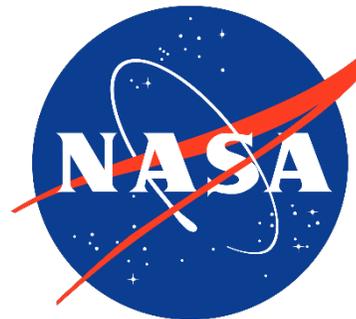


# New Perspectives on Balancing Physics with Data-Driven Models: The Case for Physics-Informed Neural Networks in Environmental Statistics

Stefano Castruccio  
University of Notre Dame

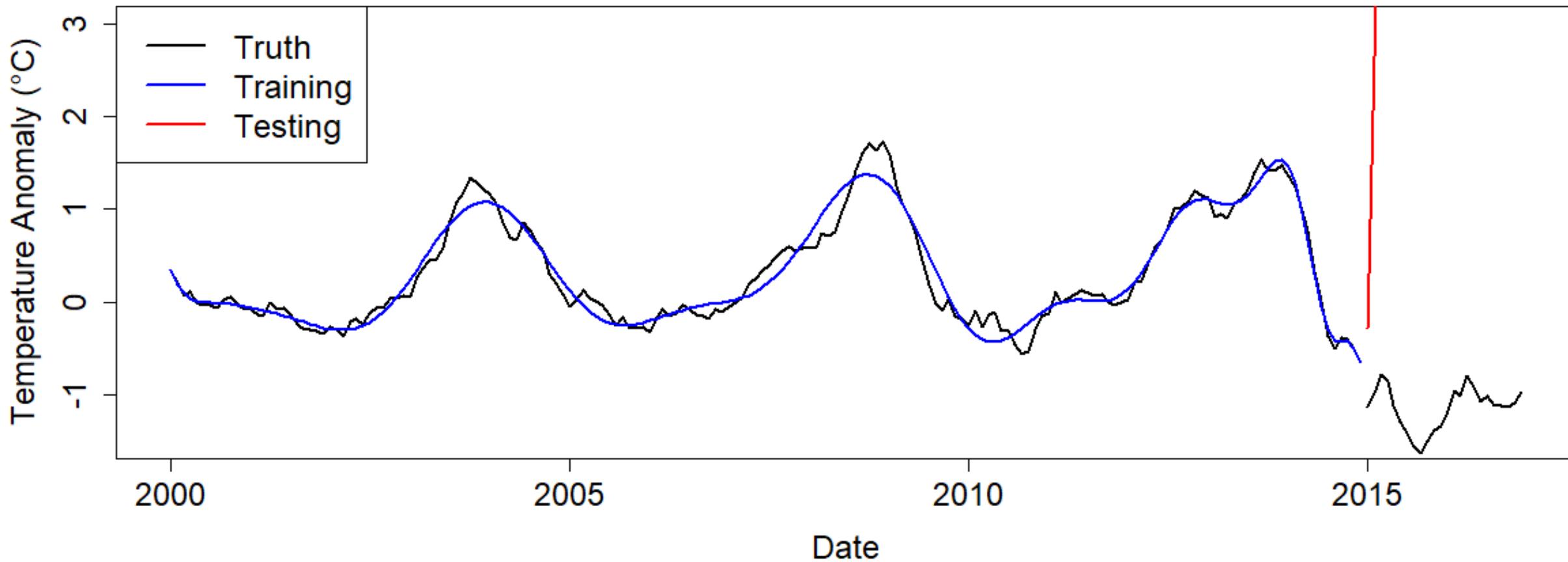
Research supported by:



# Problem- or data-driven?

- Let's start with a super-generic question
- **For a particular scientific problem, do we need more context, or should the data “speak for themselves”?**
- Data-driven: rely on **non-parametric constructs (neural networks, NNs)** have been very popular for many applications in STEM and beyond
- Shall we just ignore the context?

# Why data driven methods could fail

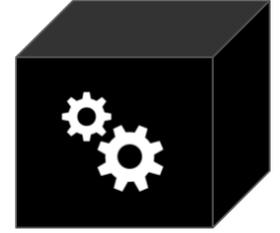


$$\frac{\partial W}{\partial t} + \nabla \cdot Q = E - P$$

$$F_\lambda = \int_{i_{West}} Q_{\lambda_{West}} \cdot dl \text{ and } F_\phi = \int_{i_{East}} Q_{\phi_{East}} \cdot dl$$

$$\frac{\partial \rho F(x)}{\partial x} + \frac{\partial \rho F(y)}{\partial y} = E - \rho P$$

# A fork in the road?



- Well known equations
- Few data (hard/expensive)
- E.g., astronomy, paleoclimate



Unknown or nonexistent equations

A lot of data (cheap and ubiquitous)

E.g., social networks, large language models, foundation models

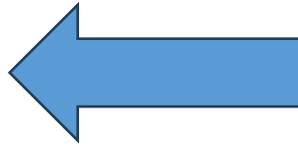
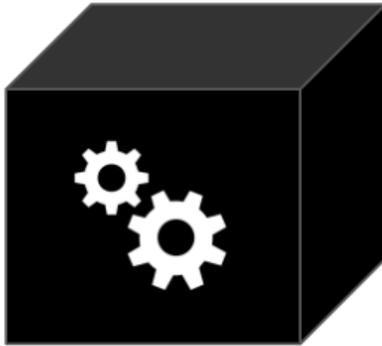
# Do we even have to choose?

- The problem of merging (assimilating, fusing, etc.) data with physical information is a very old exploit in Statistics
- These models are generally termed **physical-statistical** (PS) models
- Mark Berliner's 2003 JRG paper was the first one to formalize it (to my knowledge!)

