



# Structure-preserving learning of embedded closure models

Benjamin Sanderse, Syver Agdestein, Toby van Gastelen, Henrik Rosenberger, Hugo Melchers

20<sup>th</sup> June 2023

BIRS Workshop on Scientific Machine Learning



CWI

# Scientific Machine Learning semester programme

**JOIN OUR EVENTS!**

**RESEARCH SEMESTER PROGRAMME**  
Scientific Computing

**Scientific Machine Learning**

9 - 13 October  
**AUTUMN SCHOOL**

23 November  
**SOCIETY & INDUSTRY**

6 - 8 December  
**WORKSHOP**

Throughout the program  
**SEMINAR++**

  
[cwi.nl/semesterprogramme](http://cwi.nl/semesterprogramme)

**AUTUMN SCHOOL**

- ▶ Data-driven reduced-order models
- ▶ Neural networks and differential equations
- ▶ Data-driven multiscale modeling

9 - 13 OCTOBER 2023

**Scientific Machine Learning & Dynamical Systems**

**SPEAKERS**

Andrea Beck  
Erik Bekkers  
Steven Brunton  
Urban Fasel  
Hod Lipson  
Andrea Manzoni  
Chris Rackauckas  
Karen Veroy-Grepl  
Max Welling

Jointly organised by

  
**4TU.AMI**  
**AI INSTITUTE FOR DYNAMIC SYSTEMS**  
**NDNS+**

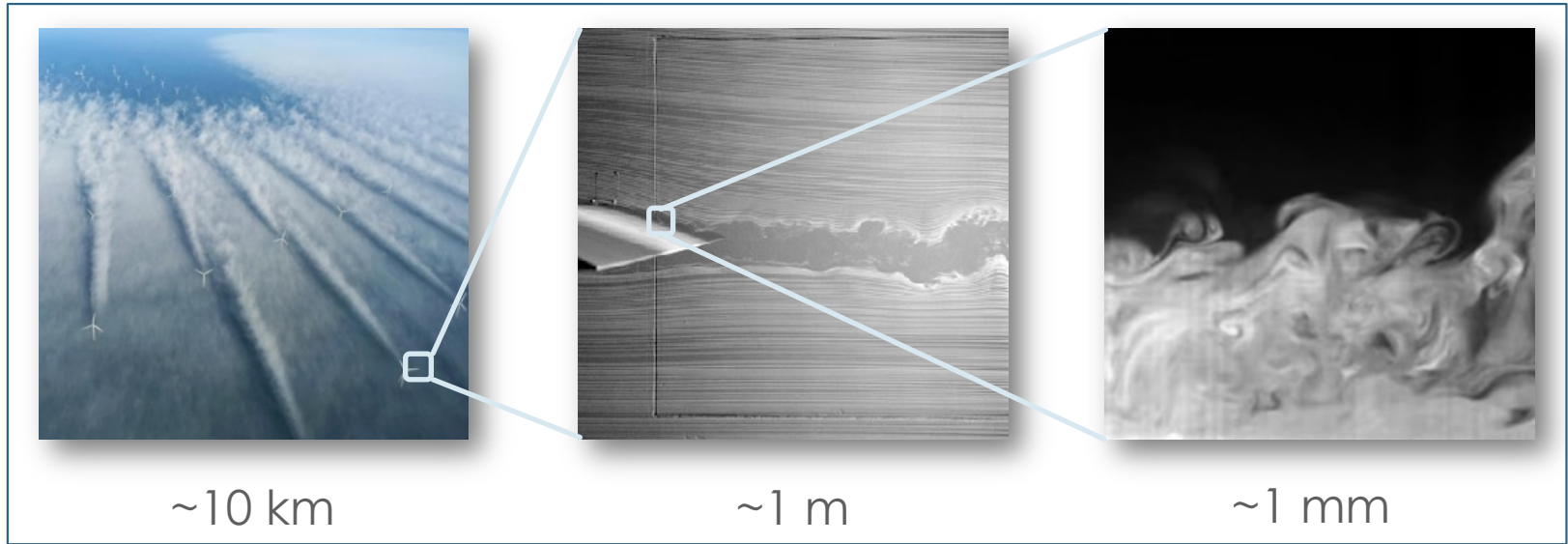
Organisers

Martina Chirilus-Brukner  
Leiden University  
Nathan Kutz  
University of Washington, USA  
Benjamin Sandese  
CINEMA - Institute for Informatics  
Rensheng Guo  
University of Twente

More information

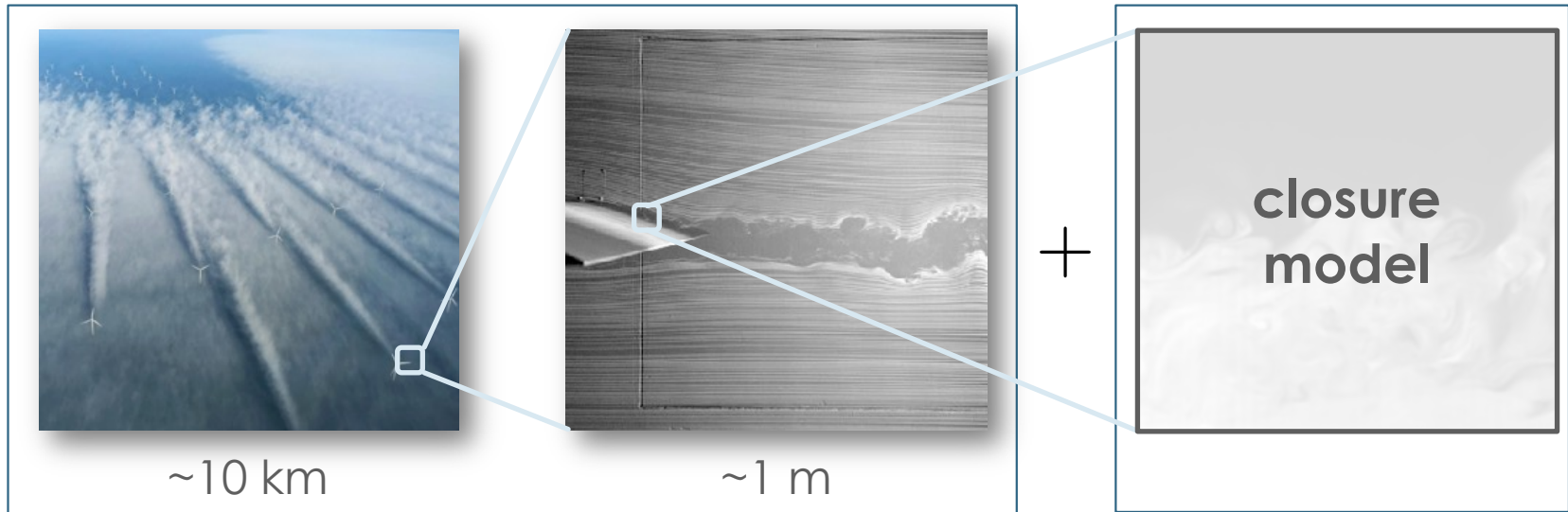


## Motivation: multiscale problems



Simulating all scales with a computational model is unfeasible

## Accurate and stable closure models needed



Closure model approximates effect of small scales on large scales

## Basics of closure modelling

- Multiscale fluid flow: Navier-Stokes equations

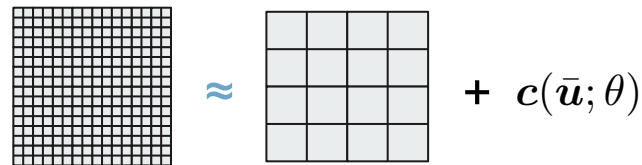
$$\frac{\partial \mathbf{u}}{\partial t} = \mathbf{f}(\mathbf{u}) \quad \mathbf{f}(\mathbf{u}) := -\nabla \cdot (\mathbf{u} \otimes \mathbf{u}) - \nabla p + \nu \nabla^2 \mathbf{u}$$

- NS describes (too) many scales of motion for small viscosity  $\nu$
- Reduce range of scales by a **filtering** operation:

$$\bar{\mathbf{u}} = \mathcal{A}(\mathbf{u}) \quad \mathcal{A}(\mathbf{u}) = \int \mathbf{u}(\xi, t) G(x, \xi) d\xi \quad \mathbf{u}' = \mathbf{u} - \bar{\mathbf{u}}$$

