial t} + NN \frac{\partial NN}{\partial x} - v \frac{\partial^2 NN}{\partial x^2} \right) (\underline{x}_i, \underline{t}_i; \theta) \right)^2$$

How do we choose λ_1 , λ_2 , and λ_3 ?

λ too small => doesn't learn unique solution
 λ too large => only learns boundary condition

Thus, there can be **competing** terms in the loss function

Hard initial / boundary conditions

 Idea: use neural network as part of a solution **ansatz**

Burgers' problem above:

$$\frac{\partial u}{\partial t} + u \frac{\partial u}{\partial x} - \nu \frac{\partial^2 u}{\partial x^2} = 0$$

$$\begin{aligned} u(x, 0) &= -\sin(\pi x) \\ u(-1, t) &= u(+1, t) = 0 \end{aligned}$$

Then let the solution approximated by

$$\begin{aligned} \hat{u}(x, t; \theta) &= (x - 1)(x + 1)(t - 0)NN(x, t; \theta) - \sin(\pi x) \\ &\approx u(x, t) \end{aligned}$$

It is easy to verify that the initial / boundary conditions are now satisfied by **construction**

Lagaris et al, Artificial neural networks for solving ordinary and partial differential equations, IEEE (1998)

Hard initial / boundary conditions



Idea: use neural network as part of a solution **ansatz**

Burgers' problem above:

$$\frac{\partial u}{\partial t} + u \frac{\partial u}{\partial x} - \nu \frac{\partial^2 u}{\partial x^2} = 0$$

$$\begin{aligned} u(x, 0) &= -\sin(\pi x) \\ u(-1, t) &= u(+1, t) = 0 \end{aligned}$$

Then let the solution approximated by

$$\begin{aligned} \hat{u}(x, t; \theta) &= (x - 1)(x + 1)(t - 0)NN(x, t; \theta) - \sin(\pi x) \\ &\approx u(x, t) \end{aligned}$$

It is easy to verify that the initial / boundary conditions are now satisfied by **construction**

Lagaris et al, Artificial neural networks for solving ordinary and partial differential equations, IEEE (1998)

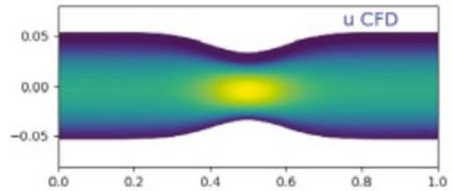
Thus, the boundary loss $L_b(\theta)$ is not needed and the problem is turned into an **unconstrained** optimisation problem, where we minimize

$$L(\theta) = L_p(\theta) = \frac{1}{N_p} \sum_i^{N_p} \left(\left(\frac{\partial \hat{u}}{\partial t} + \hat{u} \frac{\partial \hat{u}}{\partial x} - \nu \frac{\partial^2 \hat{u}}{\partial x^2} \right) (x_i, t_i; \theta) \right)^2$$

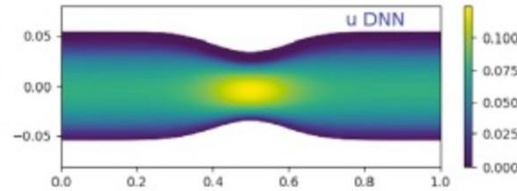
Note:

- You can think of the ansatz as a custom final layer of the network, asserting a **hard** constraint on the network's output
- Autodifferentiation lets us differentiate through ansatz
- Training is now fully **unsupervised** (only need collocation points)
- Can be challenging to use this approach for complex boundary conditions

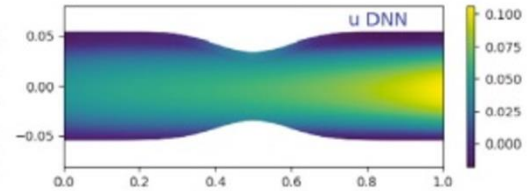
Hard vs soft PINN



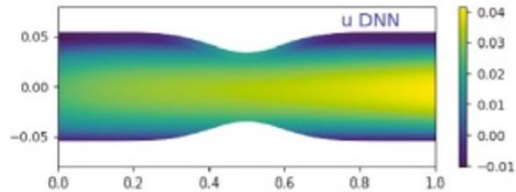
(a) CFD benchmark



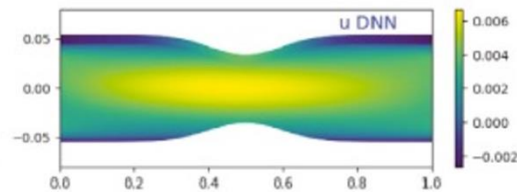
(b) DNN hard



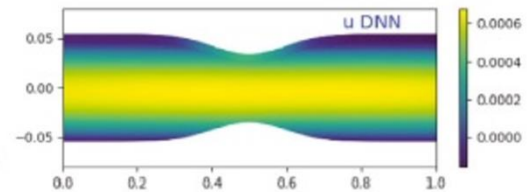
(c) DNN soft $\lambda = 1$



(d) DNN soft $\lambda = 10$



(e) DNN soft $\lambda = 100$



(f) DNN soft $\lambda = 1000$

Sun et al, Surrogate modeling for fluid flows based on physics-constrained deep learning without simulation data, Computer Methods in Applied Mechanics and Engineering (2020)

Soft PINN with varying λ vs hard PINN
When solving incompressible Navier-Stokes equations

Adaptive lambdas

Idea: Specify **separate** λ for each training point
And treat λ as **learnable**

$$L(\theta) = L_b(\theta) + L_p(\theta)$$

$$L_b(\theta, \lambda) = \sum_k \sum_j \lambda_{kj} \|\mathcal{B}_k[NN(x_{kj}; \theta)] - g_k(x_{kj})\|^2$$