*Spatio-temporal statistical models are not at odds with deterministic ones. Indeed, the most powerful models are constructed based on physical mechanisms*

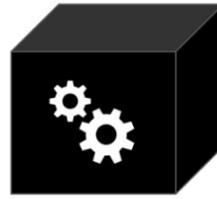
(Wikle, Zammit-Magion and Cressie, Spatio-temporal statistics with R)

# How PS models work



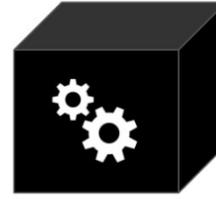
$$\frac{\partial W}{\partial t} + \nabla \cdot Q = E - P$$
$$F_\lambda = \int_{l_{West}} Q_{\lambda_{West}} \cdot dl \text{ and } F_\phi = \int_{l_{East}} Q_{\phi_{East}} \cdot dl$$
$$\frac{\partial \rho F^{(x)}}{\partial x} + \frac{\partial \rho F^{(y)}}{\partial y} = E - \rho P$$

# PS models



- PS model fits naturally in a hierarchical (Bayesian) framework
- Stage 1 (data): Suppose we observe  $Z_t(\mathbf{s})$ , we assume
$$Z_t(\mathbf{s}) = \mathcal{H} \left( Y_t(\mathbf{s}), \boldsymbol{\theta}^{\text{data}}, \varepsilon_t(\mathbf{s}) \right)$$
- $Y_t(\mathbf{s})$ : latent process
- $\mathcal{H}$ : (linear/nonlinear) mapping
- $\boldsymbol{\theta}^{\text{data}}$ : data-model parameters
- $\varepsilon_t(\mathbf{s})$ : data-model error
- Easy example  $Z_t(\mathbf{s}) \sim N(Y_t(\mathbf{s}), \sigma^2)$

# PS models



- Stage 2: The latent process  $Y_t(\mathbf{s})$  is modeled

$$Y_t(\mathbf{s}) = \mathcal{M}(\{Y_{t-1}(\mathbf{s}), \dots, Y_{t-k}(\mathbf{s})\}, \boldsymbol{\theta}^{\text{process}}, \eta_t(\mathbf{s}))$$

- $Y_t(\mathbf{s})$ : latent process
- $\mathcal{M}$ : dynamic process
- $\boldsymbol{\theta}^{\text{process}}$  : data-model parameters
- $\eta_t(\mathbf{s})$ : process-model error
- (Stage 3: Priors)



$$\begin{aligned}
 r: & \rho \left( \frac{\partial u_r}{\partial t} + u_r \frac{\partial u_r}{\partial r} + \frac{u_\phi}{r \sin(\theta)} \frac{\partial u_r}{\partial \phi} + \frac{u_\theta}{r} \frac{\partial u_r}{\partial \theta} - \frac{u_\phi^2 + u_\theta^2}{r} \right) = -\frac{\partial p}{\partial r} + \rho g_r + \\
 & \mu \left[ \frac{1}{r^2} \frac{\partial}{\partial r} \left( r^2 \frac{\partial u_r}{\partial r} \right) + \frac{1}{r^2 \sin(\theta)^2} \frac{\partial^2 u_r}{\partial \phi^2} + \frac{1}{r^2 \sin(\theta)} \frac{\partial}{\partial \theta} \left( \sin(\theta) \frac{\partial u_r}{\partial \theta} \right) - 2 \frac{u_r + \frac{\partial u_\theta}{\partial \theta} + u_\theta \cot(\theta)}{r^2} - \frac{2}{r^2 \sin(\theta)} \frac{\partial u_\phi}{\partial \phi} \right] \\
 \phi: & \rho \left( \frac{\partial u_\phi}{\partial t} + u_r \frac{\partial u_\phi}{\partial r} + \frac{u_\phi}{r \sin(\theta)} \frac{\partial u_\phi}{\partial \phi} + \frac{u_\theta}{r} \frac{\partial u_\phi}{\partial \theta} + \frac{u_r u_\phi + u_\phi u_\theta \cot(\theta)}{r} \right) = -\frac{1}{r \sin(\theta)} \frac{\partial p}{\partial \phi} + \rho g_\phi + \\
 & \mu \left[ \frac{1}{r^2} \frac{\partial}{\partial r} \left( r^2 \frac{\partial u_\phi}{\partial r} \right) + \frac{1}{r^2 \sin(\theta)^2} \frac{\partial^2 u_\phi}{\partial \phi^2} + \frac{1}{r^2 \sin(\theta)} \frac{\partial}{\partial \theta} \left( \sin(\theta) \frac{\partial u_\phi}{\partial \theta} \right) + \frac{2 \sin(\theta) \frac{\partial u_r}{\partial \phi} + 2 \cos(\theta) \frac{\partial u_\theta}{\partial \phi} - u_\phi}{r^2 \sin(\theta)^2} \right] \\
 \theta: & \rho \left( \frac{\partial u_\theta}{\partial t} + u_r \frac{\partial u_\theta}{\partial r} + \frac{u_\phi}{r \sin(\theta)} \frac{\partial u_\theta}{\partial \phi} + \frac{u_\theta}{r} \frac{\partial u_\theta}{\partial \theta} + \frac{u_r u_\theta - u_\phi^2 \cot(\theta)}{r} \right) = -\frac{1}{r} \frac{\partial p}{\partial \theta} + \rho g_\theta + \\
 & \mu \left[ \frac{1}{r^2} \frac{\partial}{\partial r} \left( r^2 \frac{\partial u_\theta}{\partial r} \right) + \frac{1}{r^2 \sin(\theta)^2} \frac{\partial^2 u_\theta}{\partial \phi^2} + \frac{1}{r^2 \sin(\theta)} \frac{\partial}{\partial \theta} \left( \sin(\theta) \frac{\partial u_\theta}{\partial \theta} \right) + \frac{2}{r^2} \frac{\partial u_r}{\partial \theta} - \frac{u_\theta + 2 \cos(\theta) \frac{\partial u_\phi}{\partial \phi}}{r^2 \sin(\theta)^2} \right].
 \end{aligned}$$