- **Aim:** use coarser meshes and larger time steps when solving for  $\bar{\mathbf{u}}$

## Basics of closure modelling



- Problem: filter and PDE operator **do not commute**
- Art: find a **closure model** with parameters  $\theta$  s.t.

$$c(\bar{\mathbf{u}}; \theta) \approx \mathcal{C}[\mathcal{A}, \mathbf{f}](\mathbf{u}) = \overline{\nabla \cdot (\mathbf{u} \otimes \mathbf{u})} - \nabla \cdot (\bar{\mathbf{u}} \otimes \bar{\mathbf{u}})$$

- Finding  $c(\bar{\mathbf{u}}; \theta)$  is an inverse problem (model discovery)
- Common form:  $\frac{\partial \bar{\mathbf{u}}}{\partial t} = \mathbf{f}(\bar{\mathbf{u}}) + c(\bar{\mathbf{u}}; \theta)$

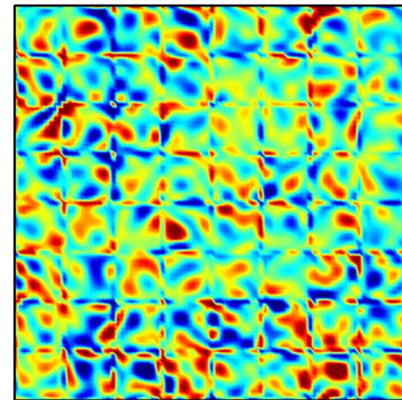
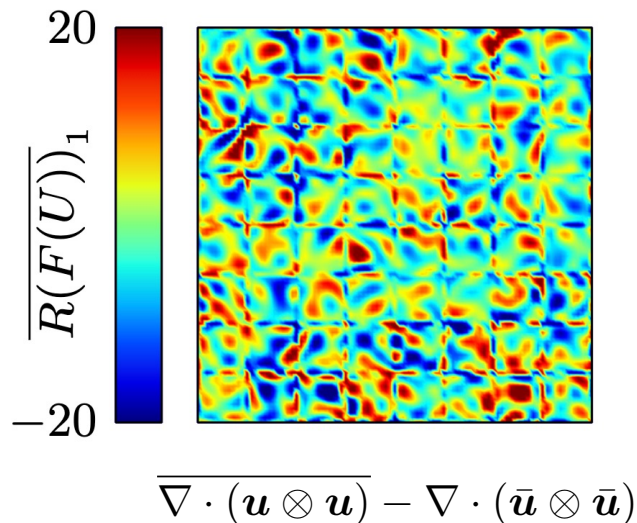
## Neural networks as closure model

$$c(\bar{\mathbf{u}}; \theta) = \text{NN}(\bar{\mathbf{u}}; \theta)$$

$$\frac{d\bar{\mathbf{u}}}{dt} = f(\bar{\mathbf{u}}) + \text{NN}(\bar{\mathbf{u}}; \theta)$$

Exact (—)

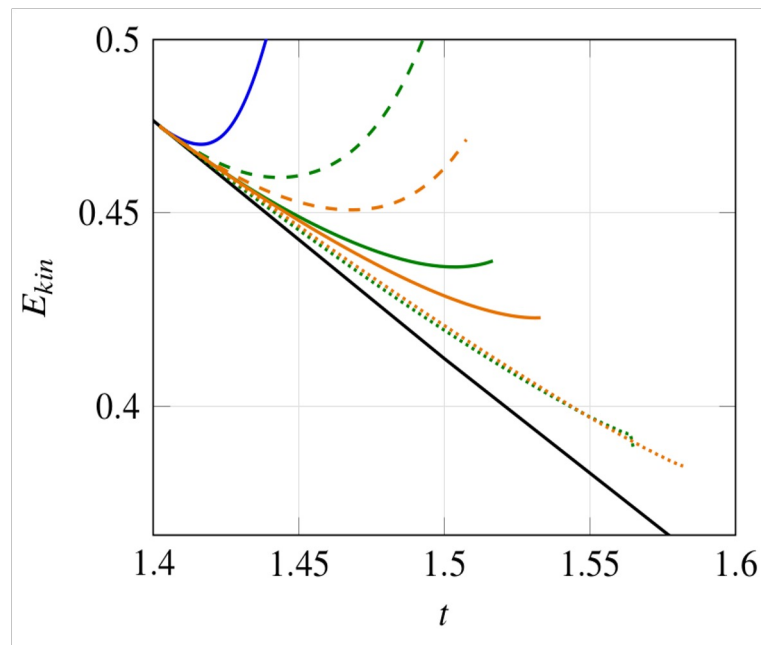
GRU3 (—)



## Issue: NNs destabilize the dynamical system

$$\frac{d\bar{\mathbf{u}}}{dt} = f(\bar{\mathbf{u}}) + \text{NN}(\bar{\mathbf{u}}; \theta)$$