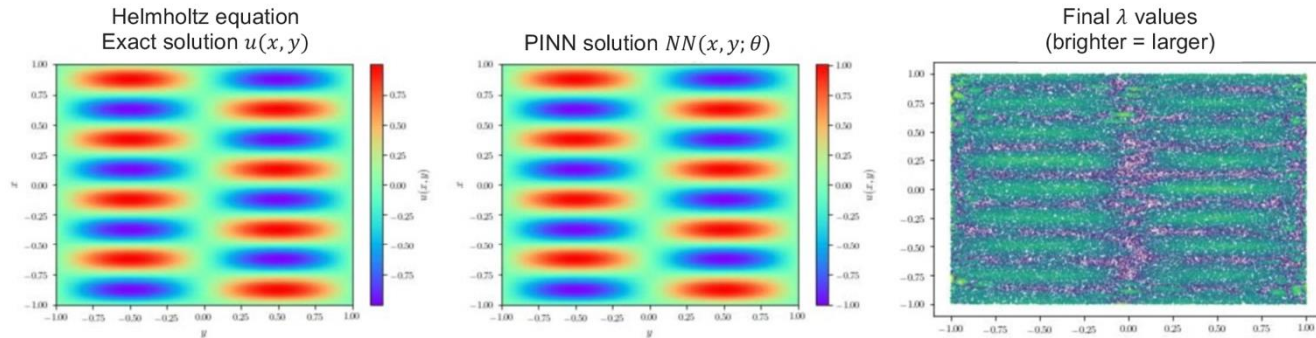
$$L_p(\theta, \lambda) = \sum_i \lambda_i \|\mathcal{D}[NN(x_i; \theta)] - f(x_i)\|^2$$

Loop:

1. Compute gradients $\frac{\partial L}{\partial \theta_l}, \frac{\partial L}{\partial \lambda_l}$
2. Update weights,
$$\theta_l \leftarrow \theta_l - \gamma \frac{\partial L}{\partial \theta_l} \text{ (minimise } L)$$
3. Update lambdas,
$$\lambda_l \leftarrow \lambda_l + \gamma \frac{\partial L}{\partial \lambda_l} \text{ (maximise } L)$$

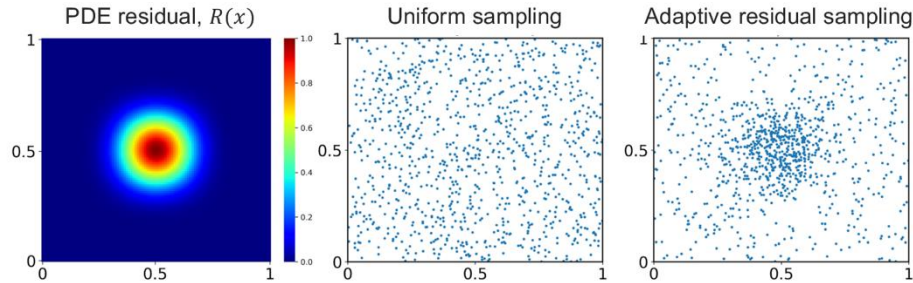
Note, for example, $\frac{\partial L}{\partial \lambda_i} = \|\mathcal{D}[NN(x_i; \theta)] - f(x_i)\|^2$

λ can be thought of as a soft attention layer



McClenny et al, Self-adaptive physics-informed neural networks, JCP (2023)

Adaptive collocation points



Idea: sample collocation points according to the **probability distribution**

$$x \sim p(x) = \frac{R(x)}{A}$$

Where R is the PDE residual, $R(x) = \|\mathcal{D}[NN(x; \theta)] - f(x)\|^2$

And A is a normalising constant

Inspired by adaptive mesh refinement in FEM

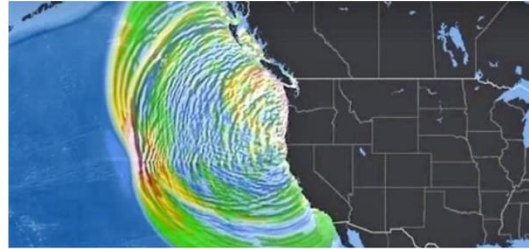
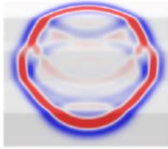
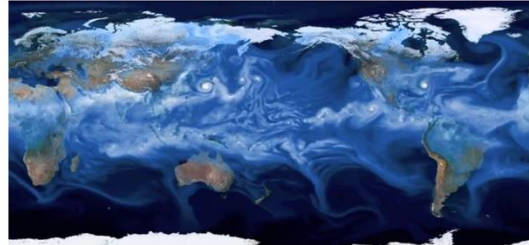
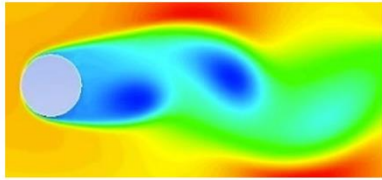
Table 2. L^2 relative error of the PINN solution for the forward problems.