# Relative merits of PS models

- PS models

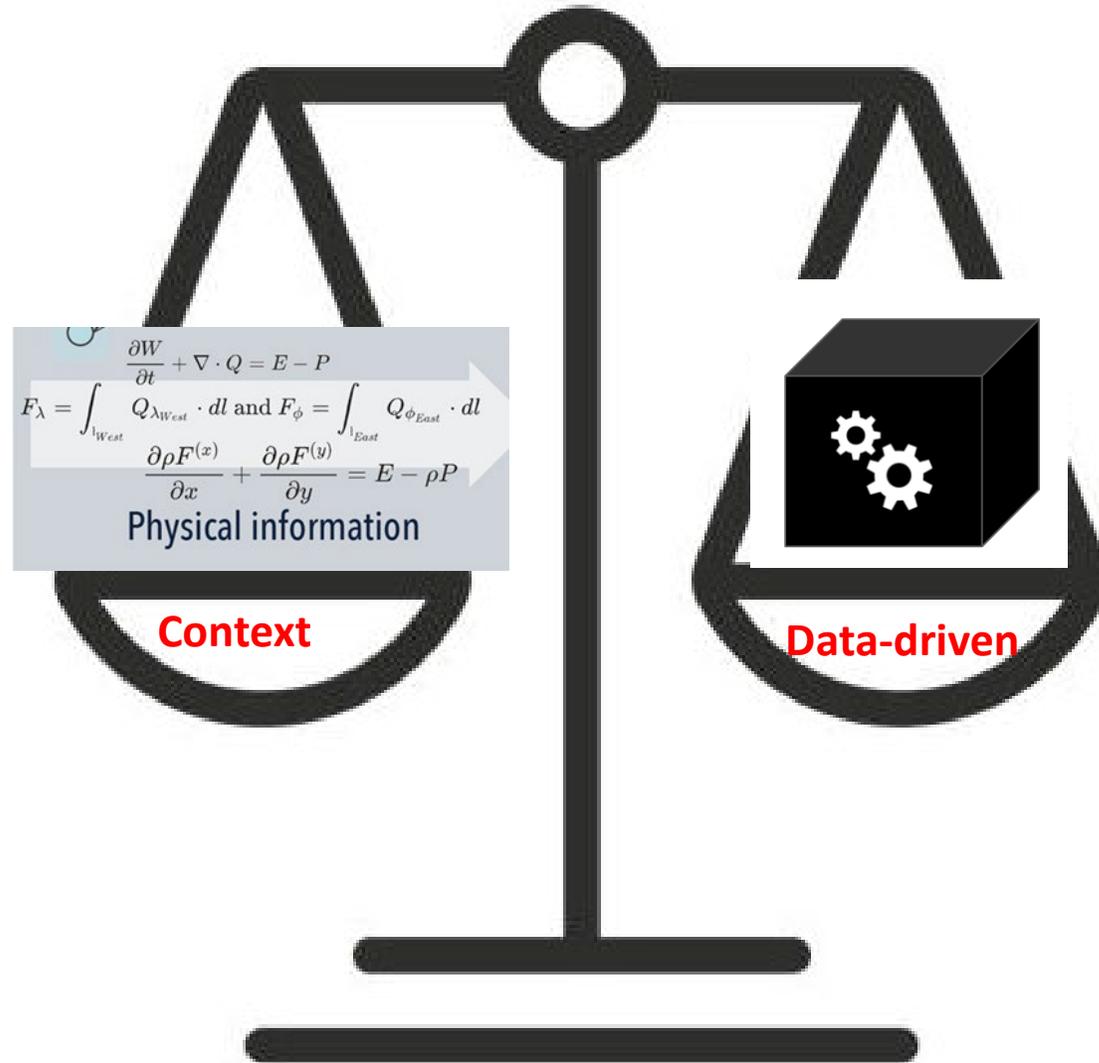
$$Y_t(\mathbf{s}) = \mathcal{M}(\{Y_{t-1}(\mathbf{s}), \dots, Y_{t-k}(\mathbf{s})\}, \boldsymbol{\theta}^{\text{process}}, \eta_t(\mathbf{s}))$$

- PS models require a definition of  $\mathcal{M}$ , so an **exact knowledge of the physics**
- If something goes wrong, we hope that  $\eta_t(\mathbf{s})$  takes care of any misspecification

# Physics-informed Neural networks

- In the last few years, the ML literature has developed a separate area of research on how to merge physics with data: **physics informed neural networks (PINNs)**
- The idea is quite different, but there is a link with PS models which was somehow lost

# The main idea: a scale



# PINNs in a nutshell

- **This is my “stats” reframing of PINN**
- The canonical PINN definition is more algorithmical
- Stage 1: We observe  $Z_t(\mathbf{s})$  and we assume

$$Z_t(\mathbf{s}) = \mathcal{H} \left( Y_t(\mathbf{s}), \boldsymbol{\theta}^{\text{data}}, \varepsilon_t(\mathbf{s}) \right)$$

- $Y_t(\mathbf{s})$ : latent process
- $\mathcal{H}$ : (linear/nonlinear) mapping
- $\boldsymbol{\theta}^{\text{data}}$ : data-model parameters
- $\varepsilon_t(\mathbf{s})$ : data-model error (noise)