- NN accurately matches closure term (operator fit)
- But: **solution is wrong**



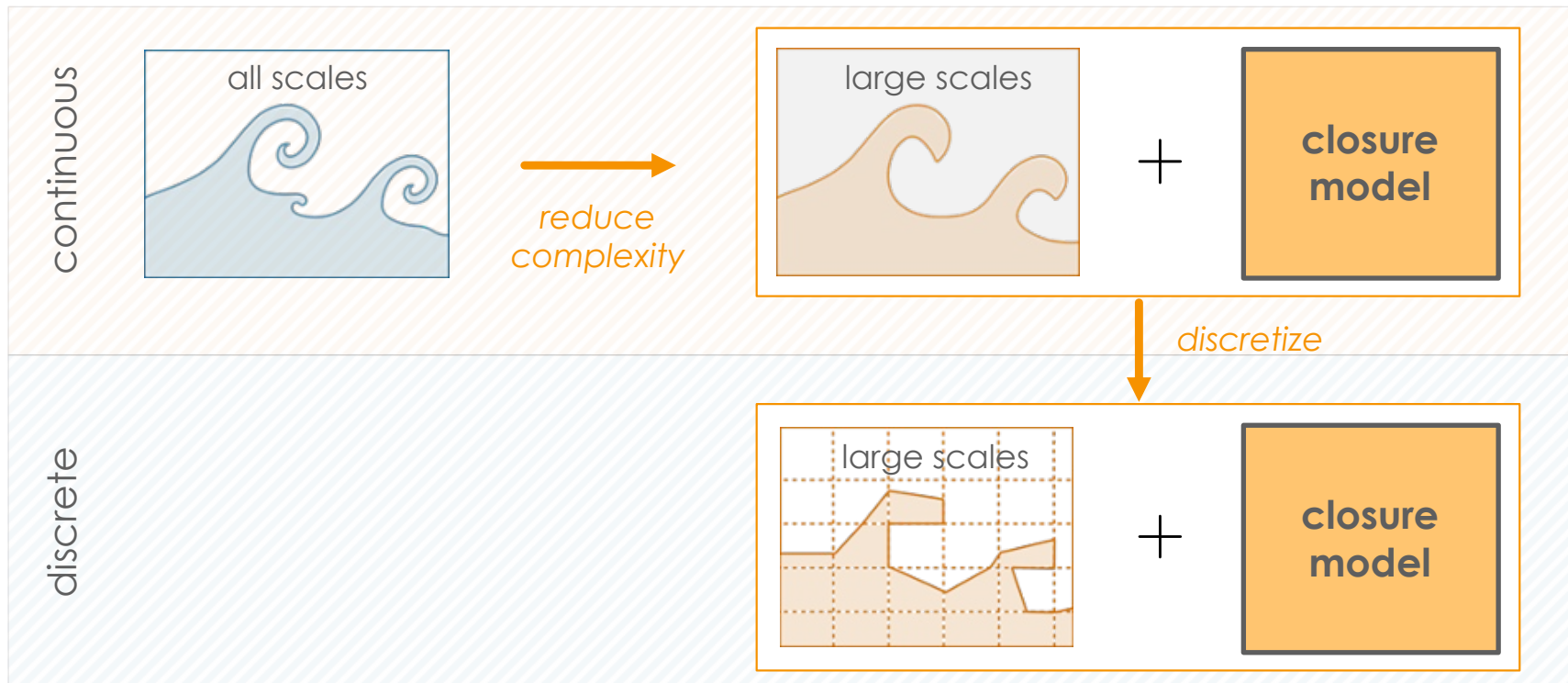


## Tackling instability in dynamical systems with NNs

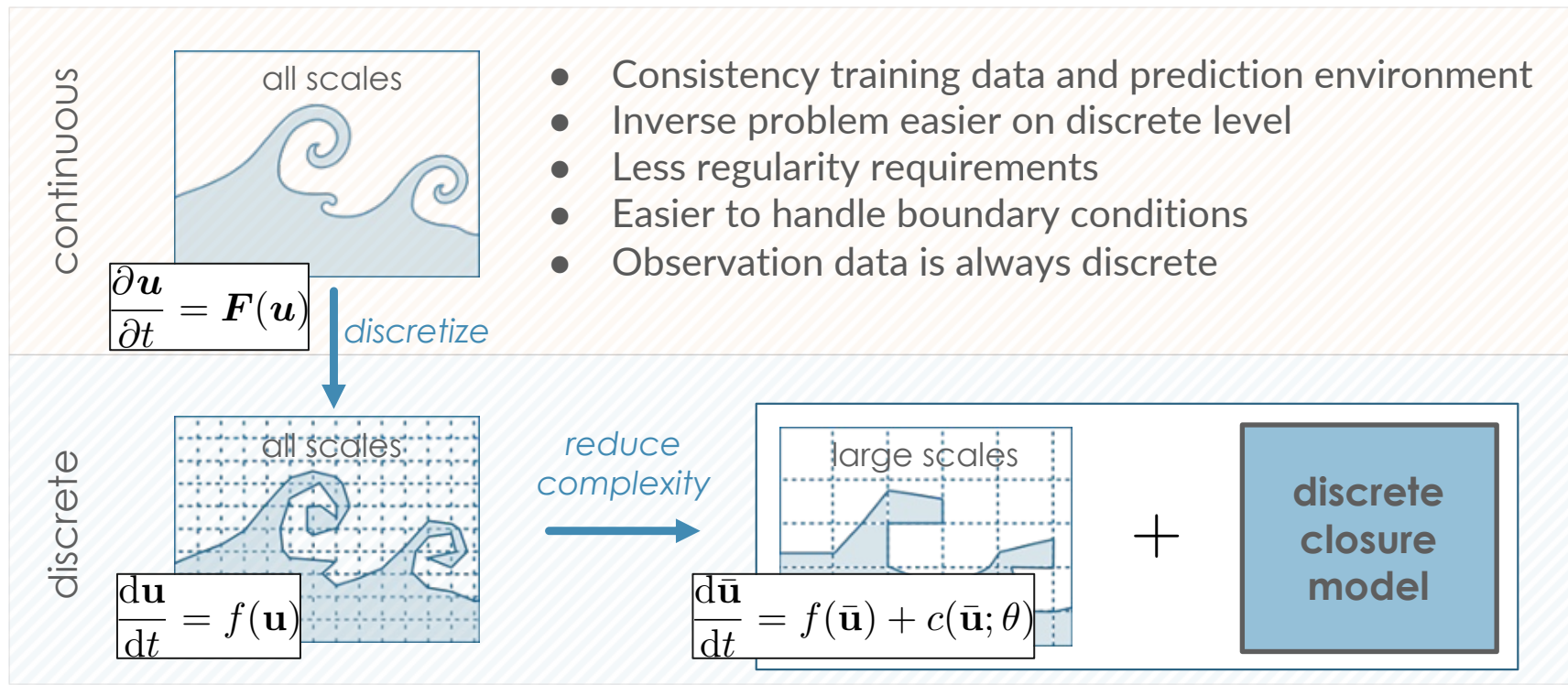
- **Instability** common problem for ML-based closure models (mismatch training environment and prediction environment)
- Recent approaches:
  - Stability training on data with artificial noise (Kurz & Beck, 2021)
  - Minimizing (or eliminating) backscatter (Park & Choi, 2021)
  - Projection onto a stable basis (Beck et al., 2019)
  - Trajectory fitting (List et al., 2022; MacArt et al., 2021)
  - Reinforcement learning (Bae & Koumoutsakos, 2022; Kurz et al. 2022)

Our approach: “discretize first” + “preserve structure” + “embedded learning”

## Common approach in closure modelling



## Alternative: discretize first



## Examples of preserving structure

- ODE formulation (“neural ODE”)
- Closure model form (“neural closure model”)
- Conservation form
- Translation invariance
- **Energy conservation**

$$\frac{d\bar{\mathbf{u}}}{dt} = \text{NN}(\bar{\mathbf{u}}; \theta)$$

$$\frac{d\bar{\mathbf{u}}}{dt} = f(\bar{\mathbf{u}}) + \text{NN}(\bar{\mathbf{u}}; \theta)$$

$$\frac{d\bar{\mathbf{u}}}{dt} = f(\bar{\mathbf{u}}) + \nabla \cdot \text{NN}(\bar{\mathbf{u}}; \theta)$$

CNN architecture

## Energy conservation implies stability

- Many PDEs possess **secondary conservation laws**, such as energy or entropy, which give a **stability bound**
- **Example: Korteweg – de Vries**

$$\frac{\partial u}{\partial t} + 3 \frac{\partial u^2}{\partial x} = - \frac{\partial^3 u}{\partial x^3} \quad \longrightarrow \quad \frac{dE}{dt} = \frac{d}{dt} \underbrace{\frac{1}{2} \int_{\Omega} u^2 d\Omega}_{=: E} = 0$$
$$E := \frac{1}{2} \int u^2 d\Omega$$

Idea: impose a similar structure on the filtered equations

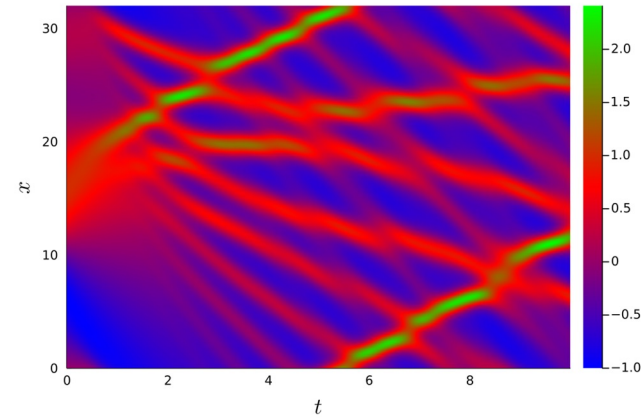
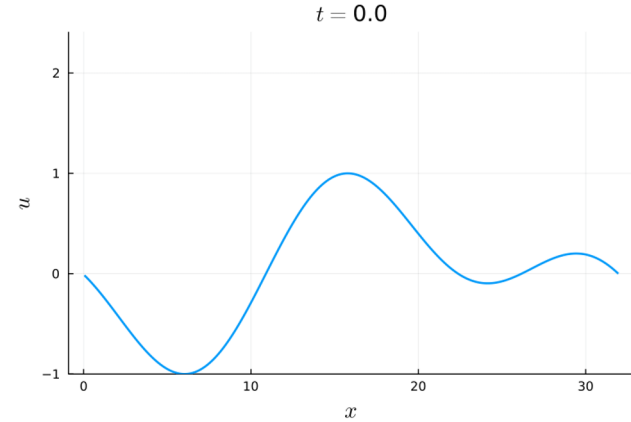
# Korteweg - de Vries equation

- KdV discretized using skew-symmetric scheme:

$$\frac{\partial u}{\partial t} + 3 \frac{\partial u^2}{\partial x} = - \frac{\partial^3 u}{\partial x^3} \quad \Longrightarrow \quad \frac{d\mathbf{u}}{dt} = -3\mathbf{G}(\mathbf{u}) - \mathbf{D}_3\mathbf{u}$$

- Energy conservation (periodic BCs):

$$\frac{dE}{dt} = \frac{d}{dt} \underbrace{\frac{1}{2} \int_{\Omega} u^2 d\Omega}_{=: E} = 0 \quad \Longrightarrow \quad \left( \mathbf{u}, \frac{d\mathbf{u}}{dt} \right) = 0$$