	Diffusion	Burgers'	Allen-Cahn	Wave
No. of residual points	30	2000	1000	2000
Grid	0.66±0.06%	13.7±2.37%	93.4±6.98%	81.3±13.7%
Random	0.74±0.17%	13.3±8.35%	22.2±16.9%	68.4±20.1%
LHS	0.48±0.24%	13.5±9.05%	26.6±15.8%	75.9±33.1%
Halton	0.24±0.17%	4.51±3.93%	0.29±0.14%	60.2±10.0%
Hammersley	0.17±0.07%	3.02±2.98%	0.14±0.14%	58.9±8.52%
Sobol	0.19±0.07%	3.38±3.21%	0.35±0.24%	57.5±14.7%
Random-R	0.12±0.06%	1.69±1.67%	0.55±0.34%	0.72±0.90%
RAR-G [3]	0.20±0.07%	0.12±0.04%	0.53±0.19%	0.81±0.11%
RAD	0.11±0.07%	0.02±0.00%	0.08±0.06%	0.09±0.04%
RAR-D	0.14±0.11%	0.03±0.01%	0.09±0.03%	0.29±0.04%

Wu et al, A comprehensive study of non-adaptive and residual-based adaptive sampling for physics-informed neural networks, Computer Methods in Applied Mechanics and Engineering (2023)

Limitation 3) – scaling to more complex problems

Scaling to more complex problems



Majority of PINN research focuses on **toy/simplified** problems,
as proof-of-principle studies

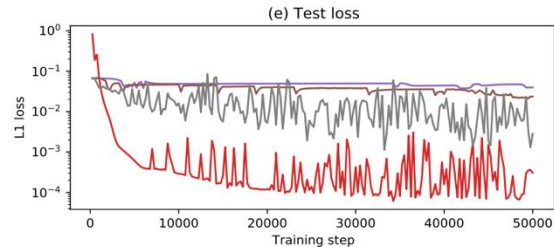
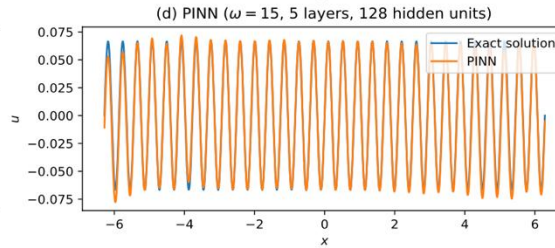
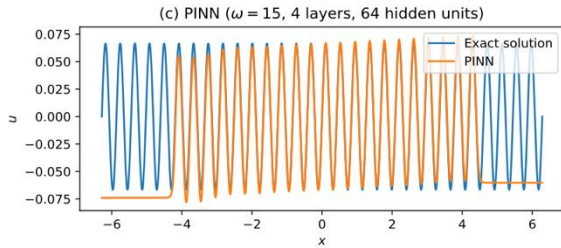
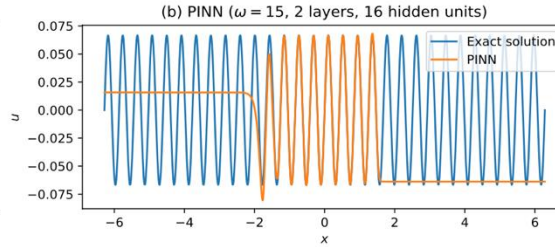
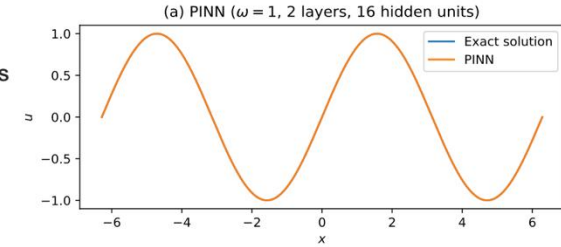
It is often challenging to **scale** traditional scientific algorithms to:

- More complex phenomena (**multi-scale, multi-physics**)
- Large domains / higher frequency solutions
- Incorporate **real, noisy** and **sparse** data

How do PINNs cope in this setting?

Scaling PINNs to higher frequencies

321 free parameters



PINN solving:

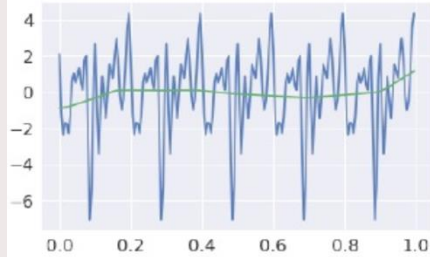
$$\frac{du}{dx} = \cos(\omega x)$$
$$u(0) = 0$$

66,433 free parameters

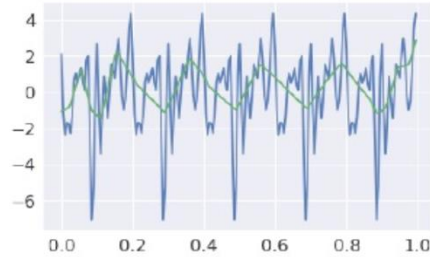
Problem: PINNs **struggle** to solve high-frequency / multiscale problems

Moseley et al, Finite Basis Physics-Informed Neural Networks (FBPINNs): a scalable domain decomposition approach for solving differential equations, ACM (2023)