# PINNs as hierarchical models

- Stage 2 (process): The latent process  $Y_t(\mathbf{s})$  is modeled
$$Y_t(\mathbf{s}) = \mathcal{M}(\{Y_{t-1}(\mathbf{s}), \dots, Y_{t-k}(\mathbf{s})\}, \boldsymbol{\theta}^{\text{process}}, \eta_t(\mathbf{s}))$$
- $Y_t(\mathbf{s})$ : latent process
- $\mathcal{M}$ : dynamic process, **highly nonparametric NN**
- $\boldsymbol{\theta}^{\text{process}}$ : data-model parameters (**huge space!**)
- $\eta_t(\mathbf{s})$ : process-model error (noise)
- (Stage 3: Priors)

# The main idea behind PINNs

- Assume  $\mathbf{Y}_t = \{Y_t(\mathbf{s})\}$  is informed by a PDE

$$\frac{\partial \mathbf{Y}_t}{\partial t} - \mathcal{N}[\mathbf{Y}_t] = 0$$

- **We are not solving the PDE**, we just want  $\mathbf{Y}_t$  to be “loosely compliant”

- So,  $g(\mathbf{Y}_t) = \frac{\partial \mathbf{Y}_t}{\partial t} - \mathcal{N}[\mathbf{Y}_t]$  should be small

- Physics-based models:

$$\sum_{t=1}^T \mathcal{C}(\mathbf{z}_t, \hat{\mathbf{z}}_t(\boldsymbol{\theta}^{\text{process}})) + \lambda g(\hat{\mathbf{Y}}_t(\boldsymbol{\theta}^{\text{process}}))^2$$

- For some cost function  $\mathcal{C}$

# The main idea behind PINNs

- We minimize

$$\sum_{t=1}^T C(\mathbf{z}_t, \hat{\mathbf{z}}_t(\boldsymbol{\theta}^{\text{process}})) + \lambda g(\hat{\mathbf{Y}}_t(\boldsymbol{\theta}^{\text{process}}))^2$$

- If we have a Gaussian process (“traditional” PINN)  $\hat{\mathbf{z}}_t = \hat{\mathbf{Y}}_t$  and we have

$$\sum_{t=1}^T (\mathbf{z}_t - \hat{\mathbf{z}}_t(\boldsymbol{\theta}^{\text{process}}))^2 + \lambda g(\hat{\mathbf{z}}_t(\boldsymbol{\theta}^{\text{process}}))^2$$

- We penalize values that departs from the PDE
- Good computational news: **this is not an inverse problem, it’s a forward problem**
- Bad computational news: **no closed form, requires gradient descent, backprop, etc.**

# PINNs and penalized functional regression

- Recall **PINN**

$$\operatorname{argmin}_{\boldsymbol{\theta}} \sum_{t=1}^T \left( \mathbf{z}_t - \widehat{\mathbf{z}}_t(\boldsymbol{\theta}^{\text{process}}) \right)^2 + \lambda g \left( \widehat{\mathbf{z}}_t(\boldsymbol{\theta}^{\text{process}}) \right)^2$$

- Even this penalized approach is actually not new in statistics
- Sangalli and collaborators at Politecnico di Milano cast this as a (spatial) functional problem

$$\operatorname{argmin}_{\widehat{\boldsymbol{Y}}} \sum_{i=1}^n \left( \mathbf{z}(s_i) - \widehat{\mathbf{z}}(s_i) \right)^2 + \lambda \int g \left( \widehat{\mathbf{z}}(s) \right)^2 ds$$

- There is no NN, but the idea is very similar and in a continuous setting

# Work 1

- We model the process

$$Z_t(\mathbf{s}) = \mathcal{H} \left( Y_t(\mathbf{s}), \boldsymbol{\theta}^{\text{data}}, \varepsilon_t(\mathbf{s}) \right)$$
$$Y_t(\mathbf{s}) = \mathcal{M}(\{Y_{t-1}(\mathbf{s}), \dots, Y_{t-k}(\mathbf{s})\}, \boldsymbol{\theta}^{\text{process}}, \eta_t(\mathbf{s}))$$

- We model  $\mathcal{M}$  in two different ways
- 1) PS:  $\mathcal{M}$  is derived from the PDE
- 2)  $\mathcal{M}$  is some very large NN and is “nudged” towards the PDE solution

**Contribution 1 (Bonas et al., JASA): develop a flexible temporal model NN for  $\mathcal{M}$  (Gaussian  $\mathcal{H}$ ) and inform it with a PDE in an experimental fluid dynamic problem (Navier-Stokes)**

# Echo State Networks

- Temporal model, call  $\mathbf{Z}_t = \{Z_t(\mathbf{s})\}$

$$\mathbf{Z}_t = \mathbf{B}\mathbf{Y}_t + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

$$\mathbf{Y}_t = (1 - \alpha)\mathbf{Y}_{t-1} + \alpha\boldsymbol{\omega}_t$$

$$\boldsymbol{\omega}_t = g(\mathbf{W}\mathbf{Y}_{t-1} + \mathbf{W}^x \mathbf{x}_t)$$

