

AI in the Sciences and Engineering

Introduction to Hybrid Workflows – Part 1

Spring Semester 2024

Siddhartha Mishra
Ben Moseley

ETH zürich

Course timeline

Tutorials

Mon 12:15-14:00 HG E 5

- 19.02.
- 26.02. Introduction to PyTorch
- 04.03. Simple DNNs in PyTorch
- 11.03. Implementing PINNs I
- 18.03. Implementing PINNs II
- 25.03. Operator learning I
- 01.04.
- 08.04. Operator learning II
- 15.04.
- 22.04. GNNs
- 29.04. Transformers
- 06.05. Diffusion models
- 13.05. Coding autodiff from scratch
- 20.05.
- 27.05. Intro to JAX / Neural ODEs

Lectures

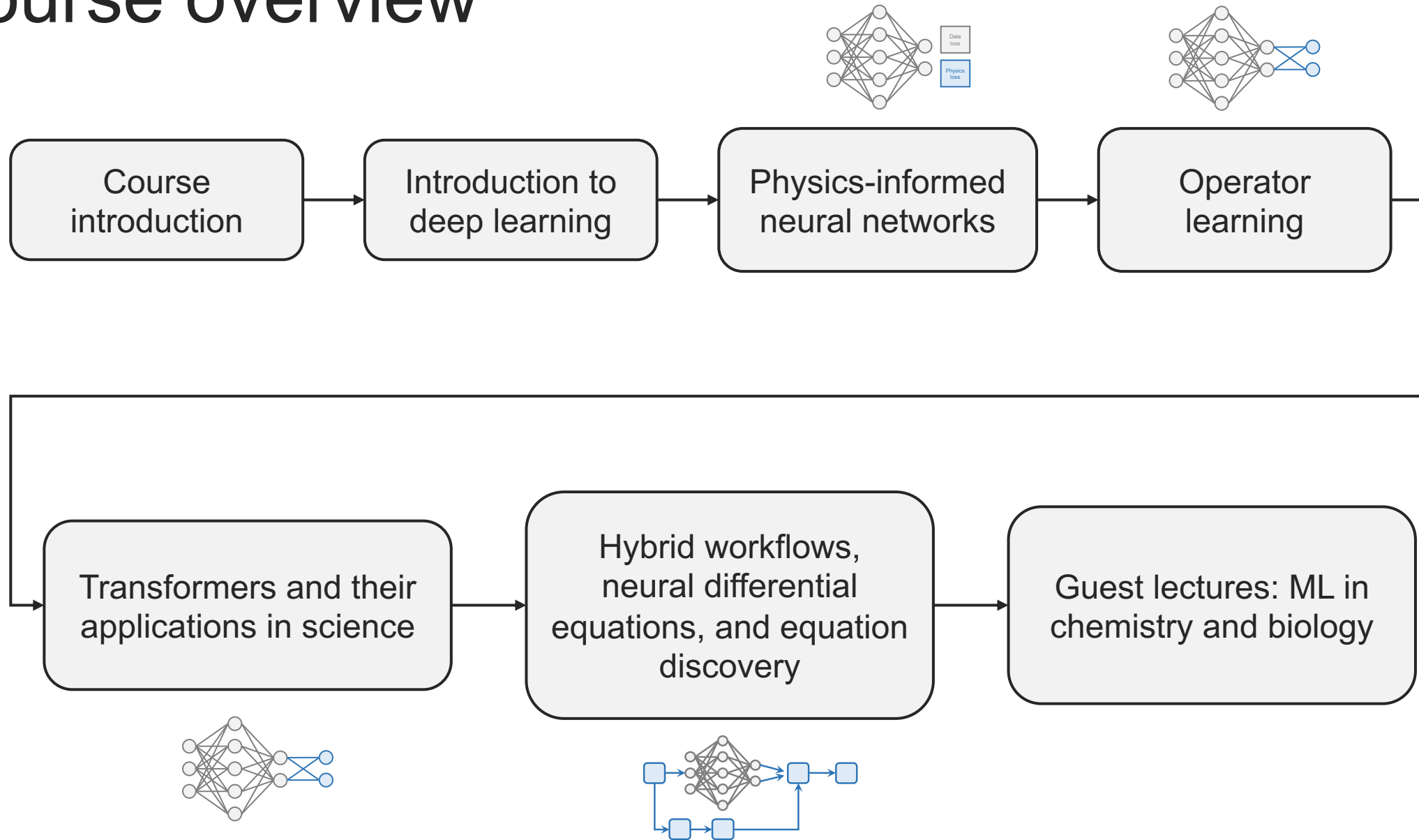
Wed 08:15-10:00 ML H 44

- 21.02. Course introduction
- 28.02. Introduction to deep learning II
- 06.03. Physics-informed neural networks – introduction
- 13.03. Physics-informed neural networks – extensions
- 20.03. Physics-informed neural networks – theory II
- 27.03. Supervised learning for PDEs II
- 03.04.
- 10.04. Introduction to operator learning I
- 17.04. Convolutional neural operators
- 24.04. Large-scale neural operators
- 01.05.
- 08.05. **Introduction to hybrid workflows I**
- 15.05. Neural differential equations
- 22.05. Symbolic regression and model discovery
- 29.05. Guest lecture: AlphaFold

Fri 12:15-13:00 ML H 44

- 23.02. Introduction to deep learning I
- 01.03. Introduction to PDEs
- 08.03. Physics-informed neural networks - limitations
- 15.03. Physics-informed neural networks – theory I
- 22.03. Supervised learning for PDEs I
- 29.03.
- 05.04.
- 12.04. Introduction to operator learning II
- 19.04. Time-dependent neural operators
- 26.04. Attention as a neural operator
- 03.05. Windowed attention and scaling laws
- 10.05. Introduction to hybrid workflows II
- 17.05. Introduction to JAX
- 24.05. Course summary
- 31.05. Guest lecture: AlphaFold

Course overview



Lecture overview

- Limitations of SciML approaches studied so far
- Hybrid SciML approaches
 - Residual modelling
 - Opening the “black-box”
 - How to train hybrid approaches
- Autodifferentiation
 - Autodifferentiation as a key enabler
 - What it is and how it works

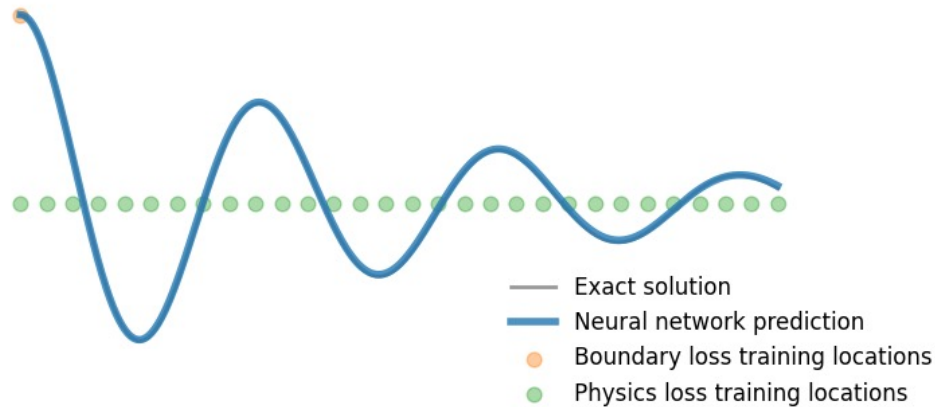
Lecture overview

- Limitations of SciML approaches studied so far
- Hybrid SciML approaches
 - Residual modelling
 - Opening the “black-box”
 - How to train hybrid approaches
- Autodifferentiation
 - Autodifferentiation as a key enabler
 - What it is and how it works

Learning objectives

- Be able to describe what a hybrid workflow is
- Understand how autodifferentiation is used to train hybrid workflows
- Understand how autodifferentiation works

Course recap - PINNs

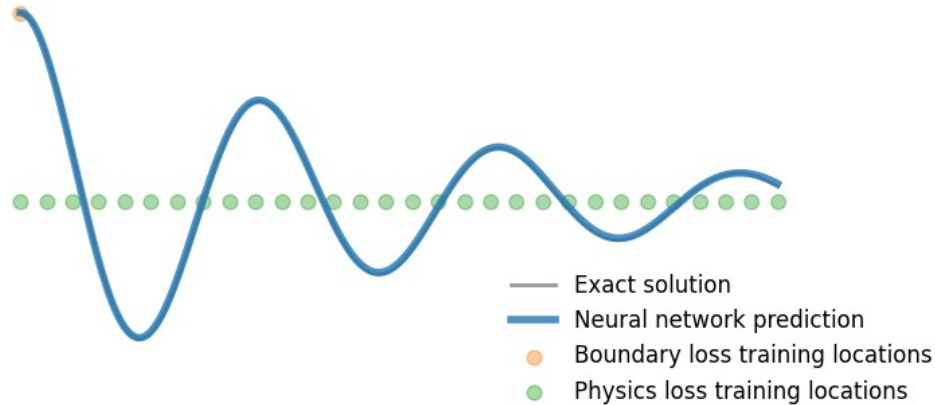


$$m \frac{d^2 u}{dt^2} + \mu \frac{du}{dt} + ku = 0$$

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

$$\begin{aligned}
 \text{Boundary loss } L_b(\theta) & \left\{ \begin{aligned}
 L(\theta) &= \lambda_1 (NN(\underline{t=0}; \theta) - \underline{1})^2 \\
 &+ \lambda_2 \left(\frac{dNN}{dt}(\underline{t=0}; \theta) - \underline{0} \right)^2
 \end{aligned} \right. \\
 \text{Physics loss } L_p(\theta) & \left\{ \begin{aligned}
 &+ \frac{1}{N_p} \sum_i^{N_p} \left(\left[m \frac{d^2}{dt^2} + \mu \frac{d}{dt} + k \right] NN(\underline{t_i}; \theta) \right)^2
 \end{aligned} \right.
 \end{aligned}$$

Course recap - PINNs



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

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

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

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

Advantages of PINNs

- **Mesh-free**
- Can jointly solve forward and inverse problems
- Often performs well on “**messy**” problems (where some observational data is available)
- Mostly **unsupervised**
- Can perform well for high-dimensional PDEs

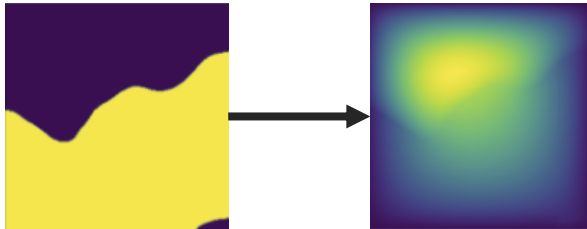
Limitations of PINNs

- **Computational cost** often high (especially for forward-only problems)
- Can be hard to **optimise**
- Challenging to **scale** to high-frequency, multi-scale problems

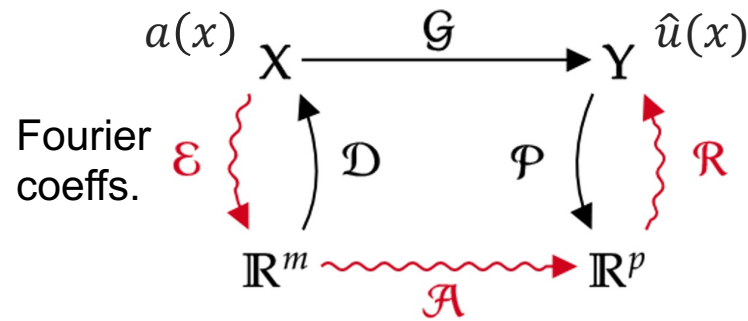
Course recap – Operator learning

Darcy PDE

$$\nabla \cdot (a(x)\nabla u(x)) = f(x)$$



Permeability, $a(x)$ Pressure, $u(x)$



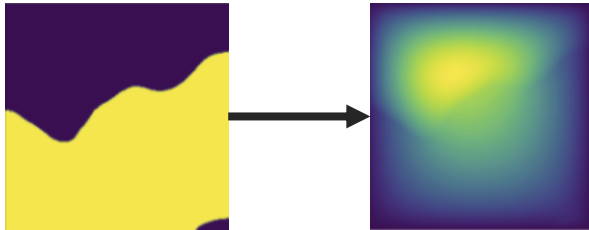
$$a(x) \rightarrow \underbrace{\{a_k\}_{k=1}^m \rightarrow NN(\{a_k\}; \theta) \rightarrow \{u_k\}_{k=1}^p}_{\mathcal{G}_\theta^*[a]} \rightarrow \hat{u}(x)$$

$$L(\theta) = \frac{1}{NM} \sum_i^N \sum_j^M \|u_i(x_j) - \mathcal{G}_\theta^*[a_i](x_j)\|^2$$

Course recap – Operator learning

Darcy PDE

$$\nabla \cdot (a(x)\nabla u(x)) = f(x)$$



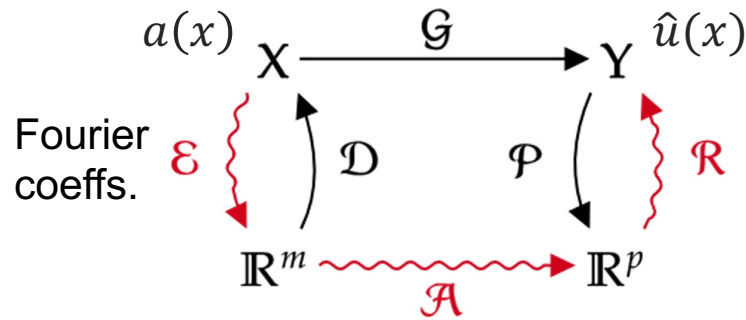
Permeability, $a(x)$ Pressure, $u(x)$

Advantages of operator learning

- Can be **orders of magnitude faster** than traditional simulation (once trained)

Limitations of operator learning

- Can require **lots** of training data, which can be expensive to obtain
- Can struggle to **generalise** to inputs outside of its training data
- Encoding / reconstruction steps require some **assumptions** about the regularity of $a(x)$ and $u(x)$



$$a(x) \rightarrow \{a_k\}_{k=1}^m \xrightarrow{\mathcal{G}_\theta^*[a]} \{u_k\}_{k=1}^p \rightarrow \hat{u}(x)$$

Fourier interpolation

$$L(\theta) = \frac{1}{NM} \sum_i^N \sum_j^M \|u_i(x_j) - \mathcal{G}_\theta^*[a_i](x_j)\|^2$$

When should I use deep neural networks for scientific problems?

Advantages of DNNs

- Usually very **fast** (once trained)
- Can represent highly **non-linear** functions

Limitations of DNNs

- Often lots of **training data required**
- Can be hard to **optimise**
- Can be hard to **interpret**
- Often struggle to **generalise**

When should I use deep neural networks for scientific problems?

Advantages of DNNs

- Usually very **fast** (once trained)
- Can represent highly **non-linear** functions

Limitations of DNNs

- Often lots of **training data required**
- Can be hard to **optimise**
- Can be hard to **interpret**
- Often struggle to **generalise**

General advice

Use DNNs to:

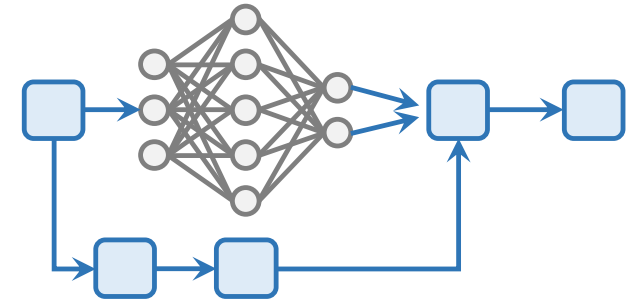
- 1) **Accelerate** your workflow, or
- 2) Learn the **parts** you are unsure of / have incomplete knowledge

Entirely replacing your existing workflow with a DNN may **not** be a good idea!