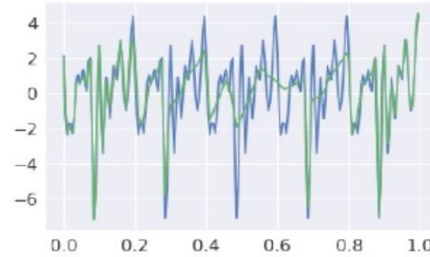
Spectral bias issue



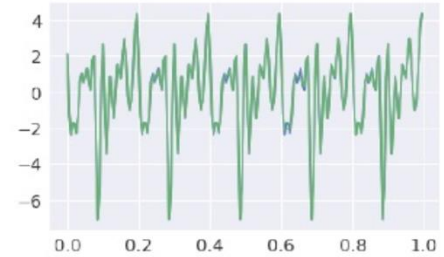
(a) Iteration 100



(b) Iteration 1000



(c) Iteration 10000

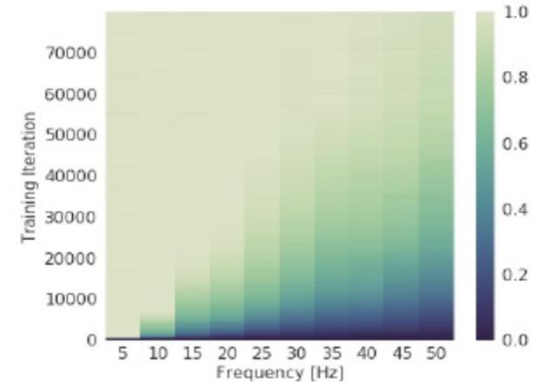


(d) Iteration 80000

NNs prioritise learning **lower** frequency functions first

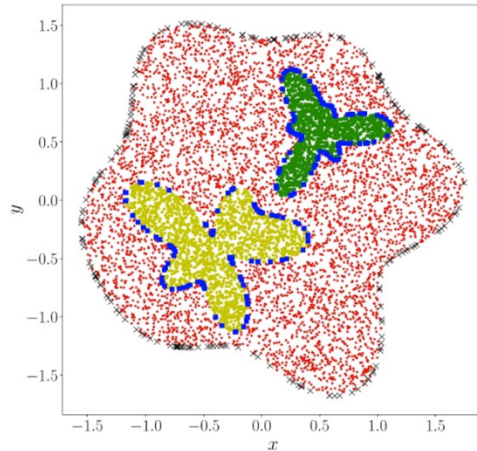
Under certain assumptions can be proved via neural tangent kernel theory

Rahaman, N., et al, On the spectral bias of neural networks. 36th International Conference on Machine Learning, ICML (2019)



This behaviour can limit the model's ability to scale effectively, especially when the goal is to approximate complex scientific data with both large-scale and fine-scale dynamics.

PINNs + domain decomposition



Idea:

Take a “**divide-and-conquer**” strategy to model more complex problems:

1. Divide modelling domain into many smaller **subdomains**
2. Use a separate neural network in each subdomain to model the solution

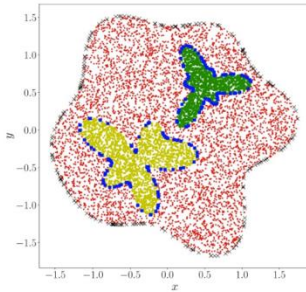
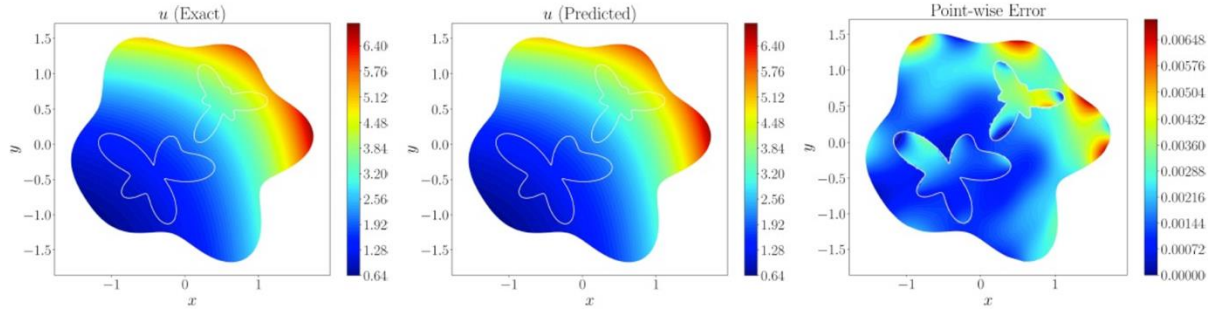
Hypothesis:

The resulting (coupled) local optimization problems are easier to solve than a single global problem

Jagtap, A., et al., Extended physics-informed neural networks (XPINNs): A generalized space-time domain decomposition based deep learning framework for nonlinear partial differential equations. Communications in Computational Physics (2020)

Extended PINNs (XPINNs)

XPINNs solving 2D
Poisson's equation

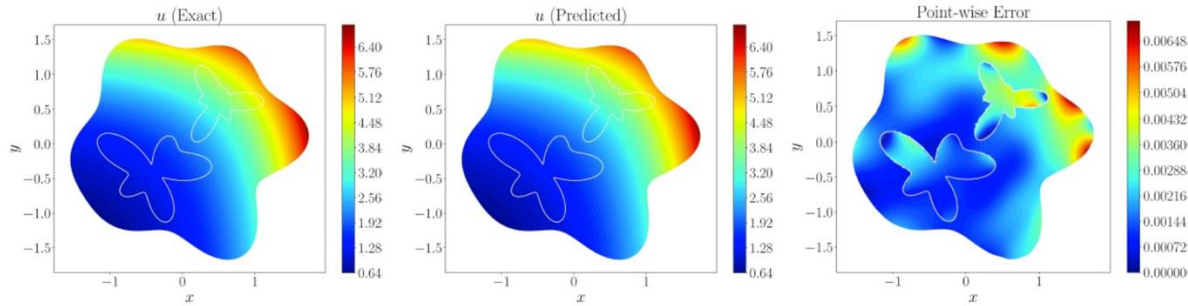


$$\begin{aligned}
 L(\theta_m) = & \frac{\lambda_b}{N_b} \sum_i^{N_b} (NN(x_i; \theta_m) - u_i)^2 && \text{Boundary loss} \\
 & + \frac{\lambda_p}{N_p} \sum_j^{N_p} \left(R \left(NN(x_j; \theta_m) \right) \right)^2 && \text{Physics loss} \\
 & + \sum_l^{N_n} \frac{\lambda_l}{N_l} \sum_k^{N_l} (NN(x_k; \theta_l) - NN(x_k; \theta_m))^2 && \text{Interface conditions} \\
 & l \in \text{Neighbours}(m)
 \end{aligned}$$