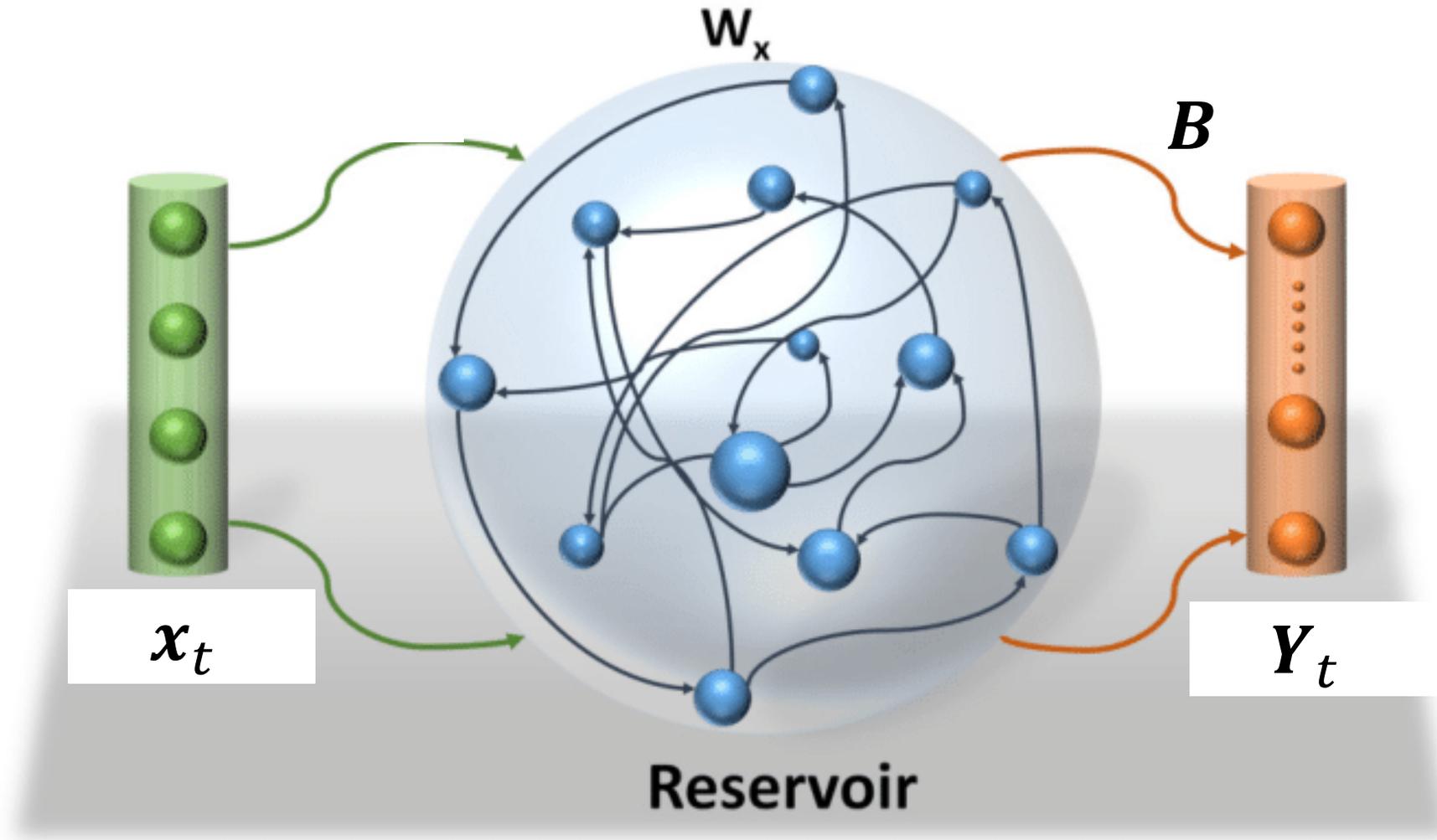
- **Sparse random matrices with spike and slab prior (reservoir)**

$$W_{ij} = p_{ij}^W f(\eta_W) + (1 - p_{ij}^W) \delta_0, \quad p_{ij}^W \sim Be(\pi_W)$$

$$W_{ij}^x = p_{ij}^{W^x} f(\eta_{W^x}) + (1 - p_{ij}^{W^x}) \delta_0, \quad p_{ij}^{W^x} \sim Be(\pi_{W^x})$$

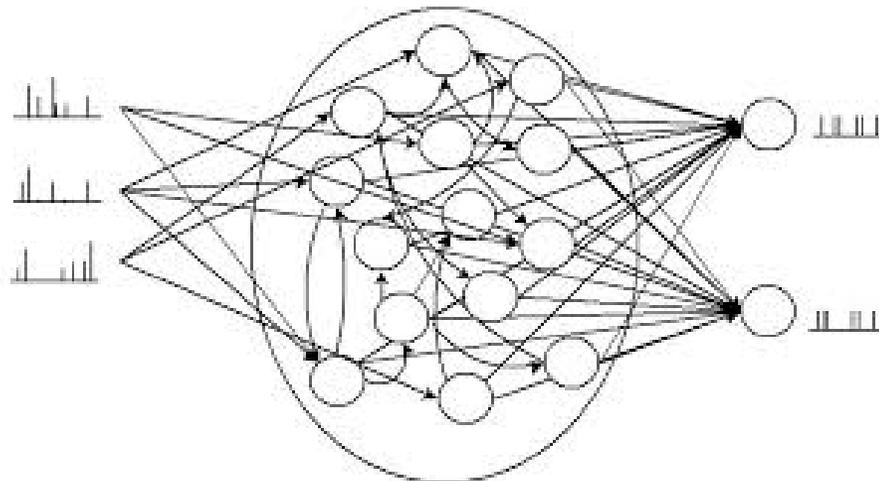
- This is called an **Echo State Network (ESN)**
- This is the shallow version, but **I will use a deep network (DESN)**

# ESN (simplified) in a picture



# Liquid State Machines

- We also use Deep Liquid State Machines (DLSTM) for time series
- This NN **mimics more closely how the brain processes information**
- The input (past data) are transformed in a **subprocess with a firing rate proportional to the data** (*spike train*)
- **Spike trains is used as input** and multiplied by synapse weights to increase (decrease) the stored energy of a neuron