Hybrid SciML approaches



What if we **directly incorporate** DNNs into a traditional algorithm instead?
= **hybrid approach**



General advice

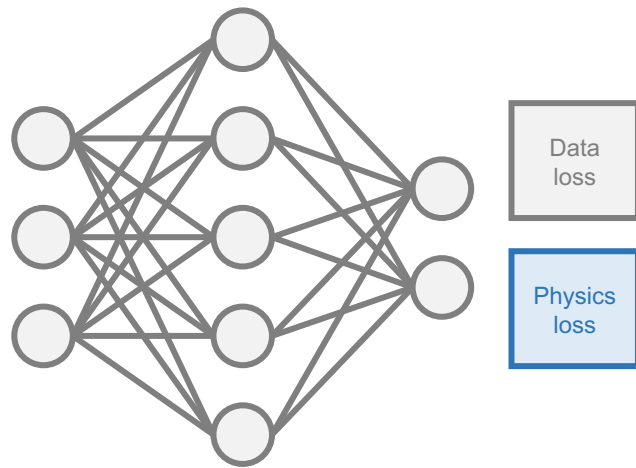
Use DNNs to:

- 1) **Accelerate** your workflow, or
- 2) Learn the **parts** you are unsure of / have incomplete knowledge

Entirely replacing your existing workflow with a DNN may **not** be a good idea!

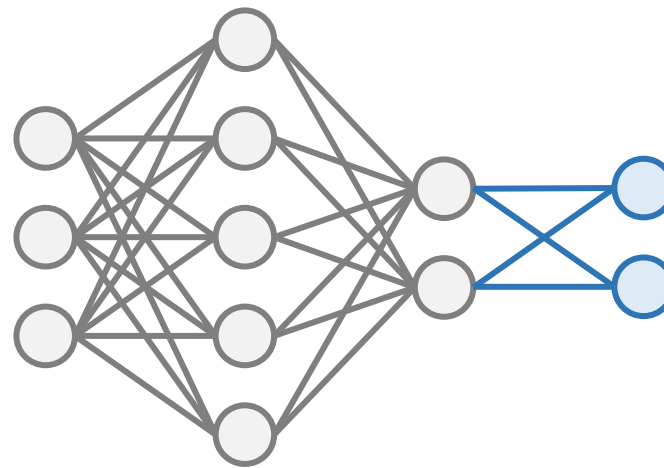
Ways to incorporate scientific principles into machine learning

Loss function



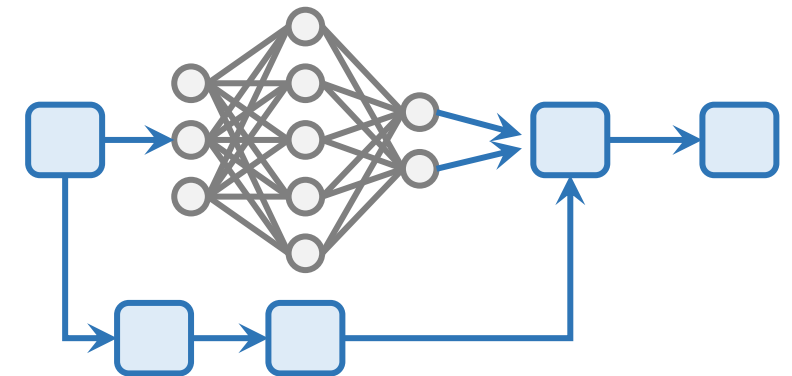
Example:
Physics-informed neural networks
(add governing equations to loss
function)

Architecture



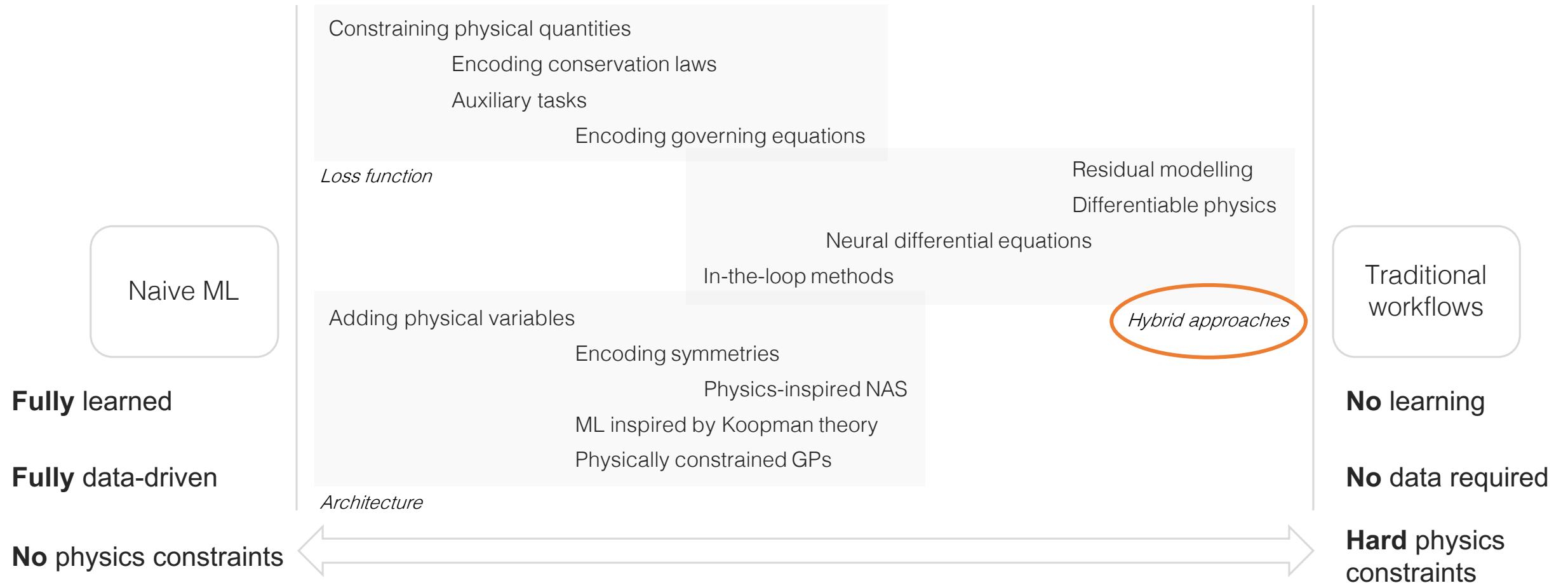
Example:
Encoding symmetries / conservation laws
(e.g. energy conservation, rotational
invariance), operator learning

Hybrid approaches



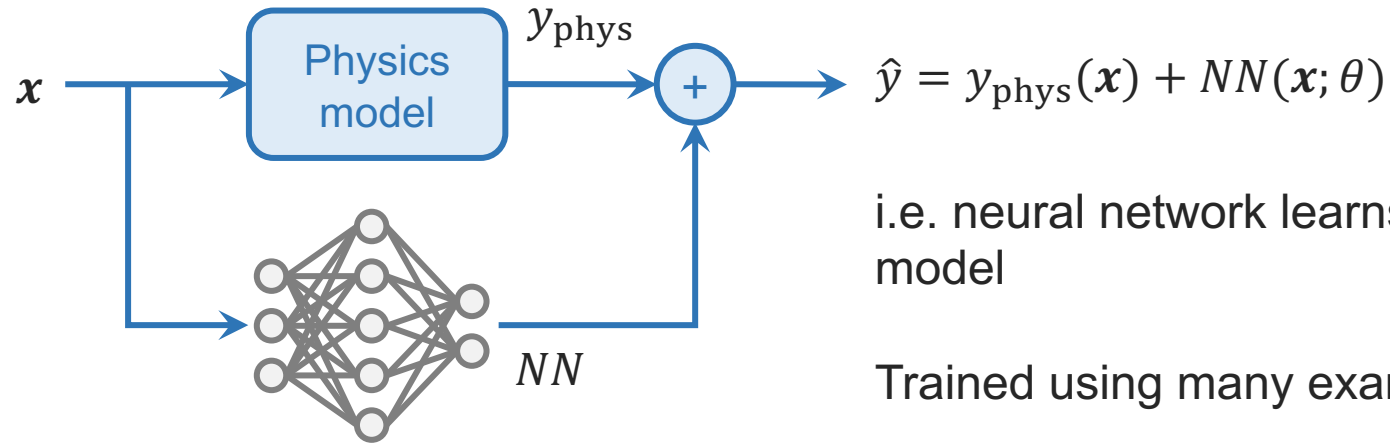
Example:
Neural differential equations
(incorporating neural networks into PDE
models)

A plethora of SciML techniques



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

A simple hybrid approach – residual modelling

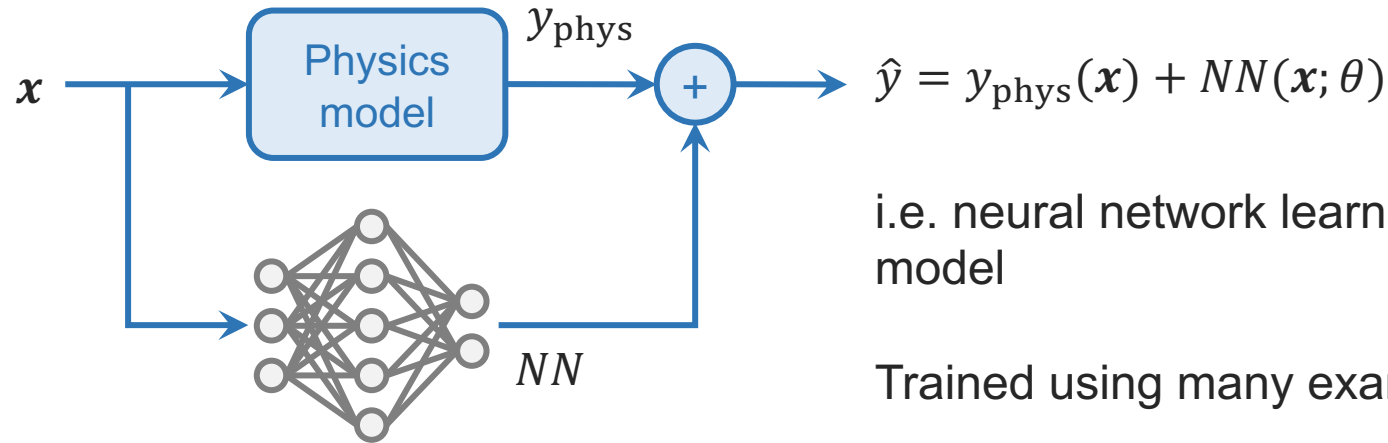


i.e. neural network learns **residual** correction to physics model

Trained using many examples of inputs/outputs

When is this useful?

A simple hybrid approach – residual modelling



i.e. neural network learns **residual** correction to physics model

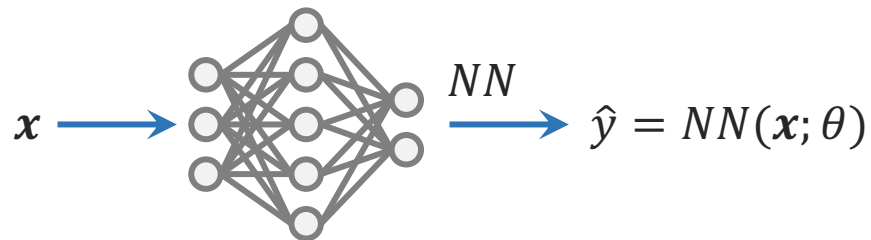
Trained using many examples of inputs/outputs

Useful when:

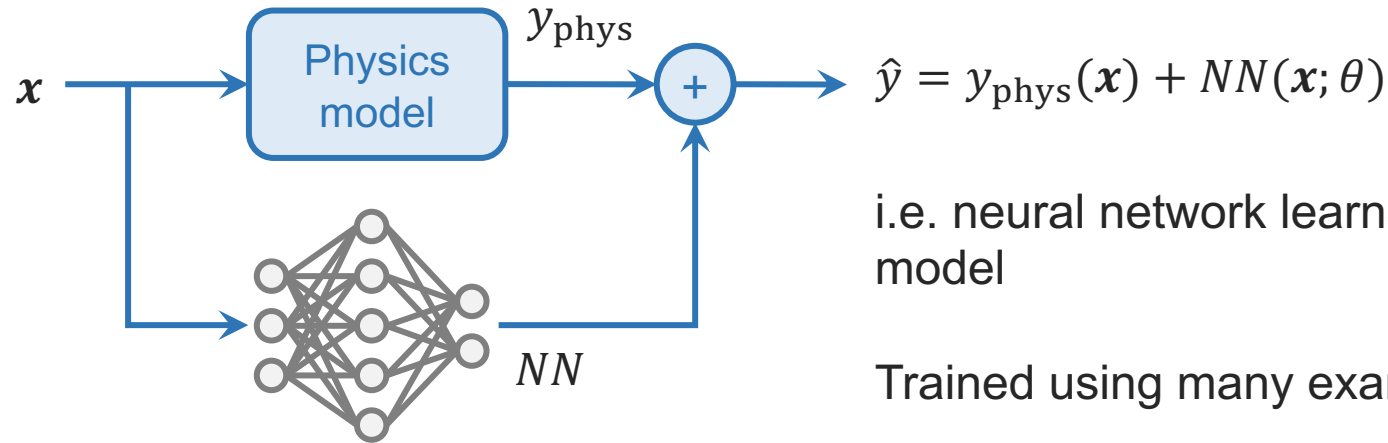
- We have incomplete understanding of physics
- More complex physical modeling is too expensive

Compared to naïve ML approach:

- **Easier** learning task: don't need to learn all the physics
- More **interpretable**



A simple hybrid approach – residual modelling



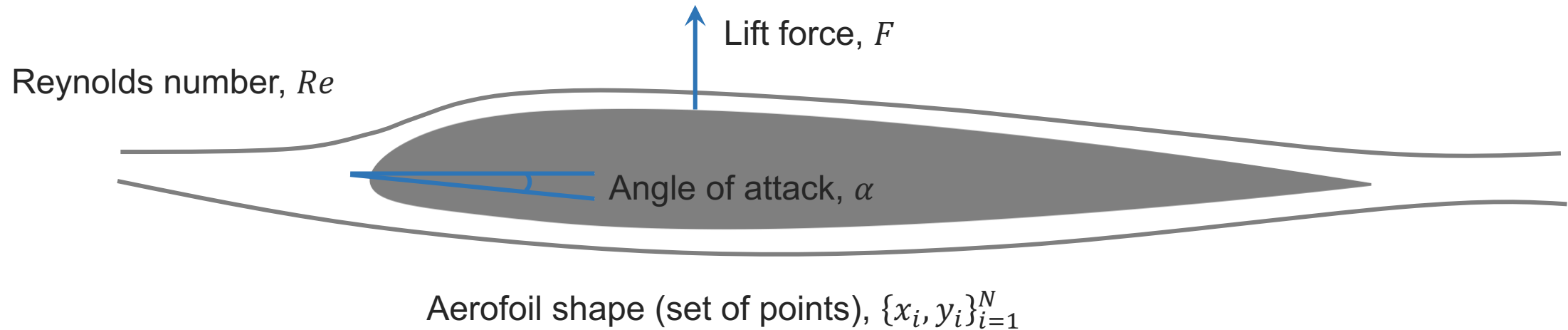
i.e. neural network learns **residual** correction to physics model

Trained using many examples of inputs/outputs

$$\begin{aligned} L(\theta) &= \sum_{i=1}^N (\hat{y}(\mathbf{x}_i; \theta) - y_i)^2 \\ &= \sum_{i=1}^N (NN(\mathbf{x}_i; \theta) - [y_i - y_{\text{phys}}(\mathbf{x}_i)])^2 \\ &\equiv \sum_{i=1}^N (NN(\mathbf{x}_i; \theta) - r(\mathbf{x}_i))^2 \end{aligned}$$