Jagtap, A., et al., Extended physics-informed neural networks (XPINNs): A generalized space-time domain decomposition based deep learning framework for nonlinear partial differential equations. Communications in Computational Physics (2020)

Extended PINNs (XPINNs)

XPINNs solving 2D
Poisson's equation



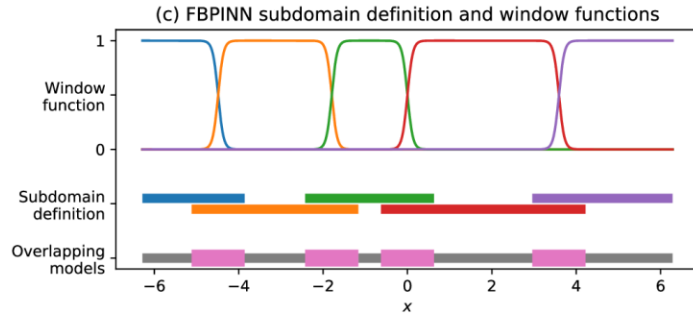
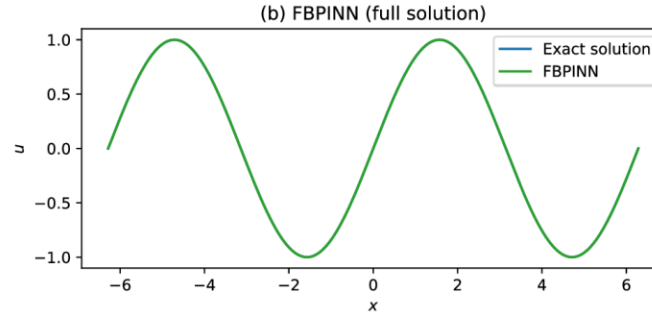
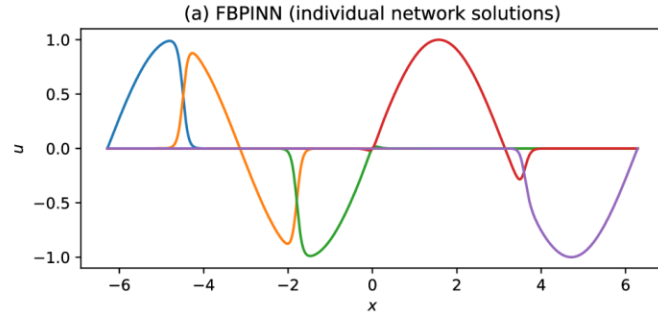
Limitations:

- Introduces discontinuities in solution at subdomain interfaces
- Requires extra loss terms

$$\begin{aligned}
 L(\theta_m) = & \frac{\lambda_b}{N_b} \sum_i^{N_b} (NN(x_i; \theta_m) - u_i)^2 && \text{Boundary loss} \\
 & + \frac{\lambda_p}{N_p} \sum_j^{N_p} \left(R \left(NN(x_j; \theta_m) \right) \right)^2 && \text{Physics loss} \\
 & + \sum_l^{N_n} \frac{\lambda_l}{N_l} \sum_k^{N_l} (NN(x_k; \theta_l) - NN(x_k; \theta_m))^2 && \text{Interface conditions} \\
 & l \in \text{Neighbours}(m)
 \end{aligned}$$

Jagtap, A., et al., Extended physics-informed neural networks (XPINNs): A generalized space-time domain decomposition based deep learning framework for nonlinear partial differential equations. *Communications in Computational Physics* (2020)

Finite basis PINNs (FBPINNs)



$$\hat{u}(x; \theta) = \mathcal{C} \left[\sum_i^n w_i(x) \cdot \text{unnorm} \circ NN_i \circ \text{norm}_i(x) \right]$$

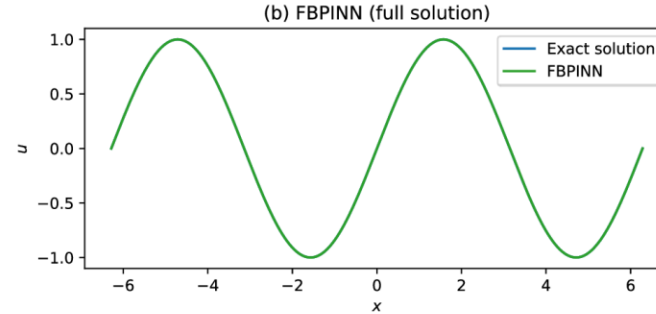
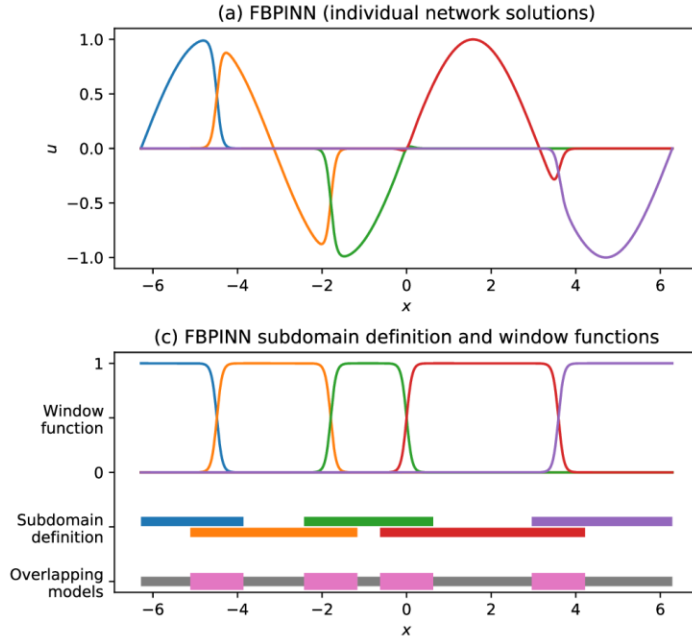
Window function
Subdomain network