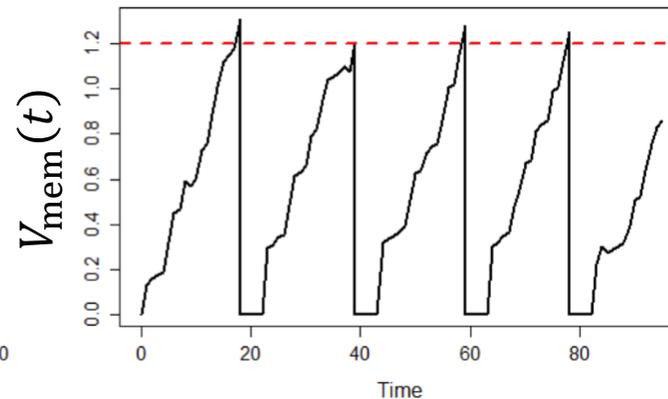
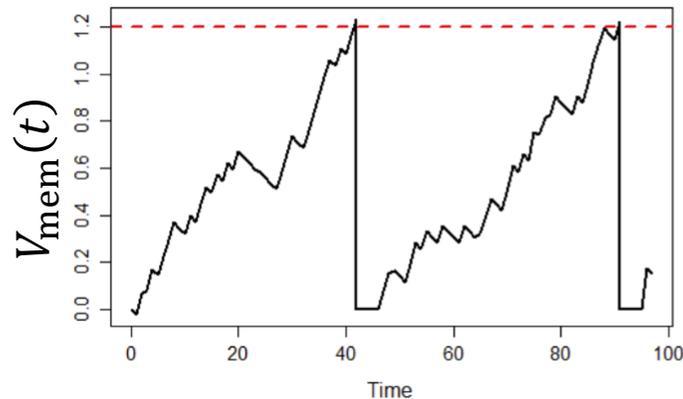


# Liquid State Machines

- Neurons 'spike' once their energy crosses a threshold
- The signal is then sent deeper into the network (like a brain synapse)
- Evolution of each neuron's  $i$  potential energy at time  $t$

$$V_{\text{mem}}(t) = V_{\text{mem}}(t - 1) + \sum_{j \in \text{synapse}} w_{j;i} I_t^{(j)} - V_{\text{leak}}$$

$$V_{\text{mem}}(t) = V_{\text{res}} \quad \text{if } V_{\text{mem}}(t) \geq V_{\text{thr}}$$
$$s_t = \begin{cases} 1, & \text{if } V_{\text{mem}}(t) \geq V_{\text{thr}} \\ 0, & \text{if } V_{\text{mem}}(t) < V_{\text{thr}} \end{cases}$$



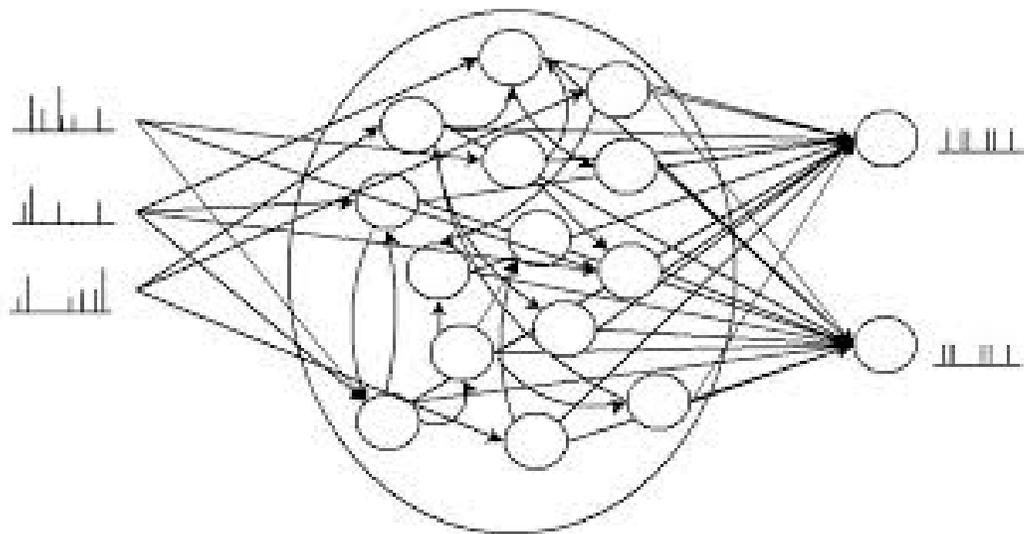
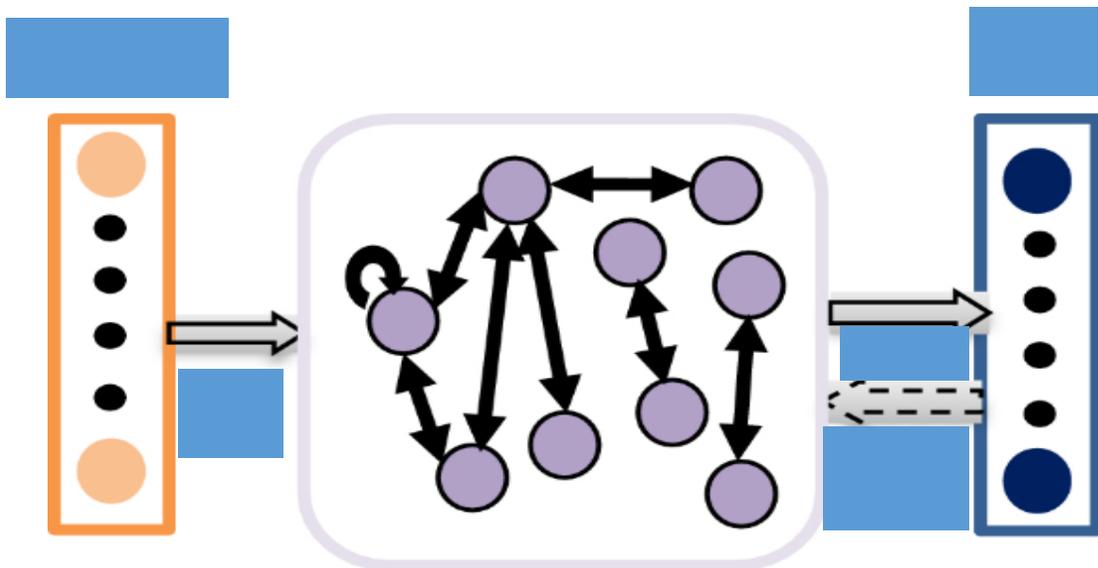
# Liquid State Machines