Note: can precompute $r(\mathbf{x}_i)$ in advance

Residual modelling – aerofoil example



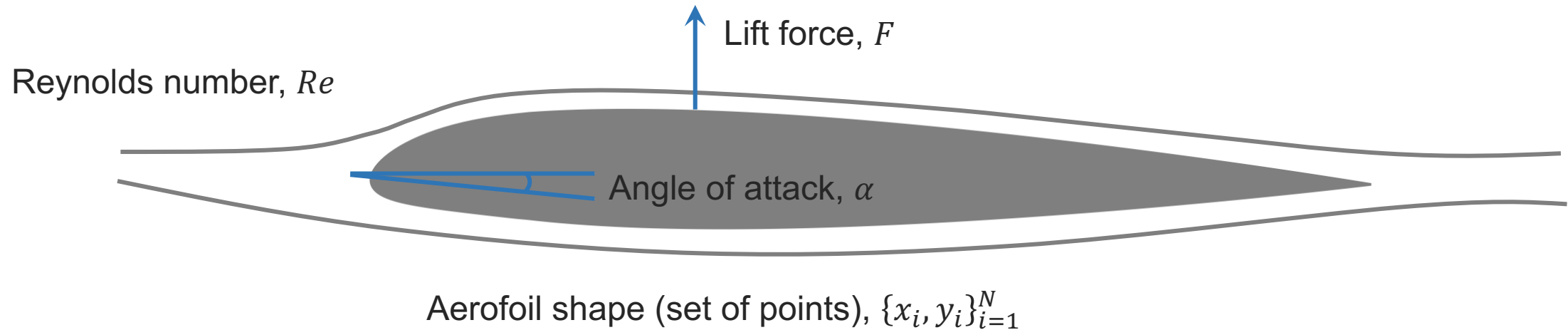
Simulation task:

Given $\{x_i, y_i\}_{i=1}^N$, Re and α

Predict F

Pawar et al, Physics guided machine learning using simplified theories, Physics of Fluids (2021)

Residual modelling – aerofoil example



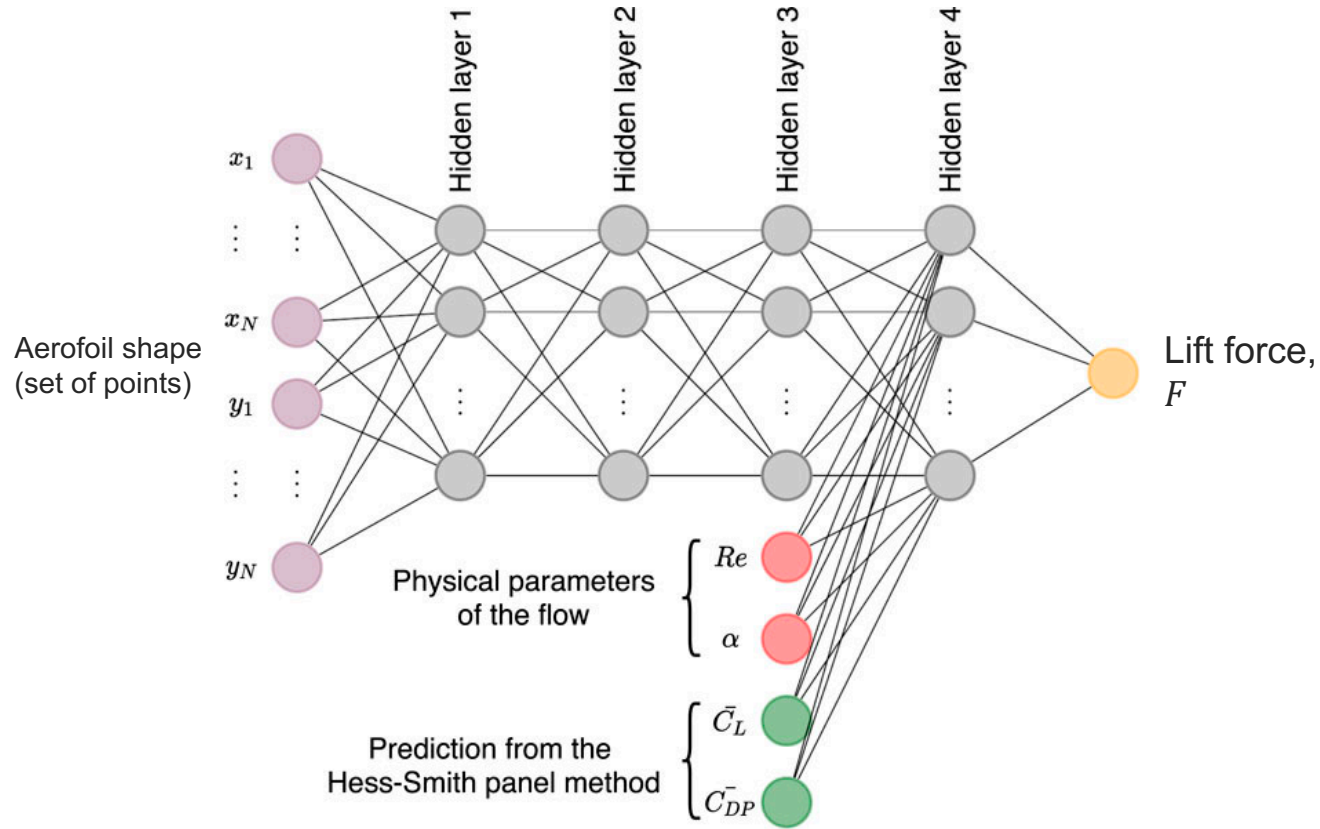
- Full CFD simulations are typically accurate, but very expensive
- Faster approximate methods exist, but are usually less accurate

Simulation task:

Given $\{x_i, y_i\}_{i=1}^N$, Re and α
Predict F

Pawar et al, Physics guided machine learning using simplified theories, Physics of Fluids (2021)

Residual modelling – aerofoil example



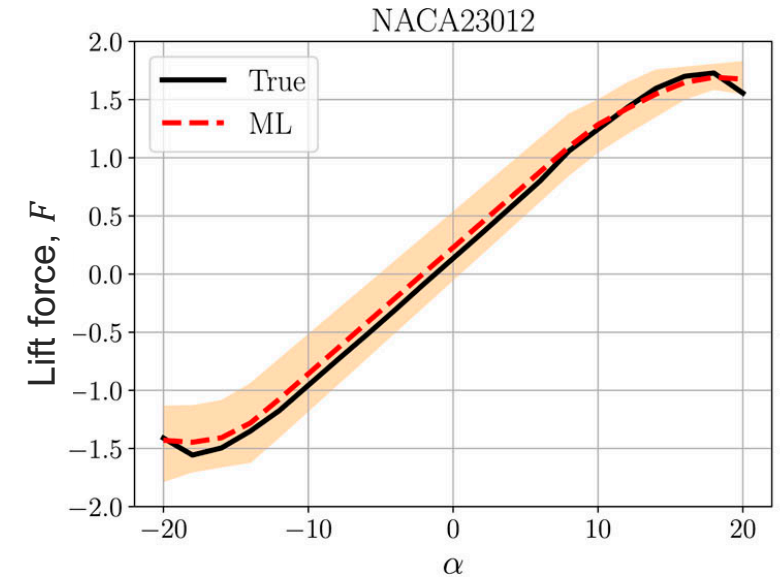
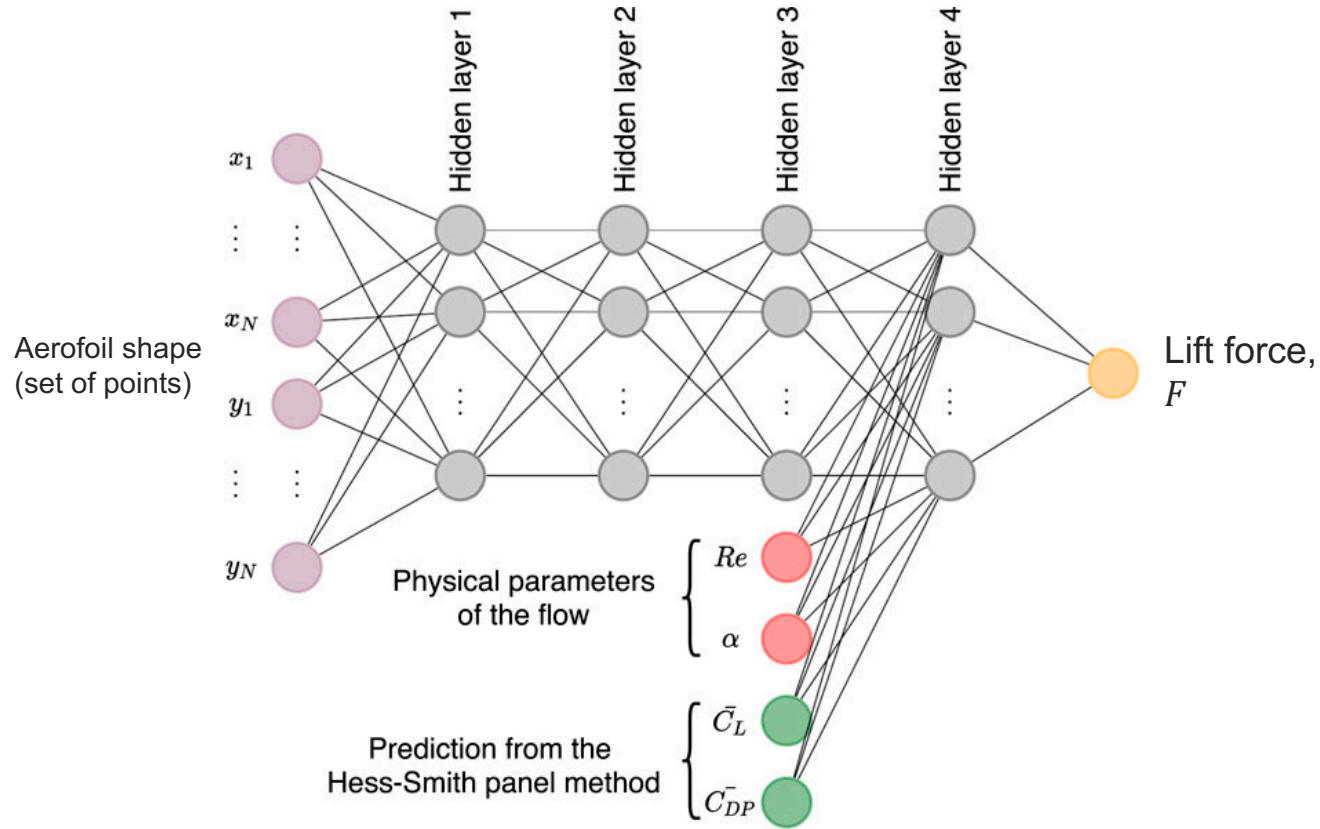
Hess-Smith panel method:
Fast **approximate** method for predicting lift force

Training data:
Many example inputs/outputs generated from (expensive) **high-fidelity** CFD modelling

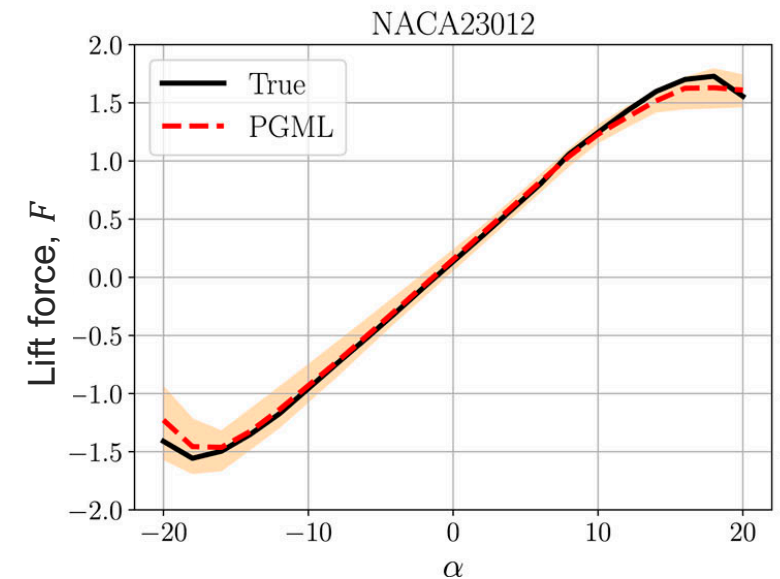
Goal:
A model which is **faster** than CFD and more **accurate** than approximate physics model

Pawar et al, Physics guided machine learning using simplified theories, Physics of Fluids (2021)

Residual modelling – aerofoil example



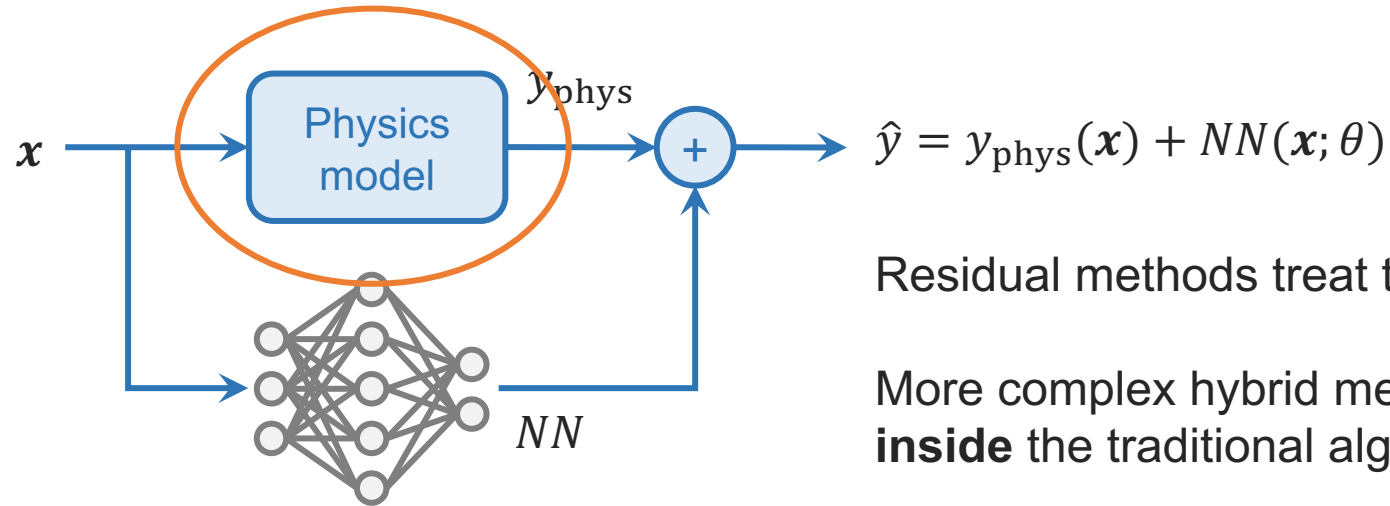
Naive NN
(no physics inputs)



NN + physics model
(hybrid approach)

Pawar et al, Physics guided machine learning using simplified theories, Physics of Fluids (2021)

Opening the black-box



Residual methods treat the physics model as a “**black-box**”

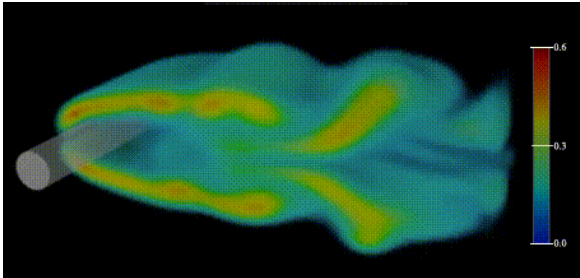
More complex hybrid methods open the box and insert ML **inside** the traditional algorithm

We insert ML where;

- 1) the algorithm is **slow**
- 2) we are **unsure** of our assumptions/ want to improve our modelling

Opening the black-box – finite difference solver

FD solver



Incompressible Navier-Stokes equation

$$\frac{\partial \mathbf{u}}{\partial t} + (\mathbf{u} \cdot \nabla) \mathbf{u} - \nu \nabla^2 \mathbf{u} = -\frac{1}{\rho} \nabla p$$
$$\nabla \cdot \mathbf{u} = 0$$

$\mathbf{u}(\mathbf{x}, t)$ is the flow velocity

$p(\mathbf{x}, t)$ is the pressure

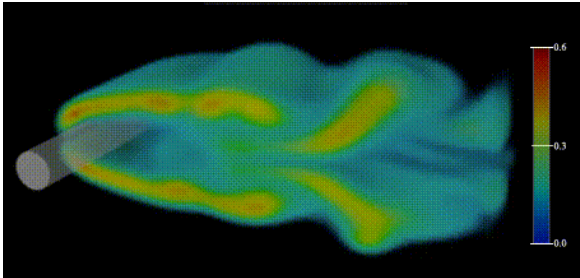
$\rho(\mathbf{x})$ is the density

ν is the viscosity

Um et al, Solver-in-the-loop: Learning from differentiable physics to interact with iterative PDE-solvers, NeurIPS (2020)