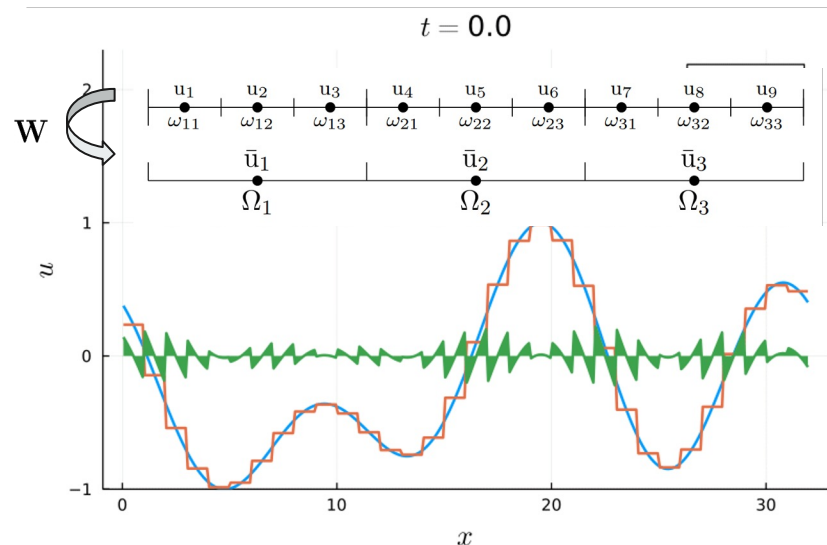
# Discrete filtering and reconstruction

- Spatial filter  $W$ :

$$\bar{u} = W u$$

- Subgrid-scales defined via reconstruction operator  $R$ :

$$u' = u - R \bar{u}$$



## Energy decomposition

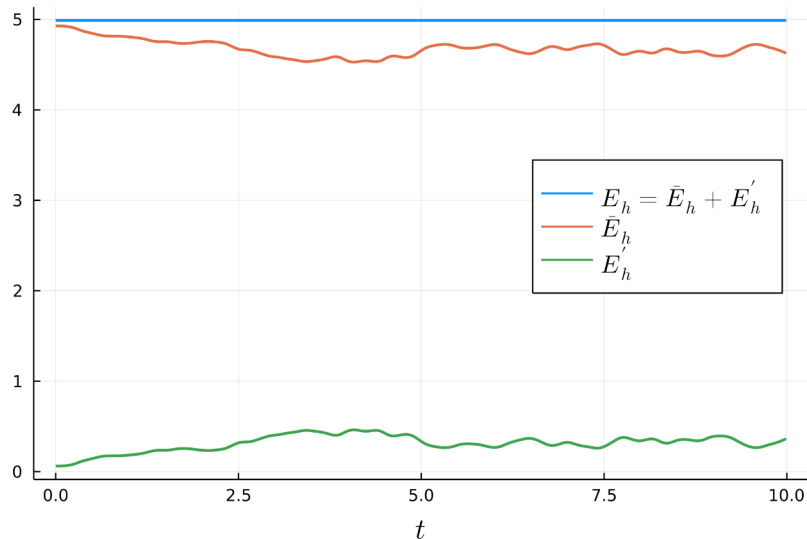
- Decompose the energy as:

$$E_h = \underbrace{\frac{1}{2}(\bar{\mathbf{u}}, \bar{\mathbf{u}})_\Omega}_{=:\bar{E}_h} + \underbrace{\frac{1}{2}(\mathbf{u}', \mathbf{u}')_\omega}_{=:E'_h}$$

- Time evolution:

$$\frac{dE_h}{dt} = \boxed{\frac{d\bar{E}_h(\bar{\mathbf{u}})}{dt}} + \boxed{\frac{dE'_h(\mathbf{u}')}{dt}} = 0$$

- To use energy conservation we need information about the small scales



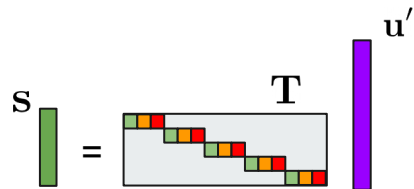
Total energy is conserved

Energy of filtered solution is not conserved



## Subgrid compression

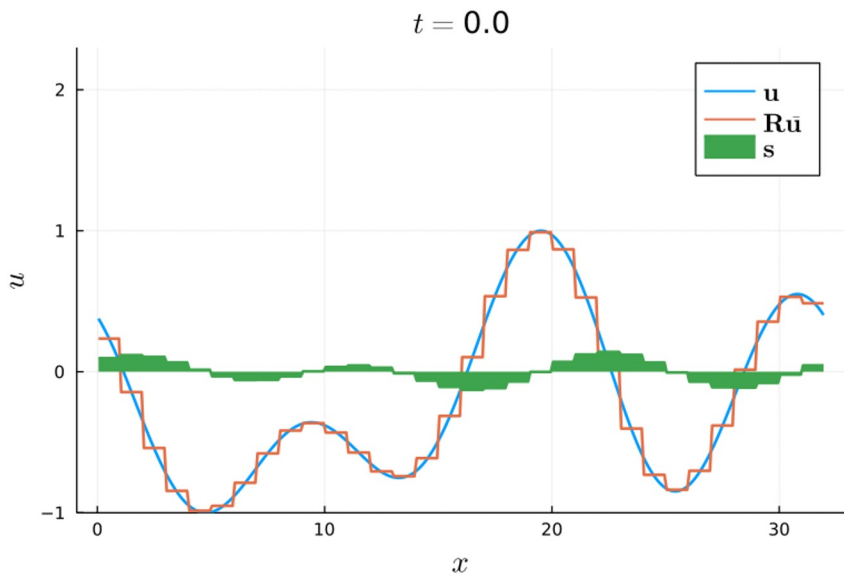
- Introduce (linear) compressed subgrid variables  $\mathbf{s}$ :
- Require  $\mathbf{s}$  to have same energy as  $\mathbf{u}'$ :
- Solve minimization problem (“local POD”):



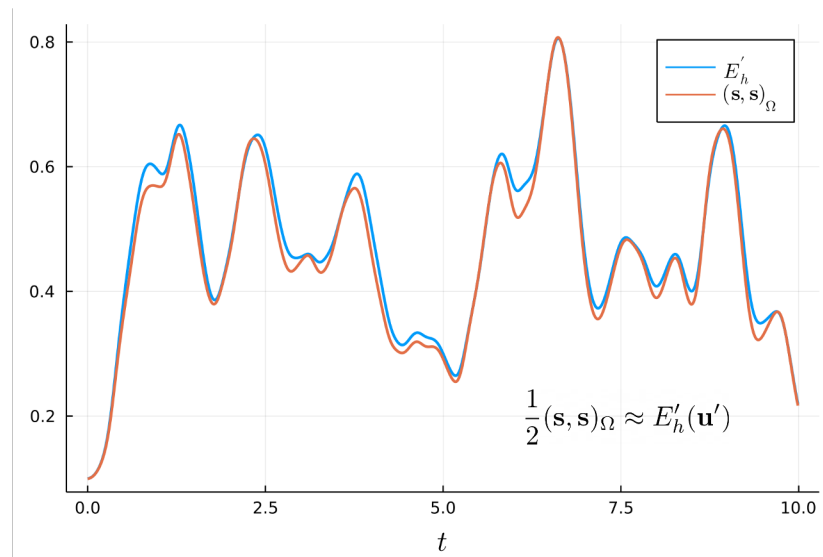
$$\frac{1}{2}(\mathbf{s}, \mathbf{s})_{\Omega} \approx E'_h(\mathbf{u}')$$

$$\begin{matrix} \text{green} & \text{orange} & \text{red} \\ \text{green} & \text{orange} & \text{red} \end{matrix} = \arg \min_{\begin{matrix} \text{green} & \text{orange} & \text{red} \\ \text{green} & \text{orange} & \text{red} \end{matrix}} \sum_{d=1}^{\mathcal{D}} \left\| \frac{1}{2} \mathbf{s}_d^2 - \frac{1}{2} \mathbf{W}(\mathbf{u}'_d)^2 \right\|_2^2$$