- We finally count how many spikes

$$\omega_t = \sum_{t^*=1}^{T^*} S_{t^*}$$

- And feed this into a reservoir
- It's actually more complicated than this
- There are two more processes to mimic synapses
- 1) **latent period** after firing
- 2) **lateral inhibition** (once a synapse fires, nearby synapses can't)

DESN



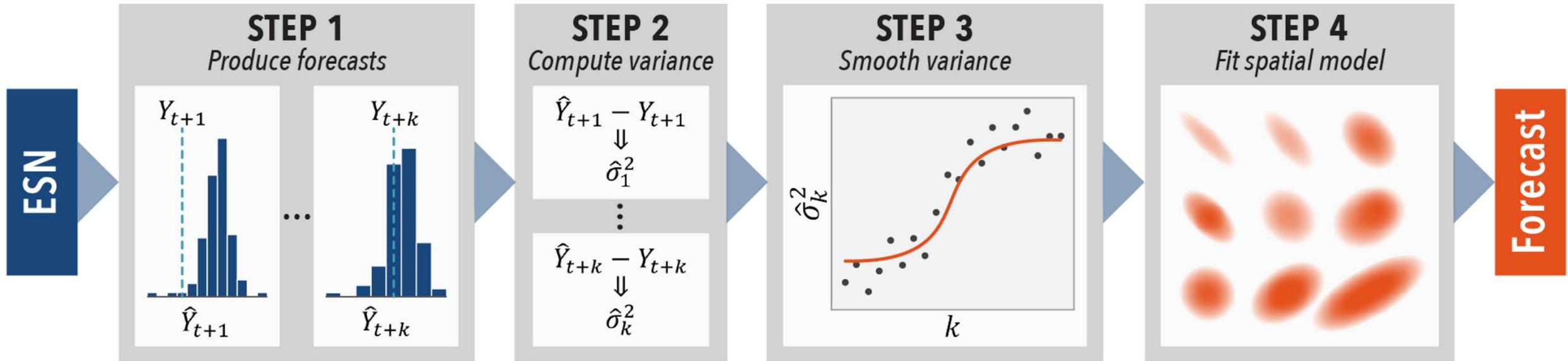
DLISM



ENSEMBLE  
RESERVOIR