Individual subdomain normalisation

Moseley et al, Finite Basis Physics-Informed Neural Networks (FBPINNs): a scalable domain decomposition approach for solving differential equations, ACM (2023)

Idea: use **overlapping** subdomains and a **globally** defined solution ansatz

Finite basis PINNs (FBPINNs)

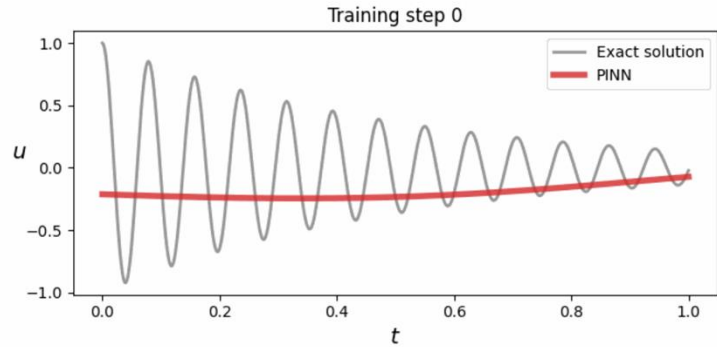


Note:

- By construction, FBPINN solution is continuous across subdomain interfaces
- Can be trained with same loss function as PINNs
- FBPINNs can simply be thought of as a custom NN architecture for PINNs

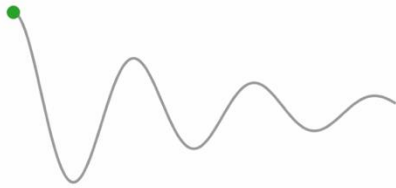
Moseley et al, Finite Basis Physics-Informed Neural Networks (FBPINNs): a scalable domain decomposition approach for solving differential equations, ACM (2023)

FBPINNs vs PINNs

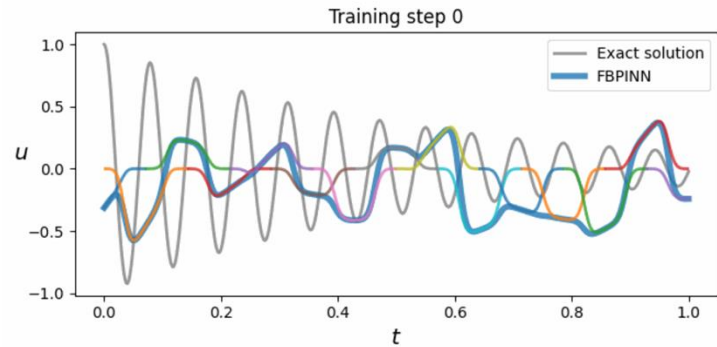


PINN solution

Network size: 2 hidden layers, 64 hidden units



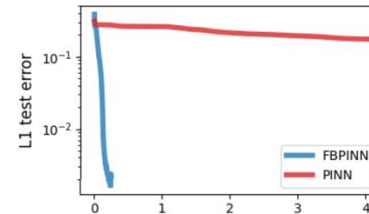
Damped harmonic oscillator



FBPINN solution

Number of subdomains: 15

Subdomain network size: 1 hidden layer, 32 hidden units



Other PINN extensions

- Different differential equations:
 - Fractional equations, stochastic PDEs, variational formulations (VPINNs), mixed formulations, ...
- Uncertainty estimation:
 - Bayesian PINNs (B-PINNs), GAN PINNs, ...
- Architectures:
 - Adaptive activation functions, extreme learning machines, SIRENs, wavelet layers, LSTMs, ...
- Other optimisation strategies:
 - NTK weighting, casual training, specialised optimisers
- PINN software libraries:
 - NVIDIA Modulus, DeepXDE, NeuroDiffEq, PyDENs, ...
- And much more!

SPRINGER LINK

Find a journal Publish with us Track your research Search

Home > Journal of Scientific Computing > Article

Scientific Machine Learning Through Physics–Informed Neural Networks: Where we are and What’s Next

Open access | Published: 26 July 2022

Volume 92, article number 88, (2022) Cite this article

Download PDF ↓

You have full access to this open access article

Salvatore Cuomo, Vincenzo Schiano Di Cola, Fabio Giampaolo, Gianluigi Rozza, Maziar Raissi & Francesco Piccialli ✉

67k Accesses 274 Citations 7 Altmetric Explore all metrics →

Abstract

Physics–Informed Neural Networks (PINN) are neural networks (NNs) that encode model equations, like Partial Differential Equations (PDE), as a component of the neural network itself. PINNs are nowadays used to solve PDEs, fractional equations, integral–differential equations, and stochastic PDEs. This novel methodology has arisen as a multi–task learning framework in which a NN must fit observed data while reducing a PDE residual. This article provides a comprehensive review of the literature on PINNs: while the primary goal of the study is to characterize these networks and their related advantages and