# Compressed variables learn effective subgrid content



compressed subgrid variable identifies sharp gradients



learned compression matches subgrid energy closely

$$\frac{d\bar{\mathbf{u}}}{dt} = f(\bar{\mathbf{u}}) + \underbrace{f(\mathbf{u}) - f(\bar{\mathbf{u}})}_{\approx c(\bar{\mathbf{u}}; \theta)}$$

## Energy-conserving closure model

- Extended neural closure model:

$$\frac{d}{dt} \begin{bmatrix} \bar{\mathbf{u}} \\ \mathbf{s} \end{bmatrix} = \begin{bmatrix} f(\bar{\mathbf{u}}) \\ \mathbf{0} \end{bmatrix} + \begin{bmatrix} c_u(\bar{\mathbf{u}}, \mathbf{s}; \theta_u) \\ c_s(\bar{\mathbf{u}}, \mathbf{s}; \theta_s) \end{bmatrix}$$

Large scale dynamics

Compressed small  
scale dynamics

- Idea: learn a skew-symmetric matrix  $\mathcal{K}$  whose entries are NN outputs

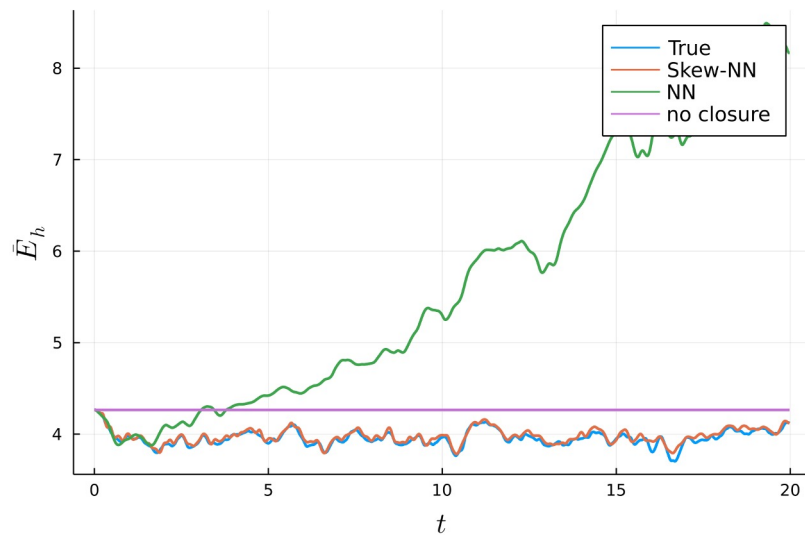
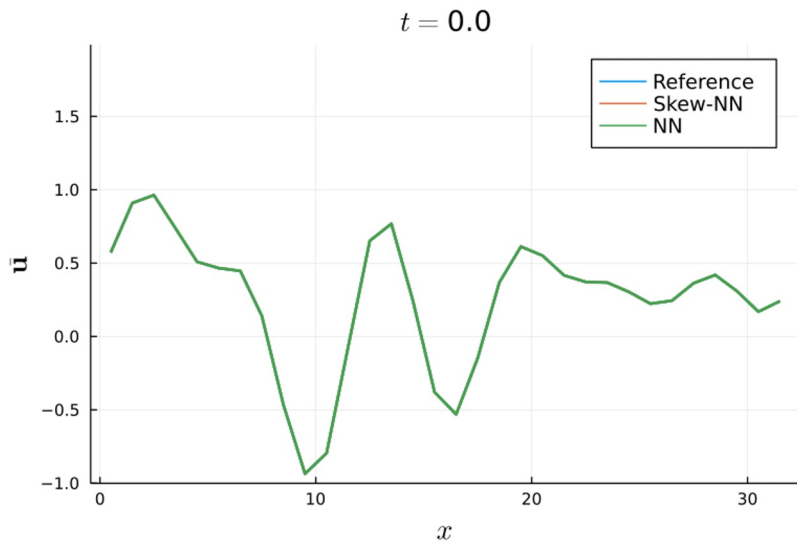
$$\begin{bmatrix} c_u(\bar{\mathbf{u}}, \mathbf{s}; \theta_u) \\ c_s(\bar{\mathbf{u}}, \mathbf{s}; \theta_s) \end{bmatrix} = \mathcal{K}(\bar{\mathbf{u}}, \mathbf{s}; \Theta) \begin{bmatrix} \bar{\mathbf{u}} \\ \mathbf{s} \end{bmatrix}$$

- Energy conservation:

$$\begin{bmatrix} \bar{\mathbf{u}} \\ \mathbf{s} \end{bmatrix}^T \mathcal{K}(\cdot) \begin{bmatrix} \bar{\mathbf{u}} \\ \mathbf{s} \end{bmatrix} = 0 \quad \Longrightarrow \quad \frac{d\bar{E}_h(\bar{\mathbf{u}})}{dt} + \frac{1}{2} \frac{d(\mathbf{s}, \mathbf{s})_\omega}{dt} = 0$$

## New closure model improves quality + stability

- Trained on different initial conditions, tested on unseen initial conditions
- Reduction from  $N = 600$  to  $N = 30$
- Compare to standard CNN

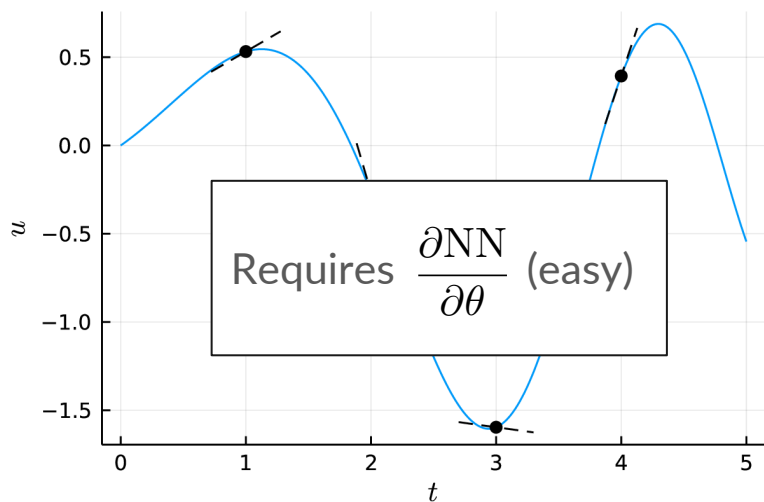




What about training neural closure models?

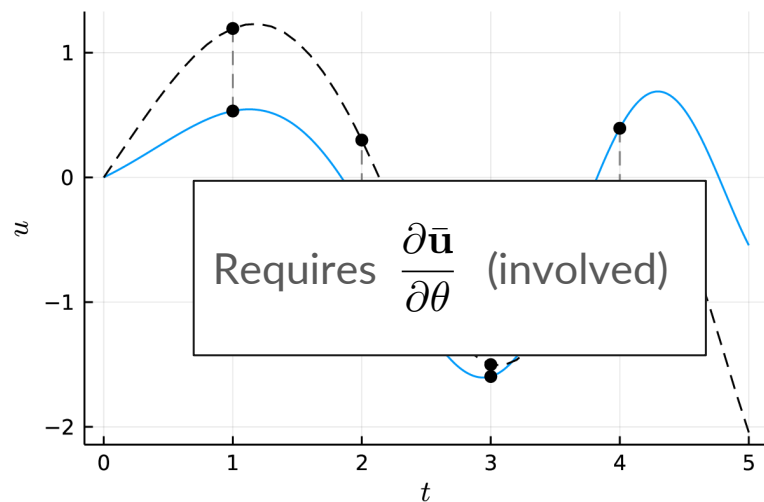
$$\frac{d\bar{\mathbf{u}}}{dt} = f(\bar{\mathbf{u}}) + \text{NN}(\bar{\mathbf{u}}; \theta)$$

## Derivative fitting



$$\text{Loss} = \left\| \left( \frac{d\bar{\mathbf{u}}}{dt} \right)_{\text{ref}} - f(\bar{\mathbf{u}}_{\text{ref}}) - \text{NN}(\bar{\mathbf{u}}_{\text{ref}}; \theta) \right\|^2$$