Opening the black-box – finite difference solver

FD solver



Incompressible Navier-Stokes equation

$$\frac{\partial \mathbf{u}}{\partial t} + (\mathbf{u} \cdot \nabla) \mathbf{u} - \nu \nabla^2 \mathbf{u} = -\frac{1}{\rho} \nabla p$$
$$\nabla \cdot \mathbf{u} = 0$$

$\mathbf{u}(\mathbf{x}, t)$ is the flow velocity
 $p(\mathbf{x}, t)$ is the pressure
 $\rho(\mathbf{x})$ is the density
 ν is the viscosity

“Operator splitting” numerical solver:

Discretise in time

$$\mathbf{u}_{t+1} = \mathbf{u}_t - \delta t (\mathbf{u}_t \cdot \nabla) \mathbf{u}_t + \delta t \nu \nabla^2 \mathbf{u}_t - \frac{\delta t}{\rho} \nabla p_{t+1} \quad (1)$$

Let

$$\mathbf{u}^* = \mathbf{u}_t - \delta t (\mathbf{u}_t \cdot \nabla) \mathbf{u}_t + \delta t \nu \nabla^2 \mathbf{u}_t - \frac{\delta t}{\rho} \nabla p_t \quad (2)$$

Then

$$\mathbf{u}_{t+1} = \mathbf{u}^* - \frac{\delta t}{\rho} \nabla (p_{t+1} - p_t)$$

Asserting $\nabla \cdot \mathbf{u}_{t+1} = 0 \Rightarrow$

$$0 = \nabla \cdot \mathbf{u}^* - \frac{\delta t}{\rho} \nabla^2 (p_{t+1} - p_t)$$

$$\nabla^2 (p_{t+1} - p_t) = \frac{\rho}{\delta t} \nabla \cdot \mathbf{u}^*$$

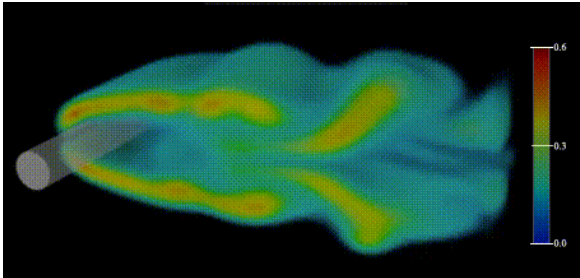
Discretise in space

$$L(p_{i,j,k,t+1} - p_{i,j,k,t}) = \frac{\rho_{i,j,k}}{\delta t} D \mathbf{u}_{i,j,k}^* \quad (3)$$

Um et al, Solver-in-the-loop: Learning from differentiable physics to interact with iterative PDE-solvers, NeurIPS (2020)

Opening the black-box – finite difference solver

FD solver



Incompressible Navier-Stokes equation

$$\frac{\partial \mathbf{u}}{\partial t} + (\mathbf{u} \cdot \nabla) \mathbf{u} - \nu \nabla^2 \mathbf{u} = -\frac{1}{\rho} \nabla p$$

$$\nabla \cdot \mathbf{u} = 0$$

$\mathbf{u}(x, t)$ is the flow velocity
 $p(x, t)$ is the pressure
 $\rho(x)$ is the density
 ν is the viscosity

“Operator splitting” numerical solver:

Discretise in time

$$\mathbf{u}_{t+1} = \mathbf{u}_t - \delta t (\mathbf{u}_t \cdot \nabla) \mathbf{u}_t + \delta t \nu \nabla^2 \mathbf{u}_t - \frac{\delta t}{\rho} \nabla p_{t+1} \quad (1)$$

Let

$$\mathbf{u}^* = \mathbf{u}_t - \delta t (\mathbf{u}_t \cdot \nabla) \mathbf{u}_t + \delta t \nu \nabla^2 \mathbf{u}_t - \frac{\delta t}{\rho} \nabla p_t \quad (2)$$

Then

$$\mathbf{u}_{t+1} = \mathbf{u}^* - \frac{\delta t}{\rho} \nabla (p_{t+1} - p_t)$$

Asserting $\nabla \cdot \mathbf{u}_{t+1} = 0 \Rightarrow$

$$0 = \nabla \cdot \mathbf{u}^* - \frac{\delta t}{\rho} \nabla^2 (p_{t+1} - p_t)$$

$$\nabla^2 (p_{t+1} - p_t) = \frac{\rho}{\delta t} \nabla \cdot \mathbf{u}^*$$

Discretise in space

$$L(p_{i,j,k,t+1} - p_{i,j,k,t}) = \frac{\rho_{i,j,k}}{\delta t} D \mathbf{u}_{i,j,k}^* \quad (3)$$

Basic algorithm:

Discretise \mathbf{u}, p and ρ

Loop:

1. Compute $\mathbf{u}_{i,j,k}^*$ using (2)
2. Solve matrix equation (3) for $p_{i,j,k,t+1}$
3. Compute $\mathbf{u}_{i,j,k,t+1}$ using (1)

```
def NS_solver(u_0, p_0, rho, nu):
    "Pseudocode for solving NS equation"

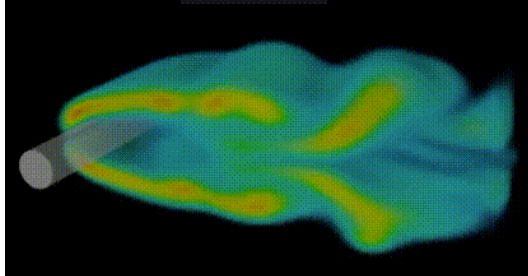
    # u_0, p_0 have shape (NX, NY, NZ)
    u_t, p_t = u_0, p_0
    for t in range(0, T):
        u_star = f(u_t, p_t, rho, nu)
        p_t = matrix_solve(u_star, p_t, rho)
        u_t = g(u_t, p_t, rho, nu)

    return u_t, p_t
```

Um et al, Solver-in-the-loop: Learning from differentiable physics to interact with iterative PDE-solvers, NeurIPS (2020)

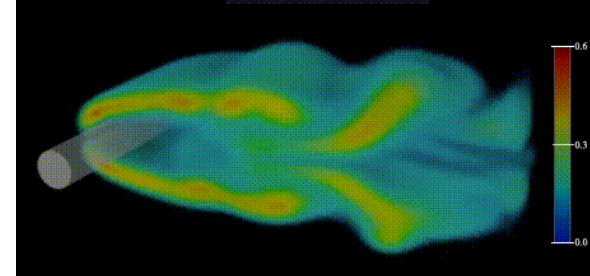
Computational cost / accuracy trade-off

Low fidelity FD solver



(32 x 32 x 64) cells
~10 seconds / 100 timesteps

High fidelity FD solver



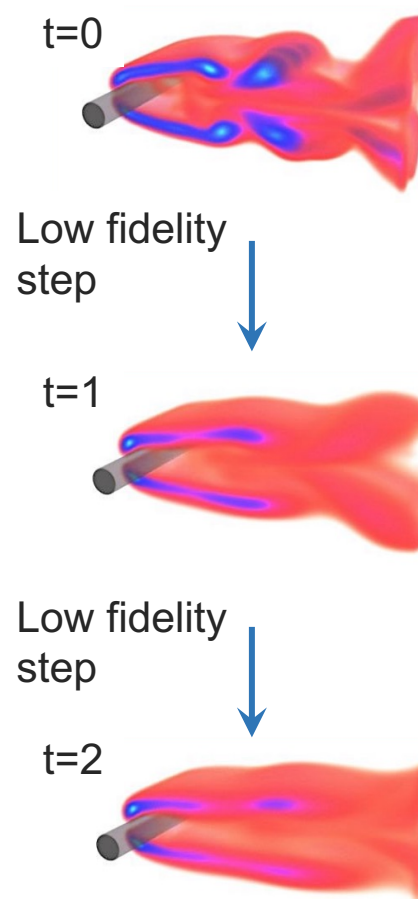
(128 x 128 x 256) cells
~1000 seconds / 100 timesteps

- Discretisation induces **errors** in the solver
- But finer grids are much more computationally expensive
- Can we use ML improve the accuracy of the **low fidelity** solver?

Um et al, Solver-in-the-loop: Learning from differentiable physics to interact with iterative PDE-solvers, NeurIPS (2020)

Traditional Navier-Stokes solver

```
def NS_solver(u_0, p_0, rho, nu):  
    "Pseudocode for solving NS equation"  
  
    # u_0, p_0 have shape (NX, NY, NZ)  
    u_t, p_t = u_0, p_0  
    for t in range(0, T):  
        u_star = f(u_t, p_t, rho, nu)  
        p_t = matrix_solve(u_star, p_t, rho)  
        u_t = g(u_t, p_t, rho, nu)  
  
    return u_t, p_t
```

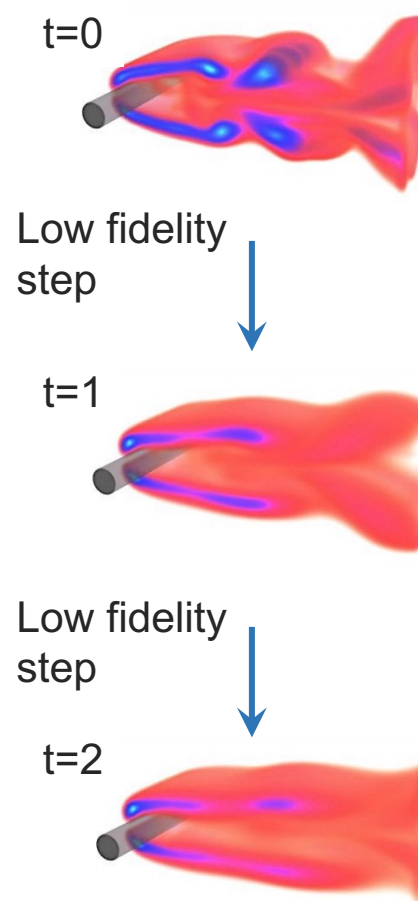


Um et al, Solver-in-the-loop: Learning from differentiable physics to interact with iterative PDE-solvers, NeurIPS (2020)

Traditional Navier-Stokes solver

```
def NS_solver(u_0, p_0, rho, nu):  
    "Pseudocode for solving NS equation"  
  
    # u_0, p_0 have shape (NX, NY, NZ)  
    u_t, p_t = u_0, p_0  
    for t in range(0, T):  
        u_star = f(u_t, p_t, rho, nu)  
        p_t = matrix_solve(u_star, p_t, rho)  
        u_t = g(u_t, p_t, rho, nu)  
  
    return u_t, p_t
```

- Where could we insert ML **inside** this workflow to improve accuracy / efficiency?

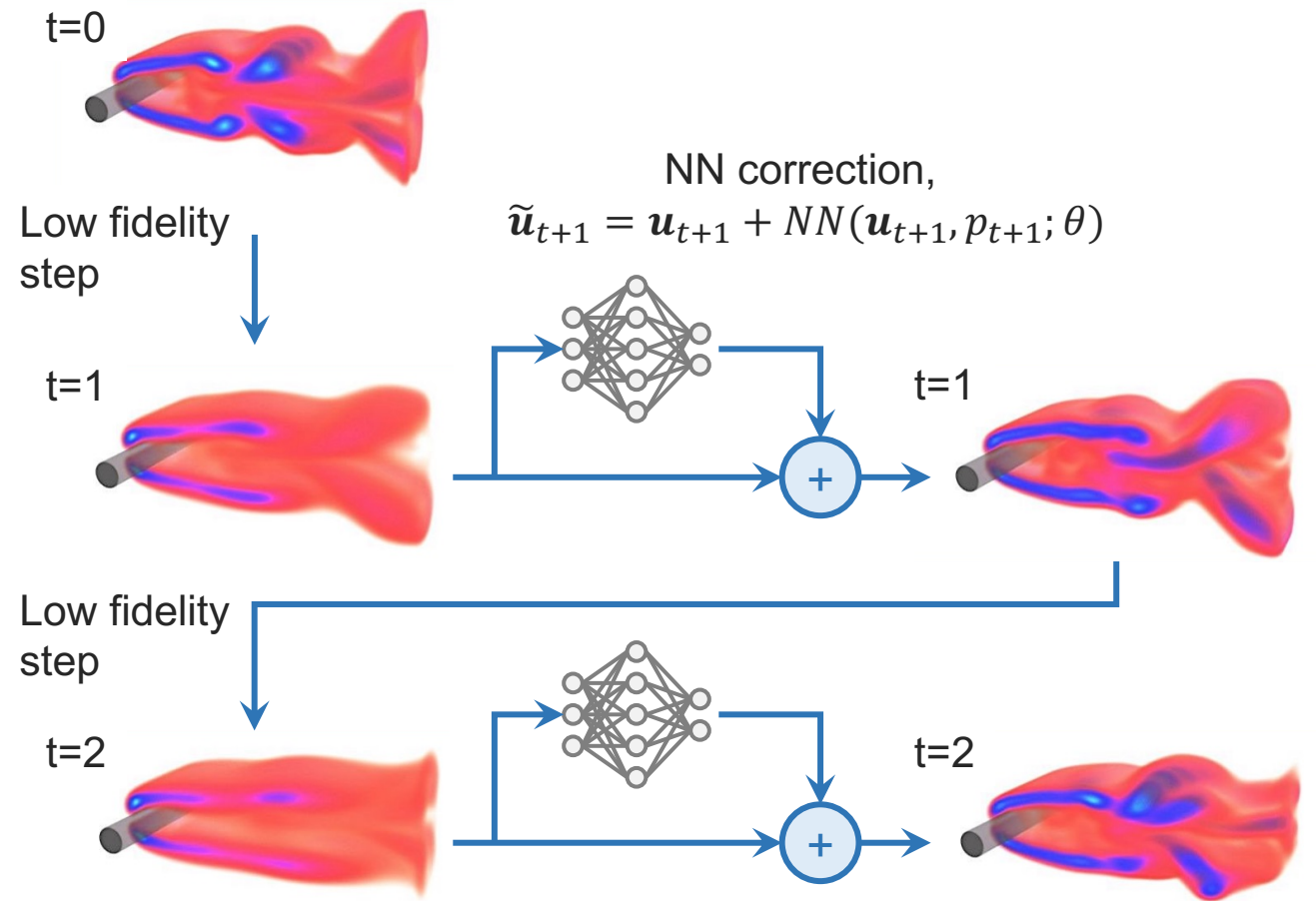


Um et al, Solver-in-the-loop: Learning from differentiable physics to interact with iterative PDE-solvers, NeurIPS (2020)

Hybrid Navier-Stokes solver

```
def NS_solver(u_0, p_0, rho, nu):  
    "Pseudocode for solving NS equation"  
  
    # u_0, p_0 have shape (NX, NY, NZ)  
    u_t, p_t = u_0, p_0  
    for t in range(0, T):  
        u_star = f(u_t, p_t, rho, nu)  
        p_t = matrix_solve(u_star, p_t, rho)  
        u_t = g(u_t, p_t, rho, nu)  
  
    return u_t, p_t
```

```
def Hybrid_NS_solver(u_0, p_0, rho, nu, theta):  
    "Pseudocode for solving NS equation, with NN correction"  
  
    # u_0, p_0 have shape (NX, NY, NZ)  
    u_t, p_t = u_0, p_0  
    for t in range(0, T):  
        u_star = f(u_t, p_t, rho, nu)  
        p_t = matrix_solve(u_star, p_t, rho)  
        u_t = g(u_t, p_t, rho, nu)  
  
        u_t, p_t = (u_t, p_t) + NN(u_t, p_t, theta)  
  
    return u_t, p_t
```