1. Training DNN is slow and do not scale

- Deep neural networks (DNNs) are used to solve learning tasks with “state-of-the-art accuracies”
...but the training time and cost is increasing!
- **GPT-4** contain hundreds of billions of parameters, far exceeding the scale of earlier models like GPT-3. Training a model like GPT-4 can take months on **thousands of GPUs**...
- All sophisticated deep learning architectures also have long training times and **models do not scale!**

The **future of AI** may depend on developing more **resource-efficient architectures** and **training paradigms!**

Why SciML models do not scale? ChatGPT says:

Scaling Scientific Machine Learning (SciML) models presents several challenges. These issues arise from the **unique characteristics of scientific data and models, which differ from those typically used in traditional machine learning**. Here are key reasons why SciML models may not scale well:

1. Complexity of Scientific Data and Models

- **High Dimensionality:** Scientific datasets are often high-dimensional, involving numerous variables and parameters. This can lead to the "curse of dimensionality," where the computational complexity grows exponentially with the number of variables.
- **Non-linear Dynamics:** Many scientific models, such as those in physics, chemistry, or biology, involve highly non-linear relationships, which are challenging for machine learning models to capture efficiently at larger scales.
- **Coupled Systems:** SciML models often represent coupled systems (e.g., climate models, fluid dynamics), where different subsystems interact with each other. This introduces additional layers of complexity that make it difficult to scale models across multiple systems.

2. Divide and conquer approaches tackle spectral bias

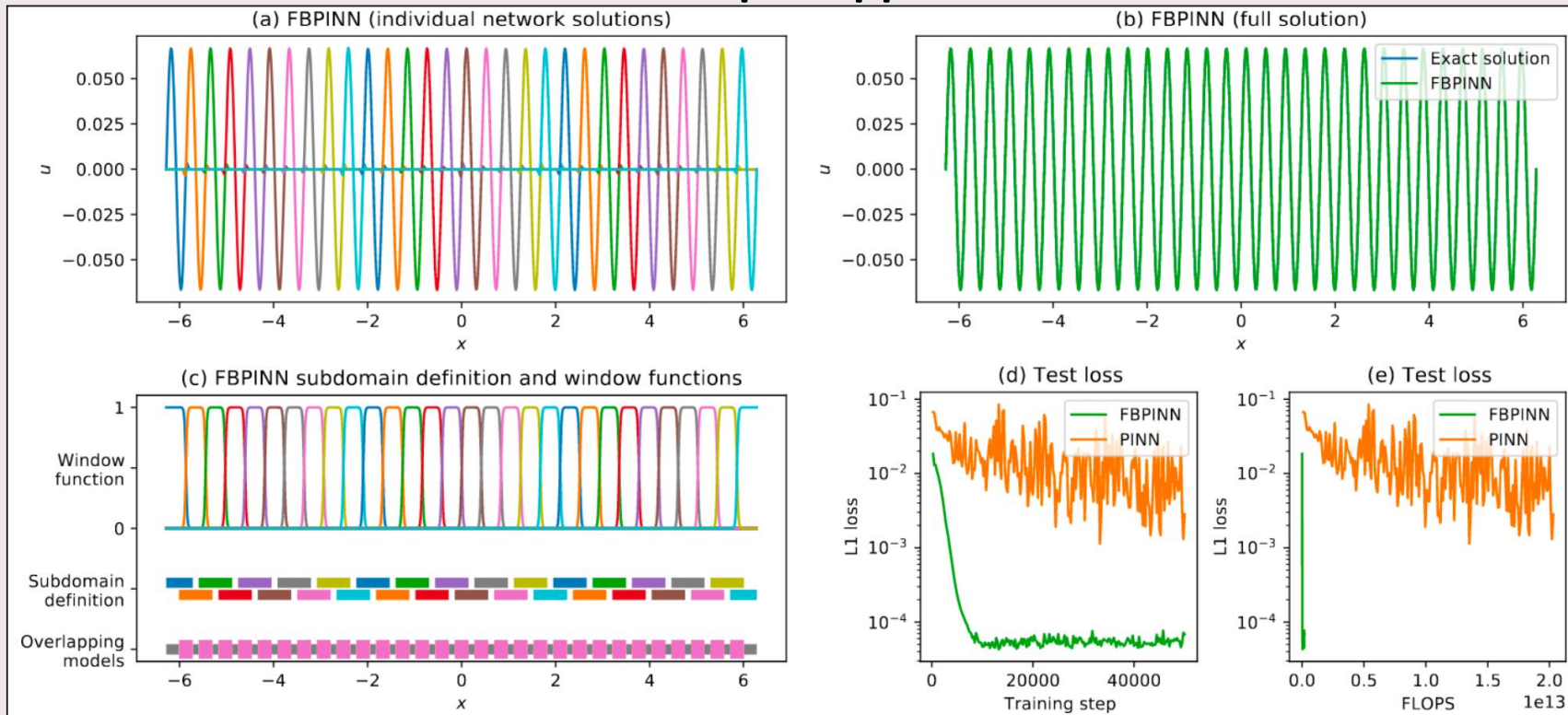
Chat GPT says: Partition of Unity Networks (POU Nets) offer a promising approach to scaling Scientific Machine Learning (SciML) models:

- **Localized learning** reduces computational complexity and improves model stability.
- **Parallel training** of local models enhances scalability.
- **Piecewise approximation** handles non-linearities more effectively.
- **Adaptive refinement** allows models to scale in specific regions where higher accuracy is needed.
- **Preservation of physical constraints** ensures model accuracy, even when scaled.