## Trajectory fitting



$$\text{Loss} = \sum_{i=1}^{N_t} \|\bar{\mathbf{u}}_{\text{ref}}(t_i) - \bar{\mathbf{u}}(t_i)\|^2, \text{ where } \frac{d\bar{\mathbf{u}}}{dt} = f(\bar{\mathbf{u}}) + \text{NN}(\bar{\mathbf{u}}; \theta)$$

$$\text{Loss} = \left\| \left( \frac{d\bar{\mathbf{u}}}{dt} \right)_{\text{ref}} - f(\bar{\mathbf{u}}_{\text{ref}}) - \text{NN}(\bar{\mathbf{u}}_{\text{ref}}; \theta) \right\|^2$$

## Derivative fitting: theoretical results

**Theorem 3.1** (Hairer et al. [11]). Let  $\mathbf{u}_{\text{ref}}(t) \in \mathbb{R}^{N_x}, t \geq 0$  be given, let  $\mathbf{u}(t) \in \mathbb{R}^{N_x}, t \geq 0$  be the solution of the ODE  $\frac{d\mathbf{u}}{dt} = g(\mathbf{u}; \vartheta)$  with initial condition  $\mathbf{u}(0) = \mathbf{u}_{\text{ref}}(0)$ , and let  $\|\cdot\|$  be a norm on  $\mathbb{R}^{N_x}$ . If the following holds:

- $\left\| \frac{d}{dt} \mathbf{u}_{\text{ref}}(t) - g(\mathbf{u}_{\text{ref}}(t); \vartheta) \right\| \leq \varepsilon,$
- $\|g(\mathbf{a}; \vartheta) - g(\mathbf{b}; \vartheta)\| \leq C \|\mathbf{a} - \mathbf{b}\|,$

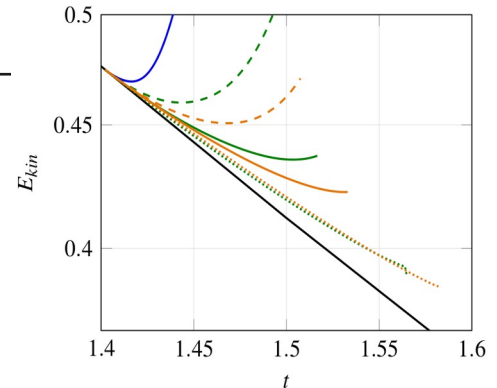
for fixed Lipschitz constant  $C > 0$  and fixed  $\varepsilon > 0$ . Then the following error bound holds:

$$\|\mathbf{u}_{\text{ref}}(t) - \mathbf{u}(t)\| \leq \frac{\varepsilon}{C} (e^{Ct} - 1).$$

If a neural ODE:

- approximates the derivative well
- is Lipschitz

Then, the resulting solution may be inaccurate.



$$\text{Loss} = \sum_{i=1}^{N_t} \|\bar{\mathbf{u}}(t_i) - \mathbf{u}(t_i)\|^2, \text{ where } \frac{d\bar{\mathbf{u}}}{dt} = f(\bar{\mathbf{u}}) + \text{NN}(\bar{\mathbf{u}}; \theta)$$

## Trajectory fitting: theoretical results

**Theorem 3.2.** Let  $\mathbf{u}_{\text{ref}}(t_i), i = 0, 1, \dots, \dots$  be a sequence of vectors in  $\mathbb{R}^{N_x}$  where  $t_i = i\Delta t$ , let  $\|\cdot\|$  be a norm on  $\mathbb{R}^{N_x}$ , and let  $G(\cdot; \vartheta) : \mathbb{R}^{N_x} \rightarrow \mathbb{R}^{N_x}$  be a function such that:

- a)  $\|\mathbf{u}_{\text{ref}}(t_{i+1}) - G(\mathbf{u}_{\text{ref}}(t_i); \vartheta)\| \leq \varepsilon$  for all  $i = 1, 2, \dots, N_t$ ,
- b)  $\|G(\mathbf{a}; \vartheta) - G(\mathbf{b}; \vartheta)\| \leq C \|\mathbf{a} - \mathbf{b}\|$  for all  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^{N_x}$ ,

for fixed Lipschitz constant  $C > 0, C \neq 1$  and fixed  $\varepsilon > 0$ . Define the sequence  $\mathbf{u}(t_{i+1}) = G(\mathbf{u}(t_i); \vartheta)$  with  $\mathbf{u}(0) = \mathbf{u}_{\text{ref}}(0)$ . Then the following error bound holds:

$$\|\mathbf{u}(t_k) - \mathbf{u}_{\text{ref}}(t_k)\| \leq \varepsilon \frac{C^k - 1}{C - 1} \text{ for } k = 0, 1, 2, \dots$$

If a discretised neural ODE:

- accurately represents single time steps
- is Lipschitz

Then the resulting ODE solution may be inaccurate.

Trajectory length should be sufficiently large



$$\text{Loss} = \sum_{i=1}^{N_t} \|\bar{\mathbf{u}}(t_i) - \mathbf{u}(t_i)\|^2, \text{ where } \frac{d\bar{\mathbf{u}}}{dt} = f(\bar{\mathbf{u}}) + \text{NN}(\bar{\mathbf{u}}; \theta)$$

## Trajectory fitting: computing $\frac{d\text{Loss}}{d\theta}$

### 1. Discretise-then-optimize (DtO):

- Need differentiable solver
- Preferably explicit

### 2. Optimize-then-discretize (OtD):

- Solve adjoint equations<sup>1</sup>
- Here: interpolating adjoint
- Need dense output

$$\begin{cases} \frac{d}{dt} \mathbf{y}^T = -\mathbf{y}^T \frac{\partial}{\partial \bar{\mathbf{u}}} (f(\bar{\mathbf{u}}) + \text{NN}(\bar{\mathbf{u}}; \theta)) \\ \frac{d}{dt} \mathbf{z}^T = -\mathbf{y}^T \frac{\partial}{\partial \theta} (f(\bar{\mathbf{u}}) + \text{NN}(\bar{\mathbf{u}}; \theta)) \\ \frac{d\text{Loss}}{d\theta} = \mathbf{z}(0) \end{cases}$$

<sup>1</sup>Chen et al, Neural ordinary differential equations, NeurIPS 2018

## Derivative fitting vs. trajectory fitting

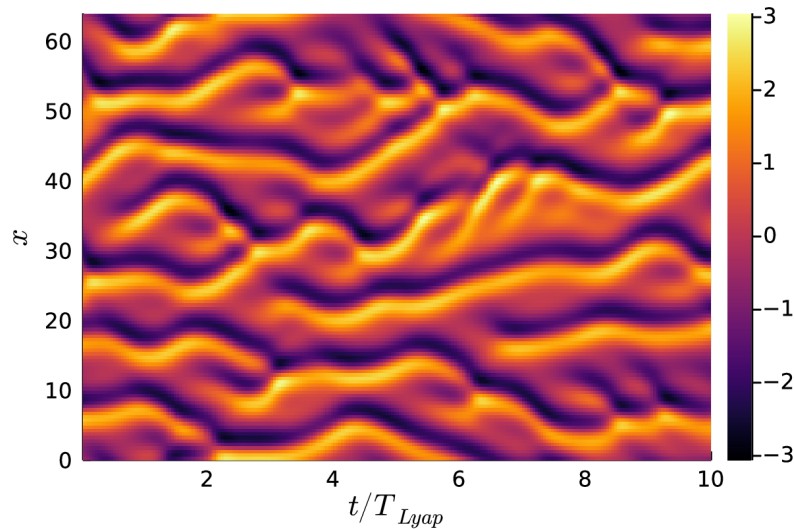
	Derivative fitting	Trajectory fitting	
		DtO	OtD
Differentiability required	NN	NN, $f$ , ODE solver	NN, $f$
Accuracy of loss function gradients	Exact	Exact	Approximate
Learns long-term accuracy	No	Yes	Yes
Requires time-derivatives of training data	Yes	No	No
Computational cost	Low	High	High

# Kuramoto-Sivashinsky equation

- **Chaotic:**
  - Trajectory lengths not too large
- **Stiff:**
  - OtD: impl/expl RK, KenCarp47<sup>1</sup>
  - DtO: expl ETDRK4 in Fourier domain<sup>2</sup>
- **Filter  $W$ :**
  - downsampling 1024 => 128

$$\frac{\partial u}{\partial t} = -\frac{1}{2} \frac{\partial}{\partial x} (u^2) - \frac{\partial^2 u}{\partial x^2} - \frac{\partial^4 u}{\partial x^4}$$