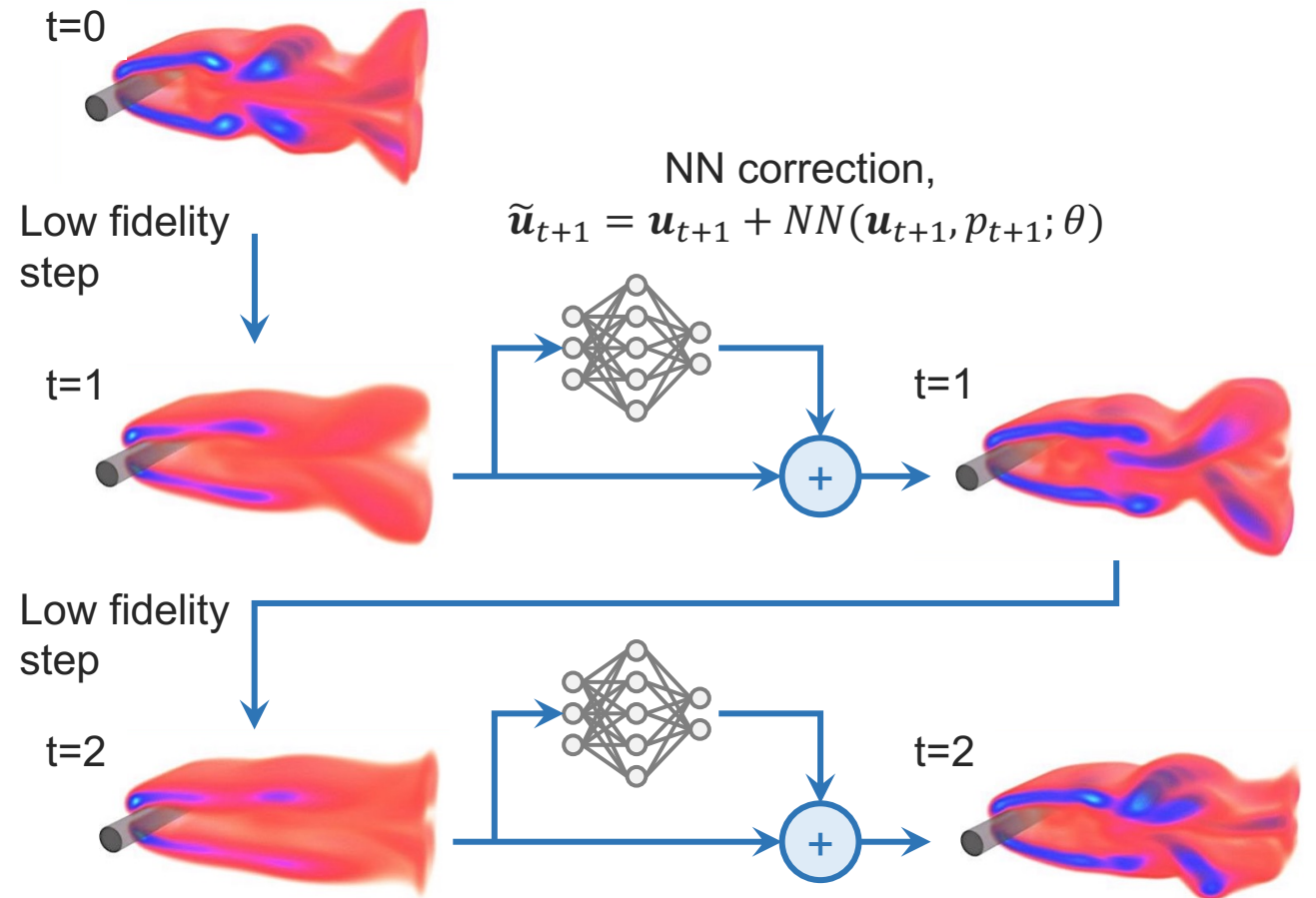


Um et al, Solver-in-the-loop: Learning from differentiable physics to interact with iterative PDE-solvers, NeurIPS (2020)

Hybrid Navier-Stokes solver

- How can we train $NN(\mathbf{u}_{t+1}, p_{t+1}; \theta)$?
- What training data do we need? (Hint: what inputs/labels do we need to train the network?)
- What loss function should we use?

```
def Hybrid_NS_solver(u_0, p_0, rho, nu, theta):  
    "Pseudocode for solving NS equation, with NN correction"  
  
    # u_0, p_0 have shape (NX, NY, NZ)  
    u_t, p_t = u_0, p_0  
    for t in range(0, T):  
        u_star = f(u_t, p_t, rho, nu)  
        p_t = matrix_solve(u_star, p_t, rho)  
        u_t = g(u_t, p_t, rho, nu)  
  
        u_t, p_t = (u_t, p_t) + NN(u_t, p_t, theta)  
  
    return u_t, p_t
```

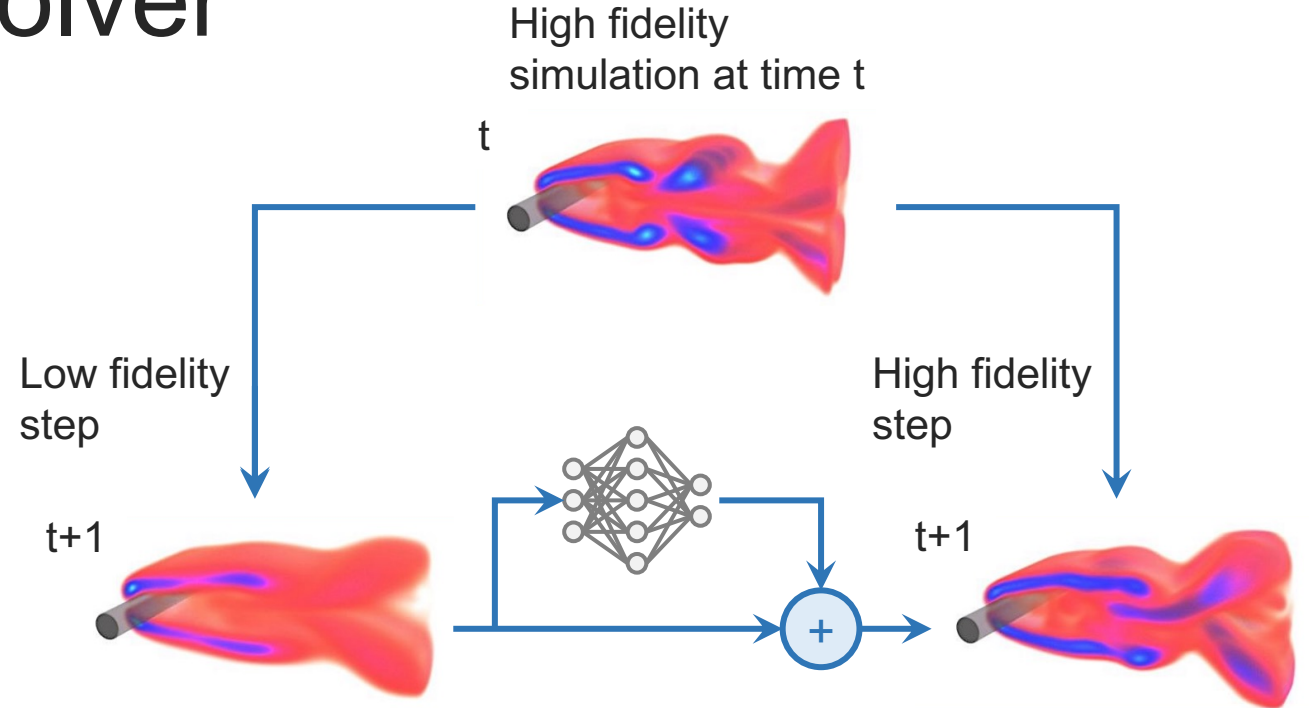


Um et al, Solver-in-the-loop: Learning from differentiable physics to interact with iterative PDE-solvers, NeurIPS (2020)

Hybrid Navier-Stokes solver

- How can we train $NN(\mathbf{u}_{t+1}, p_{t+1}; \theta)$?

Option 1: use pairs of low fidelity / high fidelity timesteps as training data



$$L(\theta) = \sum_t^N \|\mathbf{u}_{t+1}^L + NN(\mathbf{u}_{t+1}^L, p_{t+1}^L; \theta) - \mathbf{u}_{t+1}^H\|^2$$
$$= \sum_t^N \|NN(\mathbf{u}_{t+1}^L, p_{t+1}^L; \theta) - (\mathbf{u}_{t+1}^H - \mathbf{u}_{t+1}^L)\|^2$$

Note: can precompute residual in advance

Um et al, Solver-in-the-loop: Learning from differentiable physics to interact with iterative PDE-solvers, NeurIPS (2020)

Hybrid Navier-Stokes solver

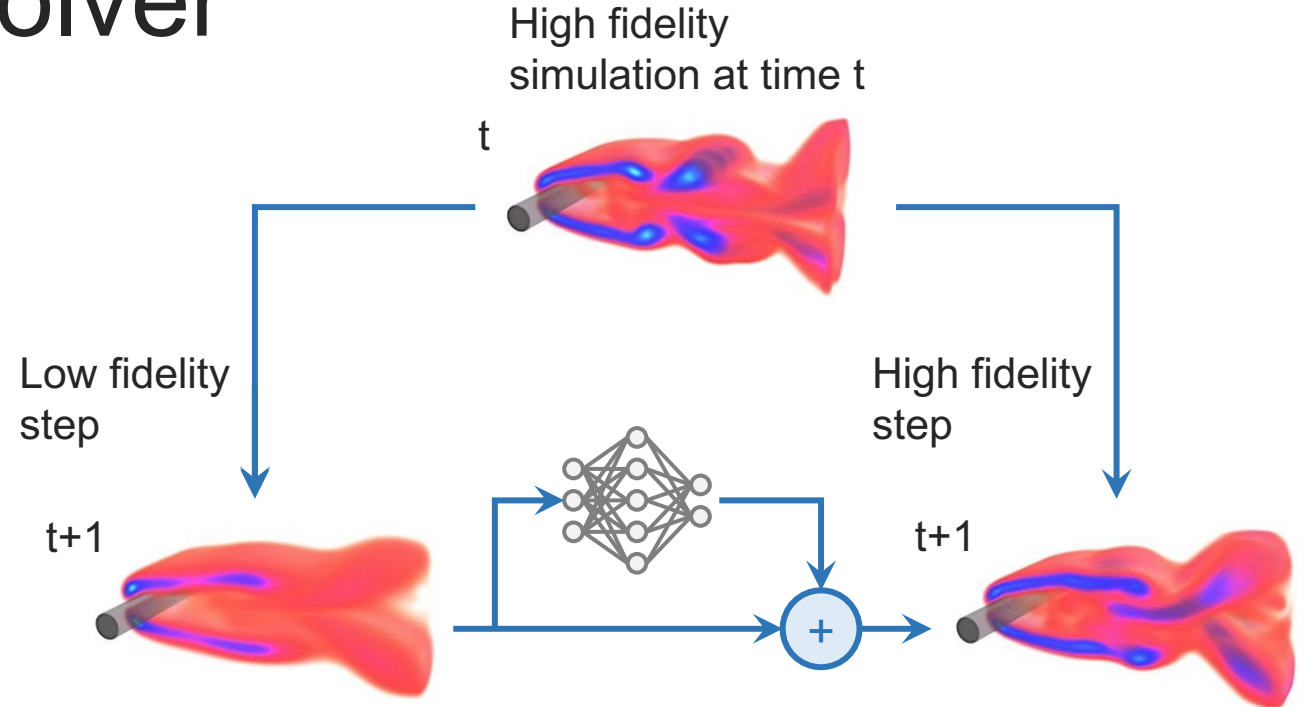
- How can we train $NN(\mathbf{u}_{t+1}, p_{t+1}; \theta)$?

Option 1: use pairs of low fidelity / high fidelity timesteps as training data

During training, neural network only sees **exact** low fidelity timesteps as input

Problem:

- But during inference, neural network sees **different** inputs (low fidelity timesteps + previous NN corrections) ⚠
- Leads to a train/test distribution shift, and error accumulation over time



$$L(\theta) = \sum_t^N \|\mathbf{u}_{t+1}^L + NN(\mathbf{u}_{t+1}^L, p_{t+1}^L; \theta) - \mathbf{u}_{t+1}^H\|^2$$

$$= \sum_t^N \|NN(\mathbf{u}_{t+1}^L, p_{t+1}^L; \theta) - (\mathbf{u}_{t+1}^H - \mathbf{u}_{t+1}^L)\|^2$$

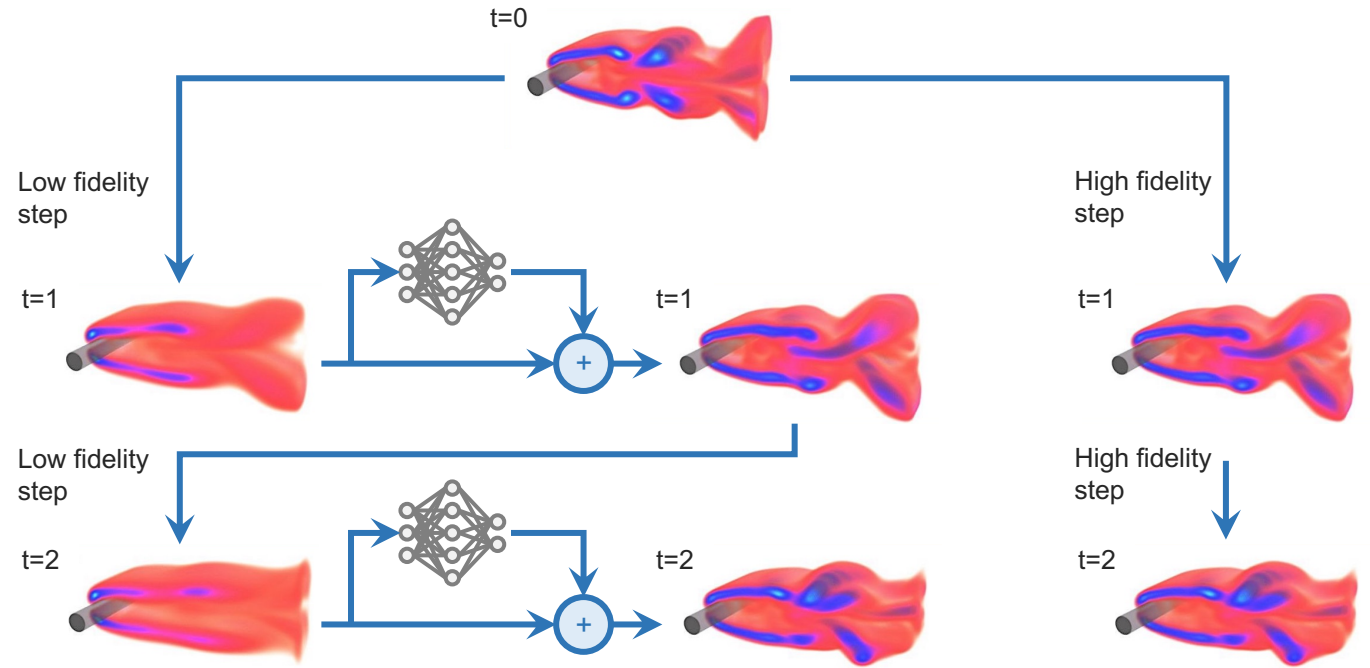
Note: can precompute residual in advance

Um et al, Solver-in-the-loop: Learning from differentiable physics to interact with iterative PDE-solvers, NeurIPS (2020)

Hybrid Navier-Stokes solver

- How can we train $NN(\mathbf{u}_{t+1}, p_{t+1}; \theta)$?