Partition of Unity Networks (POUNets) offer several mechanisms to mitigate spectral bias, especially in the context of scaling SciML models.

- **Partitioning the domain** into smaller regions,
- **Allowing localized learning** of high-frequency features.
- **Leveraging multiscale representations**, and
- **Ensuring smooth transitions across regions**

One-level divide and conquer approach based on POU nets



FBPINN by VD, B. Moseley et al.

Hu, Jagtap, Karniadakis, Kawaguchi, APINNs: A gating network-based soft DDM, 2024

3. One-level architectures are not enough

L-level network architecture

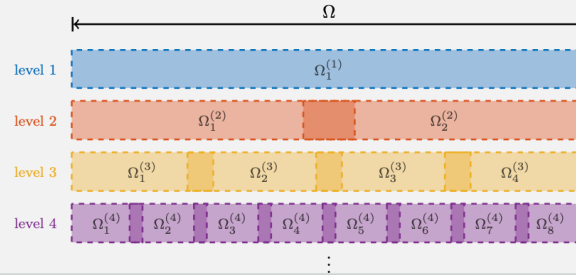
Let L overlapping domain decompositions

$\Omega = \cup_{j=1}^{J^{(l)}} \Omega_j^{(l)}$ and **window functions** $\omega_j^{(l)}$ with

$$\sum_{j=1}^{J^{(l)}} \omega_j^{(l)} \equiv 1, \text{supp}(\omega_j^{(l)}) \subset \Omega_j^{(l)}.$$

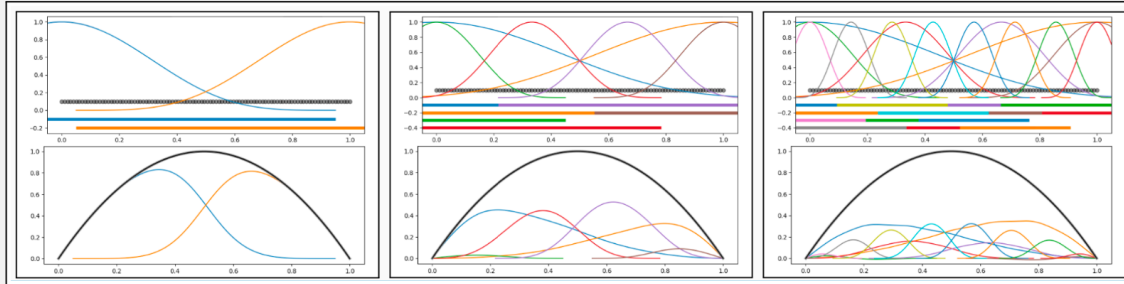
Then the multilevel architecture is

$$u(\mathbf{x}, \theta) = \mathcal{C} \sum_{l=1}^L \sum_{j=1}^{J^{(l)}} \omega_j^{(l)} u_j^{(l)}(\mathbf{x}, \theta_j^{(l)})$$



Loss function

$$\mathcal{L}(\theta) = \frac{1}{N} \sum_{i=1}^N (\mathcal{N}[\mathcal{C} \sum_{l=1}^L \sum_{x_i \in \Omega_j^{(l)}} \omega_j^{(l)} u_j^{(l)}] (\mathbf{x}_i, \theta_j^{(l)}) - f(\mathbf{x}_i))^2.$$



23

Dolean, Heinlein, Mishra, Moseley, Multilevel DDM (CMAME, 2024)

Final thoughts and conclusions

There is far more than that:

- Layer parallelism can be achieved in a similar fashion (comes in 1 and 2-levels flavours):

Kopanicacova, Kothari, Karniadakis, Krause, SISC, 2023

Lee, Kopanicakova, Karniadakis, arXiv, 2024

- Combining POU networks with ELM (extreme learning machines) and add a second level: RFM (Random Feature Method)

Chi, Chen, E, Yang, JML, 2022

Chen, E, Sun, CAMC, 2024

- (literature growing at a fast pace as we speak)

Conclusions

Home > Computational Science and Engineering > Article

Machine learning and domain decomposition methods – a survey

Research | Open access | Published: 23 September 2024
Volume 1, article number 2, (2024) [Cite this article](#)

[Download PDF](#)  You have full access to this [open access](#) article

Axel Klawonn , Martin Lanser  & Janine Weber 

 791 Accesses [Explore all metrics](#) →

Abstract

Hybrid algorithms, which combine black-box machine learning methods with experience from traditional numerical methods and domain expertise from diverse application areas, are progressively gaining importance in scientific machine learning and various industrial domains, especially in computational science and engineering. In the present survey, several promising avenues of research will be examined which focus on the combination of machine learning (ML) and domain decomposition methods (DDMs). The aim of this survey is to provide an overview of existing work within this field and to structure it into

A few take aways:

- Data parallelism is not enough
- Clever algorithms require some change of architecture

We use a partition of unity (PoU) to piece together different local representations as well as a global representation for the large- scale components of the solution. This strategy has proven to be vital in practice in order to achieve good accuracy. (Chi, Chen, E, Yang, 2022)

- POU networks **are akin to divide and conquer approaches** (like domain decomposition) -> DD+ML is also now a research field!
- Old HPC and DD wisdom says that **mere splitting won't guarantee scalability** (e.g. w.r.t frequency) or robustness (in multiscale contexts) -> **we need more levels...**

Acknowledgements



Ben Moseley (Imperial College London):

<https://benmoseley.blog>

GitHub: <https://github.com/benmoseley/>

YouTube channel (course on Deep Learning for Scientific Computing)

<https://www.youtube.com/@CAMLabETHZurich>