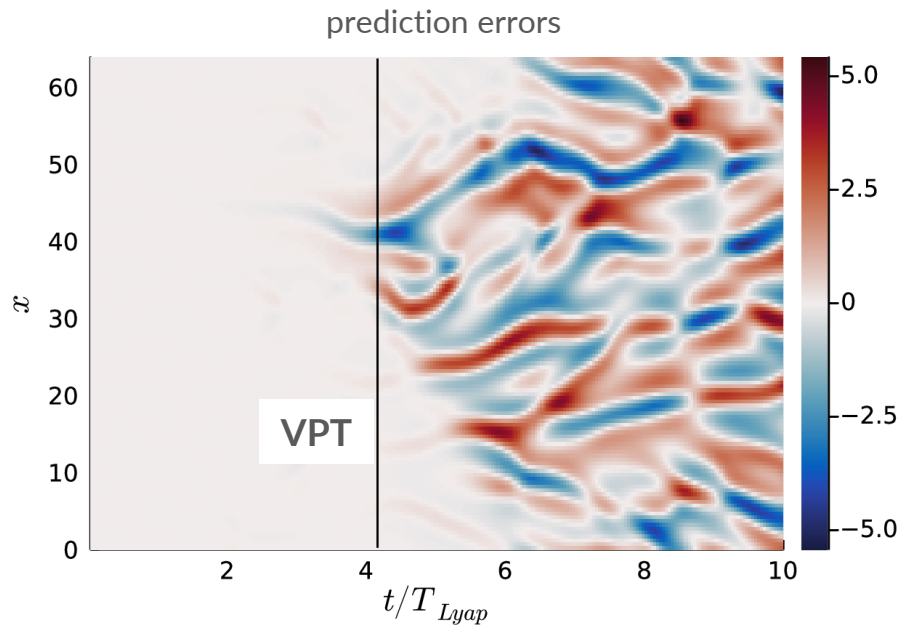
$$\frac{d\bar{u}}{dt} = f(\bar{u}) + \Delta_{\text{fwd}} \text{NN}(\bar{u}; \theta)$$



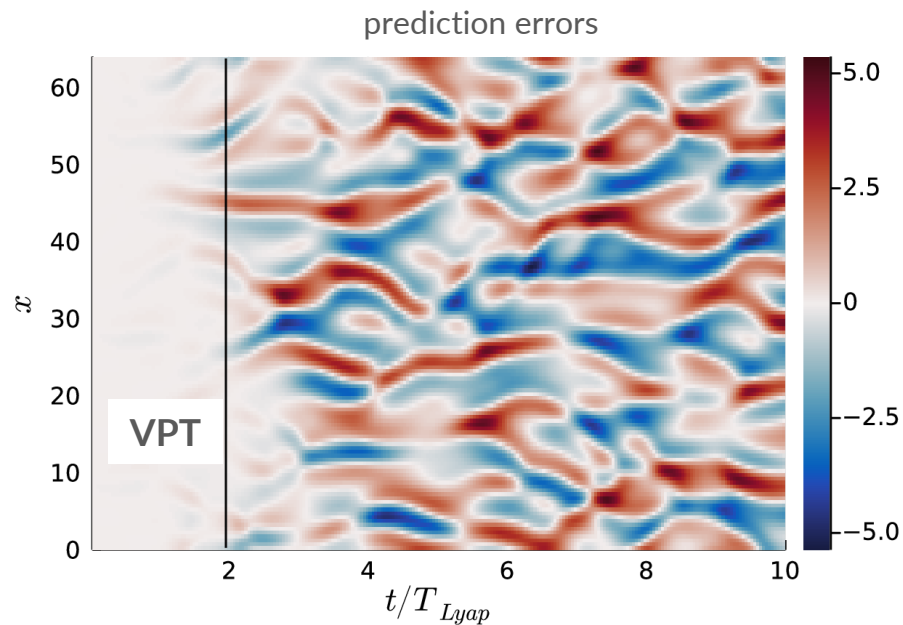
<sup>1</sup>Kennedy and Carpenter, Higher-order additive Runge-Kutta schemes for ODEs, Applied Numerical Mathematics, 2019.

<sup>2</sup>Kassam & Trefethen, Fourth-order time-stepping for stiff PDEs. SIAM Journal on Scientific Computing, 2005

## Effect of trajectory length, OtD

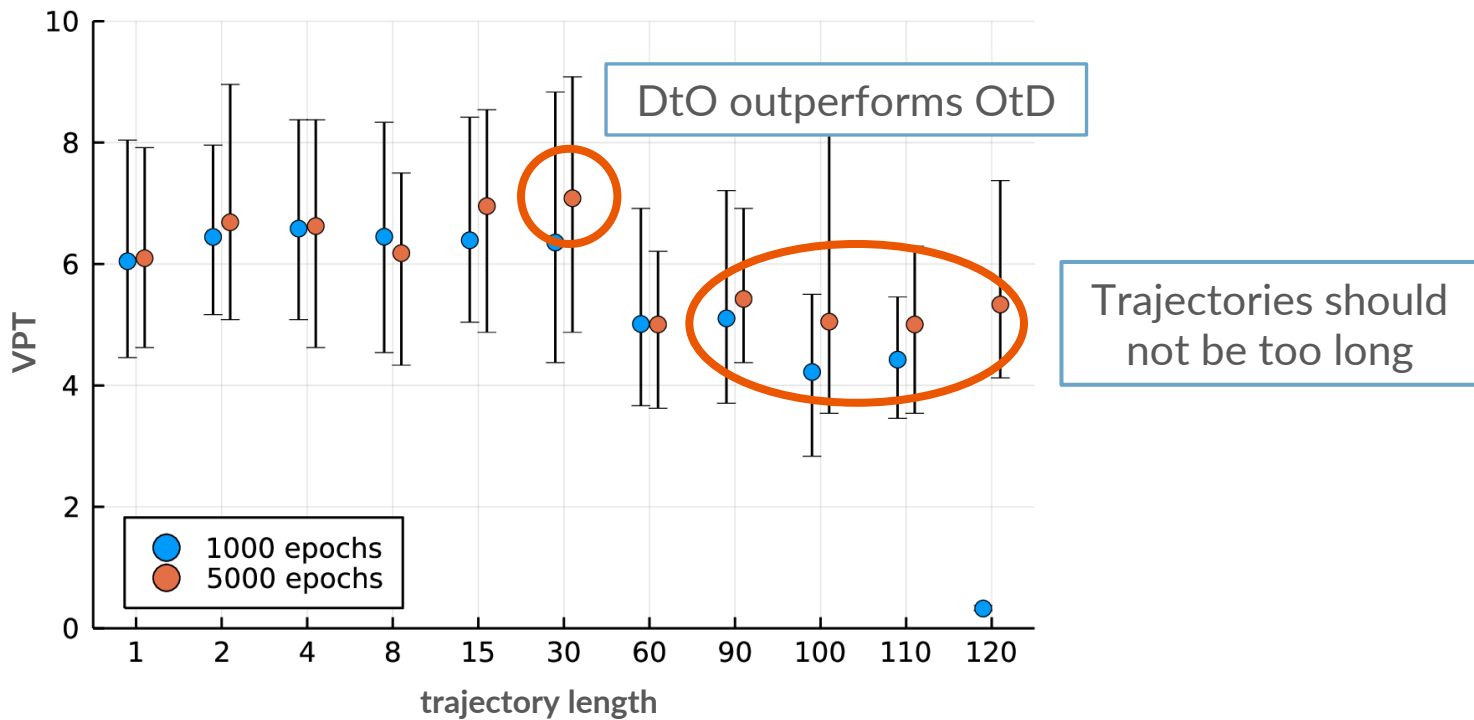


short trajectories (24 steps)