Option 2: match outputs of hybrid solver to high-fidelity simulation directly



$$L(\theta) = \sum_i^N \sum_t^T \|\text{HybridSolver}_t(\mathbf{u}_{0_i}; \theta) - \mathbf{u}_t^H(\mathbf{u}_{0_i})\|^2$$

Um et al, Solver-in-the-loop: Learning from differentiable physics to interact with iterative PDE-solvers, NeurIPS (2020)

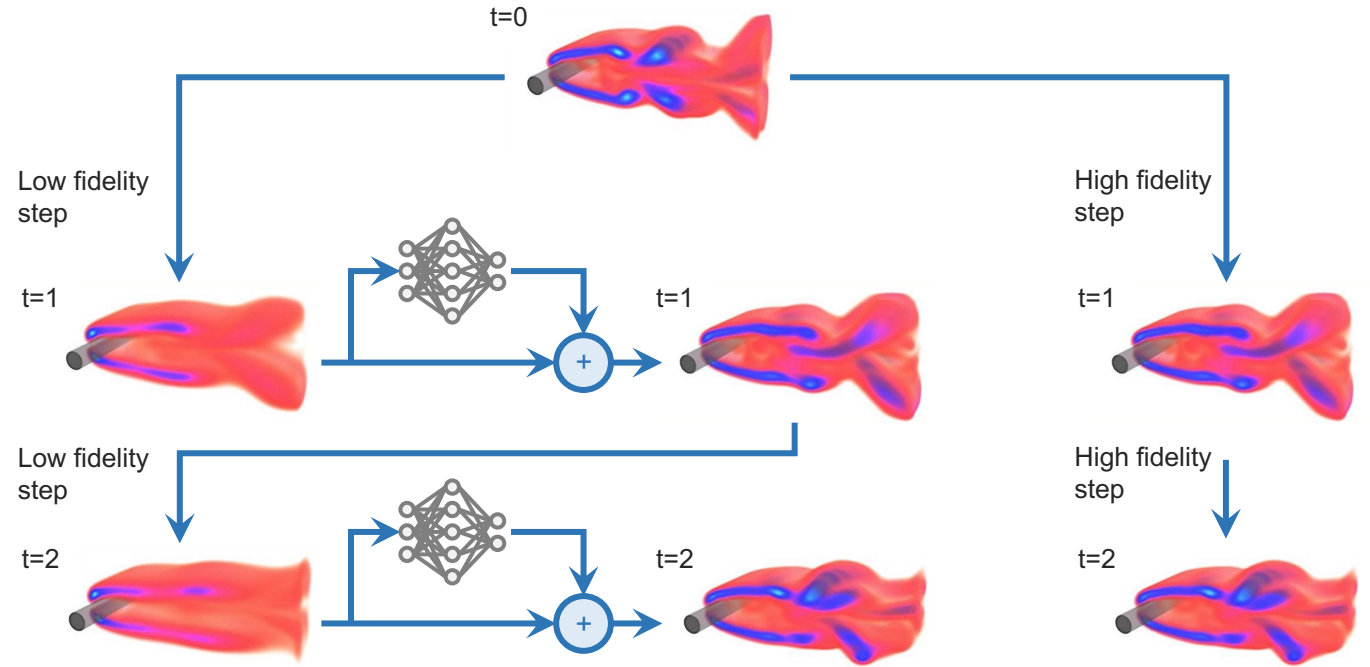
Hybrid Navier-Stokes solver

- How can we train $NN(\mathbf{u}_{t+1}, p_{t+1}; \theta)$?

Option 2: match outputs of hybrid solver to high-fidelity simulation directly

NN learns to correct its previous errors ✓
Reduces distributional shift ✓

Requires *HybridSolver* to be **differentiable!**



$$L(\theta) = \sum_i^N \sum_t^T \left\| \text{HybridSolver}_t(\mathbf{u}_{0_i}; \theta) - \mathbf{u}_t^H(\mathbf{u}_{0_i}) \right\|^2$$

Um et al, Solver-in-the-loop: Learning from differentiable physics to interact with iterative PDE-solvers, NeurIPS (2020)

Hybrid Navier-Stokes solver

- How can we train $NN(\mathbf{u}_{t+1}, p_{t+1}; \theta)$?

Option 2: match outputs of hybrid solver to high-fidelity simulation directly

NN learns to correct its previous errors ✓
Reduces distributional shift ✓

Requires *HybridSolver* to be **differentiable!**

.. we can just use autodifferentiation!

```
def Hybrid_NS_solver(u_0, p_0, rho, nu, theta):
    "Pseudocode for solving NS equation, with NN correction"

    # u_0, p_0 have shape (NX, NY, NZ)
    u_t, p_t = u_0, p_0
    for t in range(0, T):
        u_star = f(u_t, p_t, rho, nu)
        p_t = matrix_solve(u_star, p_t, rho)
        u_t = g(u_t, p_t, rho, nu)

        u_t, p_t = (u_t, p_t) + NN(u_t, p_t, theta)

    return u_t, p_t

theta.requires_grad_(True)
u_T, _ = Hybrid_NS_solver(u_0, p_0, rho, nu, theta)
loss = loss_fn(u_T, u_T_true)
dtheta = torch.autograd.grad(loss, theta)
# for learning theta (training NN)
```

$$L(\theta) = \sum_i^N \sum_t^T \left\| \text{HybridSolver}_t(\mathbf{u}_{0_i}; \theta) - \mathbf{u}_t^H(\mathbf{u}_{0_i}) \right\|^2$$

Um et al, Solver-in-the-loop: Learning from differentiable physics to interact with iterative PDE-solvers, NeurIPS (2020)

Hybrid Navier-Stokes solver

```
def NN(x, theta):  
    "Defines a FCN"  
    y = torch.tanh(theta[0]*x + theta[1])  
    return y
```

```
theta.requires_grad_(True)  
y = NN(x, theta)  
loss = loss_fn(y, y_true)  
dtheta = torch.autograd.grad(loss, theta)  
# for learning theta (training NN)
```

```
def Hybrid_NS_solver(u_0, p_0, rho, nu, theta):  
    "Pseudocode for solving NS equation, with NN correction"  
  
    # u_0, p_0 have shape (NX, NY, NZ)  
    u_t, p_t = u_0, p_0  
    for t in range(0, T):  
        u_star = f(u_t, p_t, rho, nu)  
        p_t = matrix_solve(u_star, p_t, rho)  
        u_t = g(u_t, p_t, rho, nu)  
  
        u_t, p_t = (u_t, p_t) + NN(u_t, p_t, theta)  
  
    return u_t, p_t
```

```
theta.requires_grad_(True)  
u_T,_ = Hybrid_NS_solver(u_0, p_0, rho, nu, theta)  
loss = loss_fn(u_T, u_T_true)  
dtheta = torch.autograd.grad(loss, theta)  
# for learning theta (training NN)
```

$$L(\theta) = \sum_i^N \sum_t^T \|HybridSolver_t(\mathbf{u}_{0_i}; \theta) - \mathbf{u}_t^H(\mathbf{u}_{0_i})\|^2$$

Um et al, Solver-in-the-loop: Learning from differentiable physics to interact with iterative PDE-solvers, NeurIPS (2020)

How do we train hybrid approaches?



Key idea: **autodifferentiation** allows us to differentiate and learn **arbitrary** algorithms, not just neural networks!

We train neural networks using **autodifferentiation**

But autodifferentiation = exact gradients of **arbitrary** programs

So, we can use it to differentiate (and learn) traditional algorithms too!

How do we train hybrid approaches?



Key idea: **autodifferentiation** allows us to differentiate and learn **arbitrary** algorithms, not just neural networks!

We train neural networks using **autodifferentiation**

But autodifferentiation = exact gradients of **arbitrary** programs

So, we can use it to differentiate (and learn) traditional algorithms too!



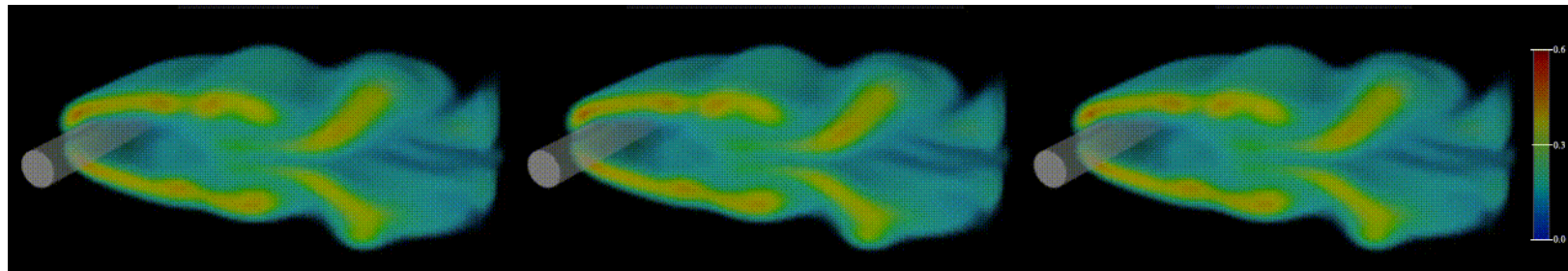
Differentiable physics = using autodifferentiation to differentiate physical algorithms

NS solver results

Low fidelity FD solver

Hybrid approach

High fidelity FD solver



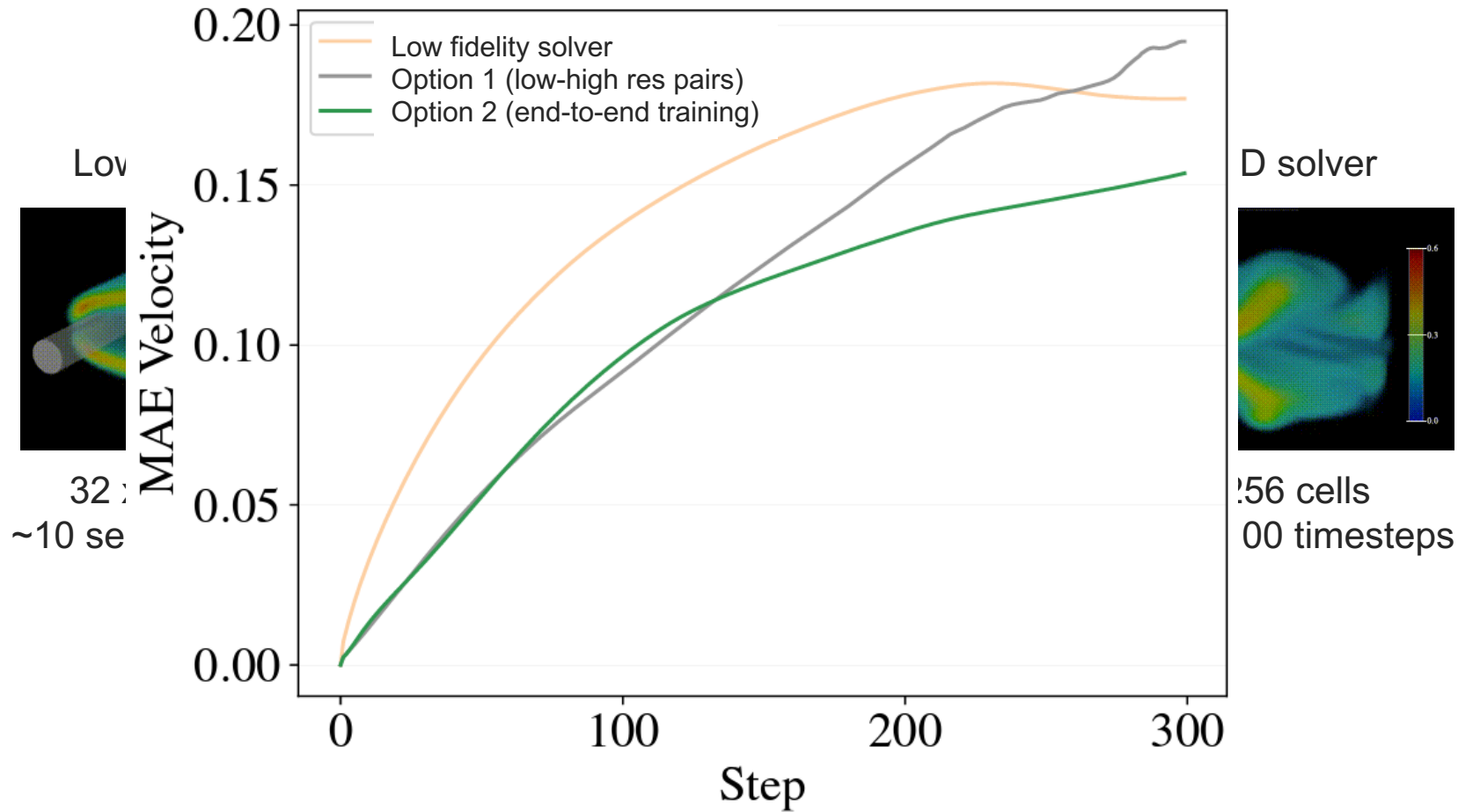
32 x 32 x 64 grid cells
~10 seconds / 100 timesteps

128 x 128 x 256 cells
~1000 seconds / 100 timesteps

32 x 32 x 64 grid cells
~15 seconds / 100 timesteps

Um et al, Solver-in-the-loop: Learning from differentiable physics to interact with iterative PDE-solvers, NeurIPS (2020)

NS solver results



Um et al, Solver-in-the-loop: Learning from differentiable physics to interact with iterative PDE-solvers, NeurIPS (2020)

Lecture overview

- Limitations of SciML approaches studied so far
- Hybrid SciML approaches
 - Residual modelling
 - Opening the “black-box”
 - How to train hybrid approaches
- Autodifferentiation
 - Autodifferentiation as a key enabler
 - What it is and how it works

Learning objectives

- Be able to describe what a hybrid workflow is
- Understand how autodifferentiation is used to train hybrid workflows
- Understand how autodifferentiation works

5 min break

Lecture overview

- Limitations of SciML approaches studied so far
- Hybrid SciML approaches
 - Residual modelling
 - Opening the “black-box”
 - How to train hybrid approaches
- Autodifferentiation
 - Autodifferentiation as a key enabler
 - What it is and how it works

Learning objectives

- Be able to describe what a hybrid workflow is
- Understand how autodifferentiation is used to train hybrid workflows
- Understand how autodifferentiation works

Autodifferentiation is a key enabler



Autodifferentiation is a **key enabler** of all the SciML techniques studied so far

It allows us to **efficiently** differentiate through complicated loss functions and get gradients of learnable parameters

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

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

Physics-informed neural network

$$a(x) \xrightarrow{\mathcal{G}_\theta^*[a]} \{a_k\}_{k=1}^m \rightarrow NN(\{a_k\}; \theta) \rightarrow \{u_k\}_{k=1}^p \rightarrow \hat{u}(x)$$

$$L(\theta) = \frac{1}{NM} \sum_i^N \sum_j^M \|u_i(x_j) - \mathcal{G}_\theta^*[a_i](x_j)\|^2$$

Operator learning

$$L(\theta) = \sum_i^N \sum_t^T \|HybridSolver_t(u_{0_i}; \theta) - u_t^H(u_{0_i})\|^2$$

Hybrid algorithms

Programs as vector functions



Many (scientific) programs can be decomposed in the following way:

```
def Hybrid_NS_solver(u_0, p_0, rho, nu, theta):
    "Pseudocode for solving NS equation, with NN correction"

    # u_0, p_0 have shape (NX, NY, NZ)
    u_t, p_t = u_0, p_0
    for t in range(0, T):
        u_star = f(u_t, p_t, rho, nu)
        p_t = matrix_solve(u_star, p_t, rho)
        u_t = g(u_t, p_t, rho, nu)

        u_t, p_t = (u_t, p_t) + NN(u_t, p_t, theta)

    return u_t, p_t

theta.requires_grad_(True)
u_T, _ = Hybrid_NS_solver(u_0, p_0, rho, nu, theta)
loss = loss_fn(u_T, u_T_true)
dtheta = torch.autograd.grad(loss, theta)
# for learning theta (training NN)
```

Program:

Input: a vector $x \in \mathbb{R}^n$

*Function: A series of **primitive operations** on the elements of x
add / multiply / trigonometric / ...*

Output: some transformed vector $y \in \mathbb{R}^m$

Mathematically, the program defines a vector function $\mathbf{y}: \mathbb{R}^n \rightarrow \mathbb{R}^m$, composed of primitive operations:

$$\mathbf{y}(x) = f_N \circ \dots \circ f_2 \circ f_1(x)$$

Chain rule for vector functions

Consider **any** vector function $\mathbf{y}: \mathbb{R}^n \rightarrow \mathbb{R}^m$, composed from many other vector functions

$$\mathbf{y}(\mathbf{x}) = \mathbf{f}_N \circ \dots \circ \mathbf{f}_2 \circ \mathbf{f}_1(\mathbf{x})$$

Then we can use the **multivariate chain rule** (= matrix multiplication of Jacobians) to evaluate its derivatives

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \frac{\partial \mathbf{f}_N}{\partial \mathbf{f}_{N-1}}, \dots, \frac{\partial \mathbf{f}_2}{\partial \mathbf{f}_1} \frac{\partial \mathbf{f}_1}{\partial \mathbf{x}}$$


where

$$J_{\mathbf{y}} \equiv \frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{pmatrix} \frac{\partial y_1}{\partial x_1} & \dots & \frac{\partial y_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \dots & \frac{\partial y_m}{\partial x_n} \end{pmatrix}$$

Autodifferentiation

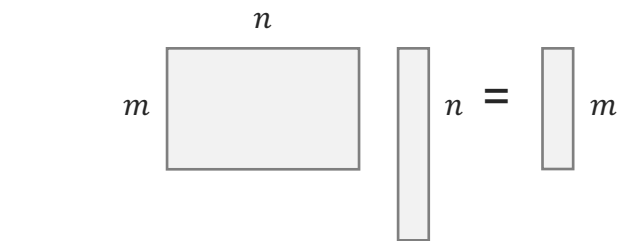
Modern autodifferentiation libraries allow us to efficiently compute:

The vector-Jacobian product (vjp):

$$\mathbf{v}^T \frac{\partial \mathbf{y}}{\partial \mathbf{x}}$$


The diagram illustrates the vector-Jacobian product (vjp). On the left, a horizontal rectangle representing a vector of size m is positioned above another horizontal rectangle representing a Jacobian matrix of size $m \times n$. The label m is placed above the first rectangle and below the second. An equals sign follows, leading to a single horizontal rectangle representing the resulting vector of size n , with the label n above it.

or the Jacobian-vector product (jvp):

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \mathbf{v}$$



The diagram illustrates the Jacobian-vector product (jvp). On the left, a horizontal rectangle representing a Jacobian matrix of size $m \times n$ is positioned above a vertical rectangle representing a vector of size n . The label m is placed to the left of the first rectangle, and the label n is placed above the second. An equals sign follows, leading to a vertical rectangle representing the resulting vector of size m , with the label m to its right.

of **arbitrary** programs.

Autodifferentiation

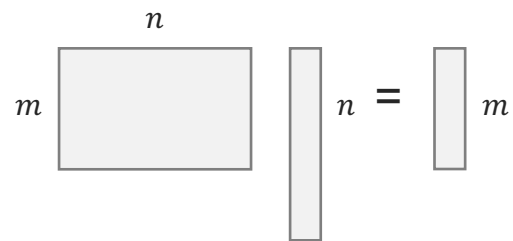
Modern autodifferentiation libraries allow us to efficiently compute:

The vector-Jacobian product (vjp):

$$\mathbf{v}^T \frac{\partial \mathbf{y}}{\partial \mathbf{x}}$$


The diagram illustrates the vector-Jacobian product (vjp). It shows a horizontal vector of size m (represented by a box with m above and below it) multiplied by a square Jacobian matrix of size $m \times n$ (represented by a box with m on the left and n on top). The result is a horizontal vector of size n (represented by a box with n above it).

or the Jacobian-vector product (jvp):

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \mathbf{v}$$


The diagram illustrates the Jacobian-vector product (jvp). It shows a square Jacobian matrix of size $m \times n$ (represented by a box with m on the left and n on top) multiplied by a vertical vector of size n (represented by a box with n to its right). The result is a vertical vector of size m (represented by a box with m to its right).

of **arbitrary** programs.

Why is it useful to evaluate the vjp / jvp?

Autodifferentiation

Modern autodifferentiation libraries allow us to efficiently compute:

The vector-Jacobian product (vjp):

$$\mathbf{v}^T \frac{\partial \mathbf{y}}{\partial \mathbf{x}}$$

or the Jacobian-vector product (jvp):

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \mathbf{v}$$

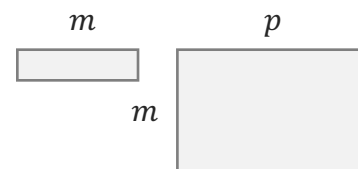
of **arbitrary** programs.

Why is it useful to evaluate the vjp / jvp?

Consider training a neural network:

$$L(\boldsymbol{\theta}): \mathbb{R}^p \rightarrow \mathbb{R}^1$$

$$\frac{\partial L}{\partial \boldsymbol{\theta}} = \frac{\partial L}{\partial \mathbf{NN}} \frac{\partial \mathbf{NN}}{\partial \boldsymbol{\theta}}$$




m = total number of network outputs
 p = total number of parameters

We need to compute a vjp!

Chain rule

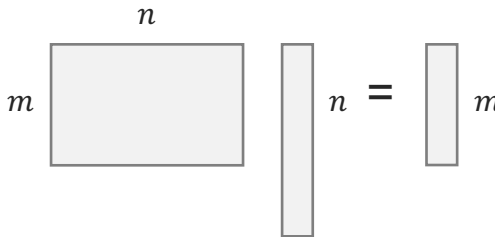
Modern autodifferentiation libraries allow us to efficiently compute:

The vector-Jacobian product (vjp):

$$\mathbf{v}^T \frac{\partial \mathbf{y}}{\partial \mathbf{x}}$$


The diagram shows a horizontal rectangle representing a row vector of size m multiplied by a square matrix of size $m \times n$. The result is a horizontal rectangle representing a row vector of size n .

or the Jacobian-vector product (jvp):

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \mathbf{v}$$


The diagram shows a square matrix of size $m \times n$ multiplied by a vertical rectangle representing a column vector of size n . The result is a vertical rectangle representing a column vector of size m .

of **arbitrary** programs.

We can evaluate the vjp / jvp using the **chain rule**, for example:

$$\mathbf{v}^T \frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \mathbf{v}^T \frac{\partial \mathbf{f}_N}{\partial \mathbf{f}_{N-1}}, \dots, \frac{\partial \mathbf{f}_2}{\partial \mathbf{f}_1} \frac{\partial \mathbf{f}_1}{\partial \mathbf{x}}$$

Dimensionality

Let's think about dimensionality. Consider a simple MLP with $m = 100$ outputs, $h = 100$ hidden units, and 10,000 inputs. Then W_1 has $100 \times 10,000 = 1\text{M}$ elements.

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

Then:

$$\frac{\partial L}{\partial W_1} = \frac{\partial L}{\partial \mathbf{f}} \frac{\partial \mathbf{f}}{\partial W_1} = \frac{\partial L}{\partial \mathbf{f}} \frac{\partial \mathbf{f}}{\partial \mathbf{g}} \frac{\partial \mathbf{g}}{\partial W_1}$$

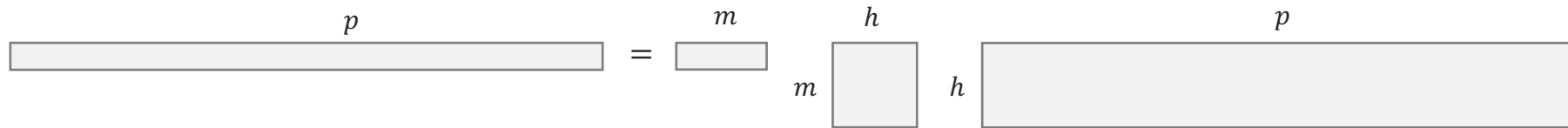
Dimensionality

Let's think about dimensionality. Consider a simple MLP with $m = 100$ outputs, $h = 100$ hidden units, and 10,000 inputs. Then W_1 has $100 \times 10,000 = 1\text{M}$ elements.

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

Then:

$$\frac{\partial L}{\partial W_1} = \frac{\partial L}{\partial \mathbf{f}} \frac{\partial \mathbf{f}}{\partial W_1} = \frac{\partial L}{\partial \mathbf{f}} \frac{\partial \mathbf{f}}{\partial \mathbf{g}} \frac{\partial \mathbf{g}}{\partial W_1}$$



$$(1 \times 1\text{M}) = (1 \times 100) (100 \times 100) (100 \times 1\text{M})$$

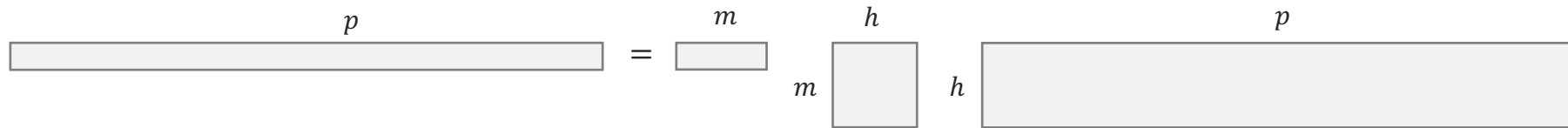
Dimensionality

Let's think about dimensionality. Consider a simple MLP with $m = 100$ outputs, $h = 100$ hidden units, and 10,000 inputs. Then W_1 has $100 \times 10,000 = 1\text{M}$ elements.

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

Then:

$$\frac{\partial L}{\partial W_1} = \frac{\partial L}{\partial \mathbf{f}} \frac{\partial \mathbf{f}}{\partial W_1} = \frac{\partial L}{\partial \mathbf{f}} \frac{\partial \mathbf{f}}{\partial \mathbf{g}} \frac{\partial \mathbf{g}}{\partial W_1}$$



Consider evaluating the chain rule (RHS): $(1 \times 1\text{M}) = (1 \times 100) (100 \times 100) (100 \times 1\text{M})$

1) From right to left (forward)

$$(100 \times 100) (100 \times 1\text{M}) \rightarrow (100 \times 1\text{M})$$

$$(1 \times 100) (100 \times 1\text{M}) \rightarrow (1 \times 1\text{M})$$

= lots of computation (large matrix-matrix multiply)

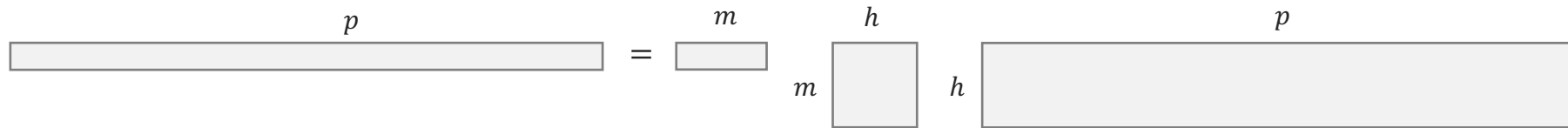
Dimensionality

Let's think about dimensionality. Consider a simple MLP with $m = 100$ outputs, $h = 100$ hidden units, and 10,000 inputs. Then W_1 has $100 \times 10,000 = 1\text{M}$ elements.

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

Then:

$$\frac{\partial L}{\partial W_1} = \frac{\partial L}{\partial \mathbf{f}} \frac{\partial \mathbf{f}}{\partial W_1} = \frac{\partial L}{\partial \mathbf{f}} \frac{\partial \mathbf{f}}{\partial \mathbf{g}} \frac{\partial \mathbf{g}}{\partial W_1}$$



Consider evaluating the chain rule (RHS):

$$(1 \times 1\text{M}) = (1 \times 100) (100 \times 100) (100 \times 1\text{M})$$

1) From right to left (forward)

$$\begin{aligned} (100 \times 100) (100 \times 1\text{M}) &\rightarrow (100 \times 1\text{M}) \\ (1 \times 100) (100 \times 1\text{M}) &\rightarrow (1 \times 1\text{M}) \end{aligned}$$

= lots of computation (large matrix-matrix multiply)

2) From left to right (reverse)

$$\begin{aligned} (1 \times 100) (100 \times 100) &\rightarrow (1 \times 100) \\ (1 \times 100) (100 \times 1\text{M}) &\rightarrow (1 \times 1\text{M}) \end{aligned}$$

= much less computation (vector-matrix multiplies)

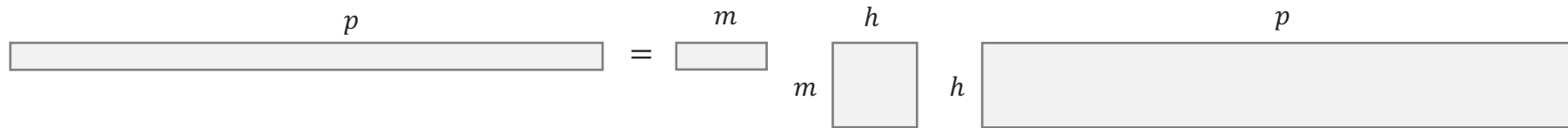
Dimensionality

Let's think about dimensionality. Consider a simple MLP with $m = 100$ outputs, $h = 100$ hidden units, and 10,000 inputs. Then W_1 has $100 \times 10,000 = 1\text{M}$ elements.

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

Then:

$$\frac{\partial L}{\partial W_1} = \frac{\partial L}{\partial \mathbf{f}} \frac{\partial \mathbf{f}}{\partial W_1} = \frac{\partial L}{\partial \mathbf{f}} \frac{\partial \mathbf{f}}{\partial \mathbf{g}} \frac{\partial \mathbf{g}}{\partial W_1}$$



Consider evaluating the chain rule (RHS): $(1 \times 1\text{M}) = (1 \times 100) (100 \times 100) (100 \times 1\text{M})$

1) From right to left (forward)

$$\begin{aligned} (100 \times 100) (100 \times 1\text{M}) &\rightarrow (100 \times 1\text{M}) \\ (1 \times 100) (100 \times 1\text{M}) &\rightarrow (1 \times 1\text{M}) \end{aligned}$$

= lots of computation (large matrix-matrix multiply)

2) From left to right (reverse)

$$\begin{aligned} (1 \times 100) (100 \times 100) &\rightarrow (1 \times 100) \\ (1 \times 100) (100 \times 1\text{M}) &\rightarrow (1 \times 1\text{M}) \end{aligned}$$

= much less computation (vector-matrix multiplies)



=> Order matters! Evaluating vjps in reverse mode is usually most efficient

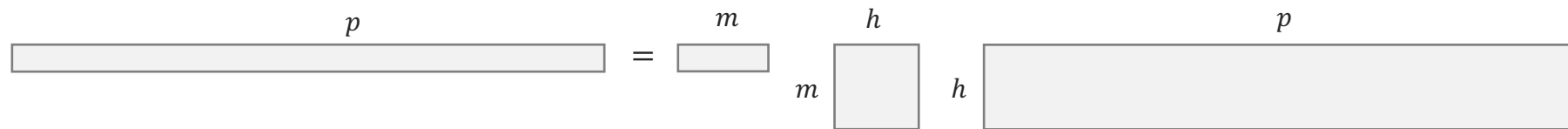
Dimensionality

Let's think about dimensionality. Consider a simple MLP with $m = 100$ outputs, $h = 100$ hidden units, and 10,000 inputs. Then W_1 has $100 \times 10,000 = 1\text{M}$ elements.

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

Then:

$$\frac{\partial L}{\partial W_1} = \frac{\partial L}{\partial \mathbf{f}} \frac{\partial \mathbf{f}}{\partial W_1} = \frac{\partial L}{\partial \mathbf{f}} \frac{\partial \mathbf{f}}{\partial \mathbf{g}} \frac{\partial \mathbf{g}}{\partial W_1}$$



Consider evaluating the chain rule (RHS):

$$(1 \times 1\text{M}) = (1 \times 100) (100 \times 100) (100 \times 1\text{M})$$

2) From left to right (reverse)

$$(1 \times 100) (100 \times 100) \rightarrow (1 \times 100)$$

$$(1 \times 100) (100 \times 1\text{M}) \rightarrow (1 \times 1\text{M})$$

= much less computation (vector-matrix multiplies)

But what about the last computation $(\frac{\partial L}{\partial \mathbf{g}} \frac{\partial \mathbf{g}}{\partial W_1})$?