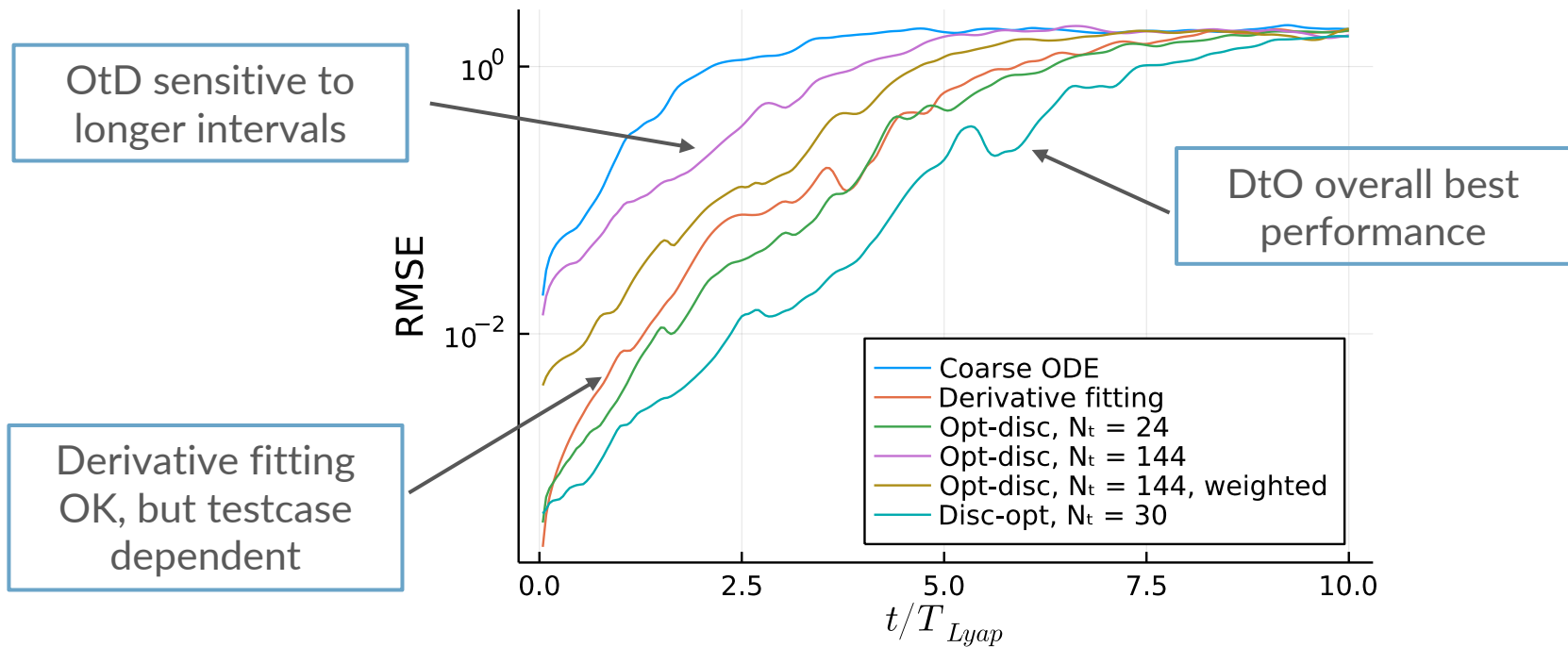


long trajectories (144 steps)

## Effect of trajectory length, DtO



## DtO outperforms OtD and derivative fitting



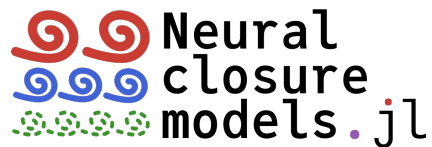


## Conclusions

- **“Discretize first”**
  - Tailor-made closure models
  - Useful framework for NNs, eases analysis
- **“Preserve structure”**
  - Accuracy improvement by adding physics knowledge
  - Non-linear stability with energy-conserving methods
- **“Embedded learning” with trajectory fitting**
  - Discretise-then-optimise with **differentiable solvers** preferred
  - Promising, but “strings attached”: problem-dependent, comparison not easy

# Julia is great for differentiable programming

- Neural closure models
  - <https://github.com/HugoMelchers/neural-closure-models>
- Incompressible, energy-conserving Navier-Stokes code
  - <https://github.com/agdestein/IncompressibleNavierStokes.jl>
- DifferentialEquations.jl by Rackauckas et al.
  - <https://sciml.ai>





# References

## Energy-Conserving Neural Network for Turbulence Closure Modeling

T. van Gastelen<sup>a</sup>, W. Edeling<sup>a</sup>, B. Sanderse<sup>a</sup>

<sup>a</sup>Centrum Wiskunde & Informatica, Science Park 123, Amsterdam, The Netherlands

### Abstract

In turbulence modeling, and more particularly in the Large-Eddy Simulation (LES) framework, we are concerned with finding closure models that represent the effect of the unresolved subgrid scales on the resolved scales. Recent approaches gravitate towards machine learning techniques to construct such models. However, the stability of machine-learned closure models and their abidance by physical structure (e.g. symmetries, conservation laws) are still open problems. To tackle both issues, we take the 'discretize first, filter next' approach, in which we apply a spatial averaging filter to existing energy-conserving (fine-grid) discretizations. The main novelty is that we extend the system of equations describing the filtered solution with a set of equations that describe the evolution of (a compressed version of) the energy of the subgrid scales. Having an estimate of the energy of the subgrid scales, we can use the concept of energy conservation and derive stability of the discrete representation. The compressed variables are determined via a data-driven technique in such a way that the energy of the subgrid scales is matched. For the extended system, the closure model should be energy-conserving, and a new skew-symmetric convolutional neural network architecture is proposed that has this property. Stability is thus guaranteed, independent of the actual weights and biases of the network. Importantly, our framework allows energy exchange between resolved scales and compressed subgrid scales and thus enables backscatter. To model dissipative systems (e.g. viscous flows), the framework is extended with a diffusive component. The introduced neural network architecture is constructed such that it also satisfies momentum conservation. We apply the new methodology to both the viscous Burgers' equation and the Korteweg-De Vries equation in 1D and show superior stability properties when compared to a vanilla convolutional neural network.

Computers and Mathematics with Applications 143 (2023) 94–107

Contents lists available at ScienceDirect



Computers and Mathematics with Applications

journal homepage: [www.elsevier.com/locate/camwa](http://www.elsevier.com/locate/camwa)



## Comparison of neural closure models for discretised PDEs

Hugo Melchers<sup>a,c,1</sup>, Daan Crommelin<sup>a,b</sup>, Barry Koren<sup>c</sup>, Vlado Menkovski<sup>c</sup>, Benjamin Sanderse<sup>a,4</sup>

<sup>a</sup>Centrum Wiskunde & Informatica, Science Park 123, 1098 XG, Amsterdam, the Netherlands

<sup>b</sup>Korteweg-de Vries Institute for Mathematics, University of Amsterdam, Science Park 105-107, 1098 XG, Amsterdam, the Netherlands

<sup>c</sup>Eindhoven University of Technology, De Zaal, 5600 MB, Eindhoven, the Netherlands

### ARTICLE INFO

Dataset link: <https://github.com/HugoMelchers/neural-closure-models>

Keywords:  
Ordinary differential equations  
Neural networks  
Neural ODE  
Partial differential equations  
Multiscale modelling  
Closure model

### ABSTRACT

Neural closure models have recently been proposed as a method for efficiently approximating small scales in multiscale systems with neural networks. The choice of loss function and associated training procedure has a large effect on the accuracy and stability of the resulting neural closure model. In this work, we systematically compare three distinct procedures: "derivative fitting", "trajectory fitting" with discrete-then-optimize, and "trajectory fitting" with optimize-then-discrete. Derivative fitting is conceptually the simplest and computationally the most efficient approach and is found to perform reasonably well on one of the test problems (Kuramoto-Sivashinsky) but poorly on the other (Burgers). Trajectory fitting is computationally more expensive but is more robust and is therefore the preferred approach. Of the two trajectory fitting procedures, the discrete-then-optimize approach produces more accurate models than the optimize-then-discrete approach. While the optimize-then-discrete approach can still produce accurate models, care must be taken in choosing the length of the trajectories used for training, in order to train the models on long-term behaviour while still producing reasonably accurate gradients during training. Two existing theorems are interpreted in a novel way that gives insight into the long-term accuracy of a neural closure model based on how accurate it is in the short term.