This is still expensive! ($\frac{\partial \mathbf{g}}{\partial W_1}$ has 100M elements, or ~0.5 GB)



=> Order matters! Evaluating vjps in reverse mode is usually most efficient

Dimensionality

Note:

$$\frac{\partial \mathbf{g}}{\partial W_1} = \frac{\partial (W_1 \mathbf{x} + \mathbf{b}_1)}{\partial W_1} = \begin{pmatrix} x_1 & x_2 & \cdots & x_n & 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \cdots & 0 & x_1 & x_2 & \cdots & x_n \end{pmatrix}$$

Then

$$\begin{aligned} \frac{\partial L}{\partial W_1} &= \frac{dL}{d\mathbf{g}} \frac{\partial \mathbf{g}}{\partial W_1} = \left(\frac{dL}{dg_1}(x_1, \dots, x_{10,000}), \dots, \frac{dL}{dg_{100}}(x_1, \dots, x_{10,000}) \right) \\ &= \frac{\partial L}{\partial \mathbf{g}} \otimes \mathbf{x} \end{aligned}$$

This is just the (flattened) **outer product** of two vectors $(100 \times 1) \otimes (10,000 \times 1)$

- ⇒ We don't have to fully populate the last Jacobian $\left(\frac{\partial \mathbf{g}}{\partial W_1}\right)$ when computing its vector-Jacobian product
- ⇒ Often, vjps (and jvps) can be computed efficiently **without** needing to populate the full Jacobian



Dimensionality

Another example:

Consider:

$$\mathbf{y} = \sin(\mathbf{x})$$

Then

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{pmatrix} \cos(x_1) & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \cos(x_n) \end{pmatrix}$$

And

$$\begin{aligned} \mathbf{v}^T \frac{\partial \mathbf{y}}{\partial \mathbf{x}} &= (v_1 \cos(x_1), \dots, v_n \cos(x_n)) \\ &= \mathbf{v} \cdot \cos(\mathbf{x}) \end{aligned}$$

Requires $\mathcal{O}(n)$ operations



- ⇒ We don't have to fully populate the last Jacobian ($\frac{\partial g}{\partial w_1}$) when computing its vector-Jacobian product
- ⇒ Often, vjps (and jvps) can be computed efficiently **without** needing to populate the full Jacobian

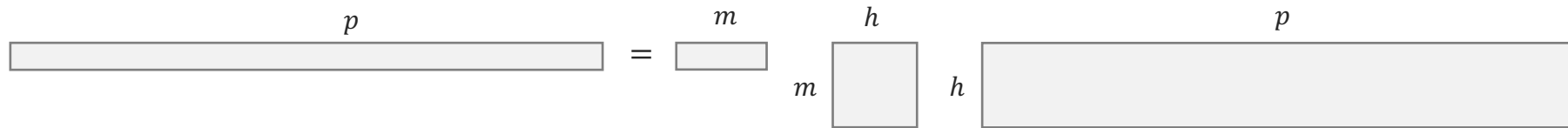
Dimensionality

Let's think about dimensionality. Consider a simple MLP with $m = 100$ outputs, $h = 100$ hidden units, and 10,000 inputs. Then W_1 has $100 \times 10,000 = 1\text{M}$ elements.

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

Then:

$$\frac{\partial L}{\partial W_1} = \frac{\partial L}{\partial \mathbf{f}} \frac{\partial \mathbf{f}}{\partial W_1} = \frac{\partial L}{\partial \mathbf{f}} \frac{\partial \mathbf{f}}{\partial \mathbf{g}} \frac{\partial \mathbf{g}}{\partial W_1}$$



Consider evaluating the chain rule (RHS):

$$(1 \times 1\text{M}) = (1 \times 100) (100 \times 100) (100 \times 1\text{M})$$

= Efficient training code

Allows us to train neural networks with **billions** of parameters

2) From left to right (reverse)

$$(1 \times 100) (100 \times 100) \rightarrow (1 \times 100)$$

$$(1 \times 100) (100 \times 1\text{M}) \rightarrow (1 \times 1\text{M})$$

$$(100 \times 1) \otimes (10,000 \times 1) \rightarrow (1 \times 1\text{M})$$

= much less computation (vector-matrix multiplies)

Vector-Jacobian product

vjp:

$$\mathbf{v}^T \frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \mathbf{v}^T \frac{\partial \mathbf{f}_N}{\partial \mathbf{f}_{N-1}}, \dots, \frac{\partial \mathbf{f}_2}{\partial \mathbf{f}_1} \frac{\partial \mathbf{f}_1}{\partial \mathbf{x}}$$

We can compute $\mathbf{v}^T \frac{\partial \mathbf{y}}{\partial \mathbf{x}}$ by iteratively computing vector-Jacobian products, from left to right (reverse-mode):

Starting with \mathbf{v}^T ,

$$\mathbf{v}^T \leftarrow \mathbf{v}^T \frac{\partial \mathbf{f}_N}{\partial \mathbf{f}_{N-1}}$$

$$\mathbf{v}^T \leftarrow \mathbf{v}^T \frac{\partial \mathbf{f}_{N-1}}{\partial \mathbf{f}_{N-2}}$$

...

$$\mathbf{v}^T \leftarrow \mathbf{v}^T \frac{\partial \mathbf{f}_1}{\partial \mathbf{x}}$$

- We only need to define the **vjp** for each **primitive operation** to compute $\mathbf{v}^T \frac{\partial \mathbf{y}}{\partial \mathbf{x}}$
- Usually, we do not need to explicitly compute the full intermediate Jacobians $\frac{\partial \mathbf{f}_i}{\partial \mathbf{f}_{i-1}}$

Jacobian-vector product

jvp:

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \mathbf{v} = \frac{\partial f_N}{\partial f_{N-1}}, \dots, \frac{\partial f_2}{\partial f_1} \frac{\partial f_1}{\partial \mathbf{x}} \mathbf{v}$$

We can compute $\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \mathbf{v}$ by iteratively computing Jacobian-vector products, from right to left (forward-mode):

Starting with \mathbf{v} ,

$$\mathbf{v} \leftarrow \frac{\partial f_1}{\partial \mathbf{x}} \mathbf{v}$$

$$\mathbf{v} \leftarrow \frac{\partial f_2}{\partial f_1} \mathbf{v}$$

$$\mathbf{v} \leftarrow \frac{\overset{\dots}{\partial f_N}}{\partial f_{N-1}} \mathbf{v}$$

- We only need to define the **jvp** for each **primitive operation** to compute $\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \mathbf{v}$
- Usually, we do not need to explicitly compute the full intermediate Jacobians $\frac{\partial f_i}{\partial f_{i-1}}$

Full Jacobian

- What if we want the full Jacobian? $J_y = \frac{\partial y}{\partial x}$

Full Jacobian

- What if we want the full Jacobian? $J_y = \frac{\partial \mathbf{y}}{\partial \mathbf{x}}$
- We can combine vjps / jvps to compute the full Jacobian row by row / column by column if necessary

Let

$$\mathbf{v}^T = (1, 0, \dots, 0)$$

Then

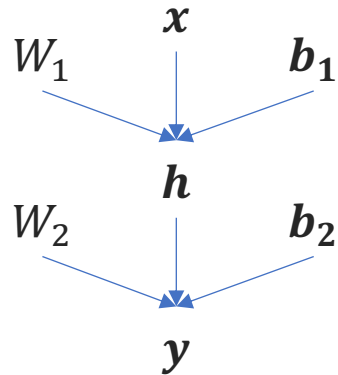
$$\mathbf{v}^T \frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \left(\frac{\partial y_1}{\partial x_1}, \dots, \frac{\partial y_1}{\partial x_n} \right)$$

= First row of Jacobian

- Note jvps are usually more efficient for “tall” Jacobians, whilst vjps are more efficient for “wide” Jacobians

Autodiff in practice

$$y = W_2 \sigma(W_1 x + b_1) + b_2$$



- 1) Decompose given function into its **primitive** operations
- 2) Build a **directed graph** of these operations

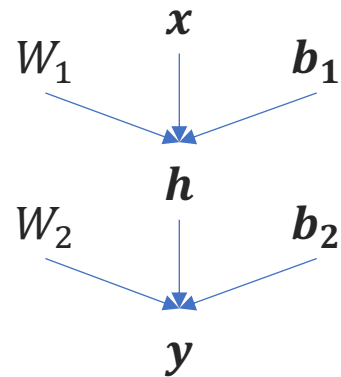
 PyTorch



 TensorFlow

Autodiff in practice

$$y = W_2 \sigma(W_1 x + b_1) + b_2$$



- 1) Decompose given function into its **primitive** operations
- 2) Build a **directed graph** of these operations
- 3) For each primitive operation, define
 - 1) Forward operation
 - 2) vector-Jacobian product
 - 3) Jacobian-vector product

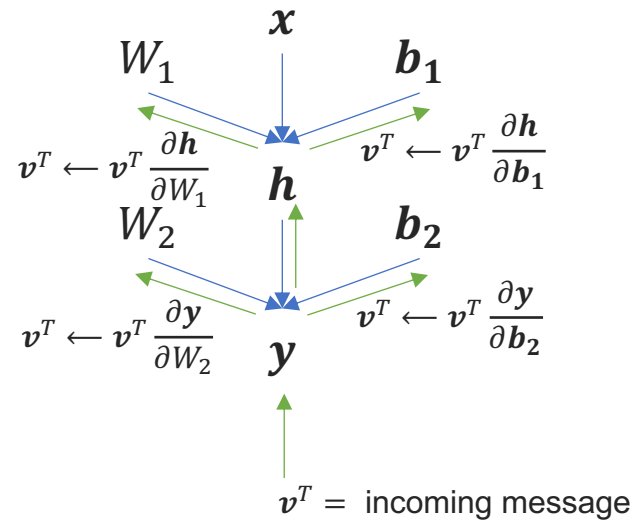
 PyTorch



 TensorFlow

Autodiff in practice

$$y = W_2 \sigma(W_1 x + b_1) + b_2$$



- 1) Decompose given function into its **primitive** operations
- 2) Build a **directed graph** of these operations
- 3) For each primitive operation, define
 - 1) Forward operation
 - 2) vector-Jacobian product
 - 3) Jacobian-vector product
- 4) Evaluate the vjp or jvp of the function by applying the **chain rule (=message passing)** through the graph
 - 1) Forwards for jvp
 - 2) Backwards for vjp

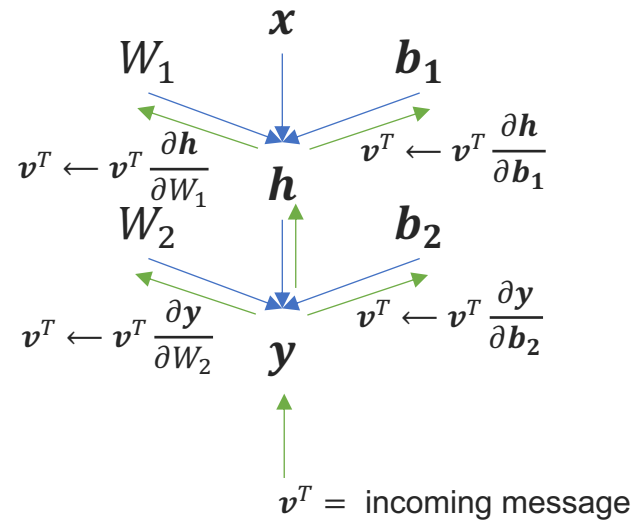


TensorFlow



Autodiff in practice

$$y = W_2 \sigma(W_1 x + b_1) + b_2$$



- 1) Decompose given function into its **primitive** operations
- 2) Build a **directed graph** of these operations
- 3) For each primitive operation, define
 - 1) Forward operation
 - 2) vector-Jacobian product
 - 3) Jacobian-vector product
- 4) Evaluate the vjp or jvp of the function by applying the **chain rule (=message passing)** through the graph
 - 1) Forwards for jvp
 - 2) Backwards for vjp

PyTorch



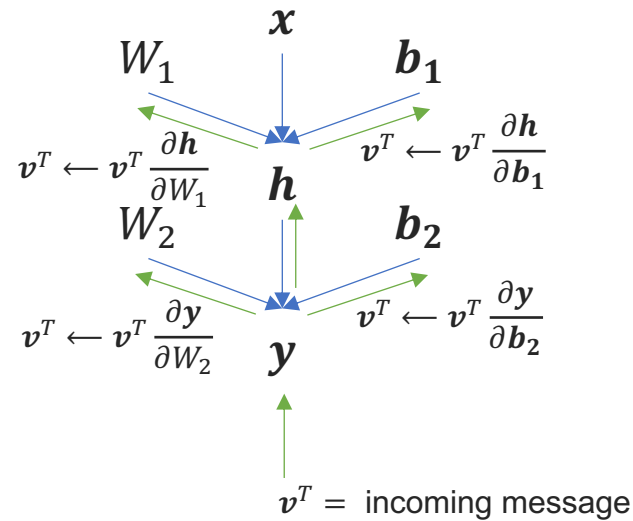
TensorFlow



- How does required memory scale with depth of computation for vjp vs jvp?

Autodiff in practice

$$y = W_2 \sigma(W_1 x + b_1) + b_2$$



- 1) Decompose given function into its **primitive** operations
- 2) Build a **directed graph** of these operations
- 3) For each primitive operation, define
 - 1) Forward operation
 - 2) vector-Jacobian product
 - 3) Jacobian-vector product
- 4) Evaluate the vjp or jvp of the function by applying the **chain rule (=message passing)** through the graph
 - 1) Forwards for jvp
 - 2) Backwards for vjp

PyTorch

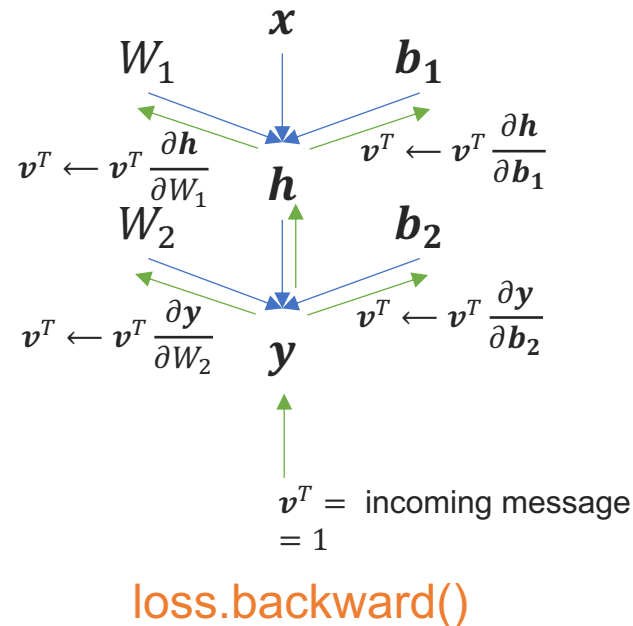


TensorFlow

- How does required memory scale with depth of computation for vjp vs jvp?
- vjp: memory scales linearly with depth (need to store forward computations)
- jvp: memory independent of depth (can compute jvp alongside forward pass)

Autodiff in practice

$$y = W_2 \sigma(W_1 x + b_1) + b_2$$



```
torch.autograd.grad(outputs, inputs, grad_outputs=None, retain_graph=None,
create_graph=False, only_inputs=True, allow_unused=None, is_grads_batched=False,
materialize_grads=False) [SOURCE]
```

Computes and returns the sum of gradients of outputs with respect to the inputs.

`grad_outputs` should be a sequence of length matching `output` containing the “vector” in vector-Jacobian product, usually the pre-computed gradients w.r.t. each of the outputs. If an output doesn’t `require_grad`, then the gradient can be `None`).

Note autodiff is **not**

- Symbolic differentiation
- Finite differences

It is a way of efficiently computing exact gradients!

Hybrid workflows in practice

Step 1: **rewrite** your traditional scientific algorithm in an autodifferentiation framework (e.g. PyTorch/JAX)

Hybrid workflows in practice

Step 1: **rewrite** your traditional scientific algorithm in an autodifferentiation framework (e.g. PyTorch/JAX)

Step 2: make parts of this algorithm **learnable** (either to accelerate it, or to improve accuracy)

Hybrid workflows in practice

Step 1: **rewrite** your traditional scientific algorithm in an autodifferentiation framework (e.g. PyTorch/JAX)

Step 2: make parts of this algorithm **learnable** (either to accelerate it, or to improve accuracy)

Step 3: get some **training examples** of what you want the input/output of the algorithm to be

Hybrid workflows in practice

Step 1: **rewrite** your traditional scientific algorithm in an autodifferentiation framework (e.g. PyTorch/JAX)

Step 2: make parts of this algorithm **learnable** (either to accelerate it, or to improve accuracy)

Step 3: get some **training examples** of what you want the input/output of the algorithm to be

Step 4: train your algorithm by **(auto)differentiating** through it and using gradient descent

Hybrid workflows in practice

Step 1: **rewrite** your traditional scientific algorithm in an autodifferentiation framework (e.g. PyTorch/JAX)

Step 2: make parts of this algorithm **learnable** (either to accelerate it, or to improve accuracy)

Step 3: get some **training examples** of what you want the input/output of the algorithm to be

Step 4: train your algorithm by **(auto)differentiating** through it and using gradient descent

Bonus: your code now runs on the GPU!

Summary

- Hybrid approaches insert learnable components **inside** traditional algorithms
- Autodifferentiation is **the key enabler** for SciML
 - Allows hybrid approaches to be trained end-to-end
 - Is an incredibly general and powerful tool