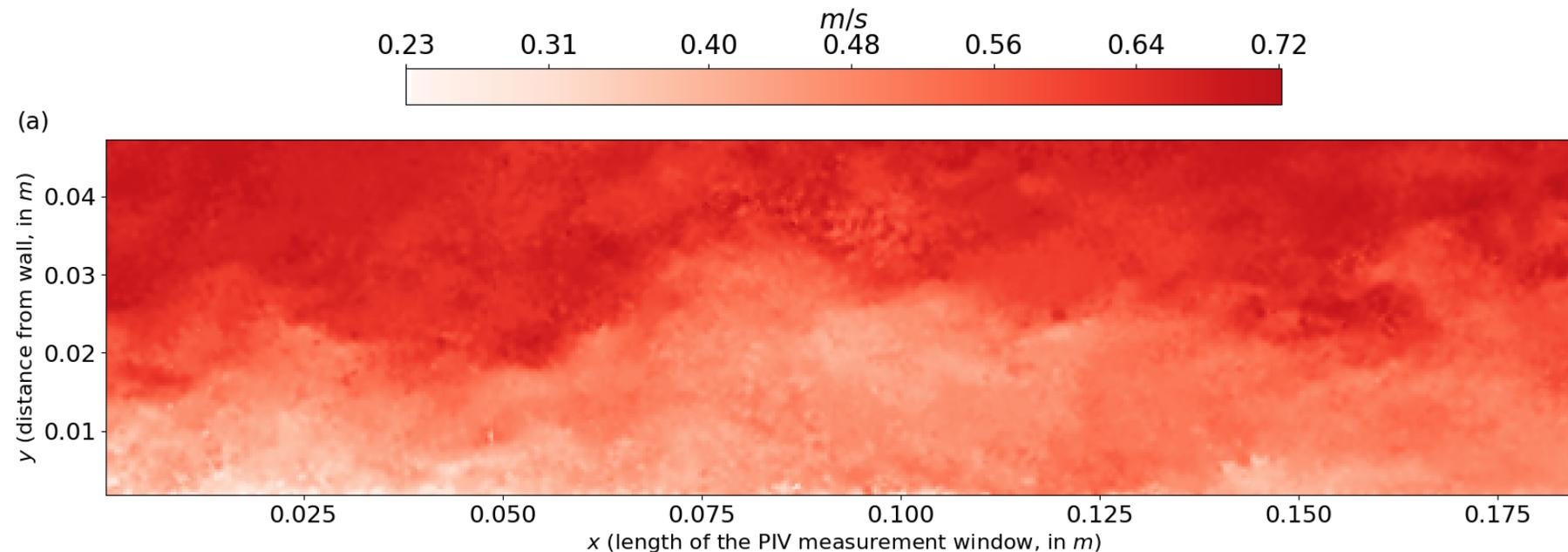


# What about uncertainty?

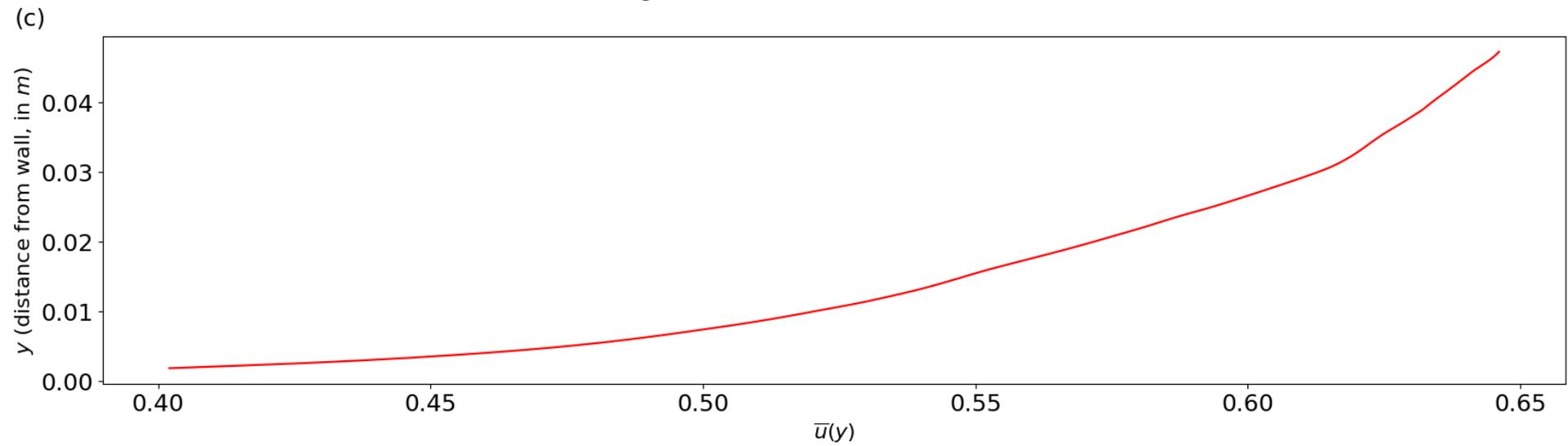
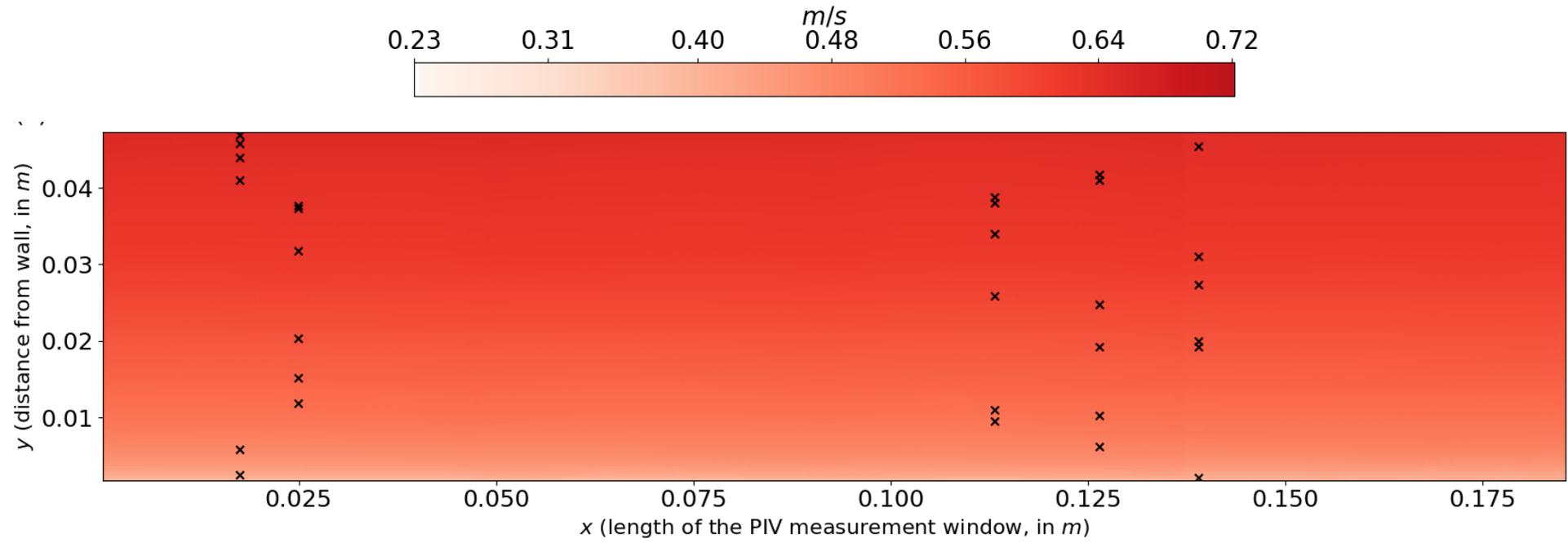


# Data

- **Velocity data of water flow in a tunnel**
- Bottom: wall, top: free flow
- 2,500 time steps at frequency of 1,000 Hz: total time 2.5s
- **Goal: provide physically consistent forecasting**



# Average velocity profile



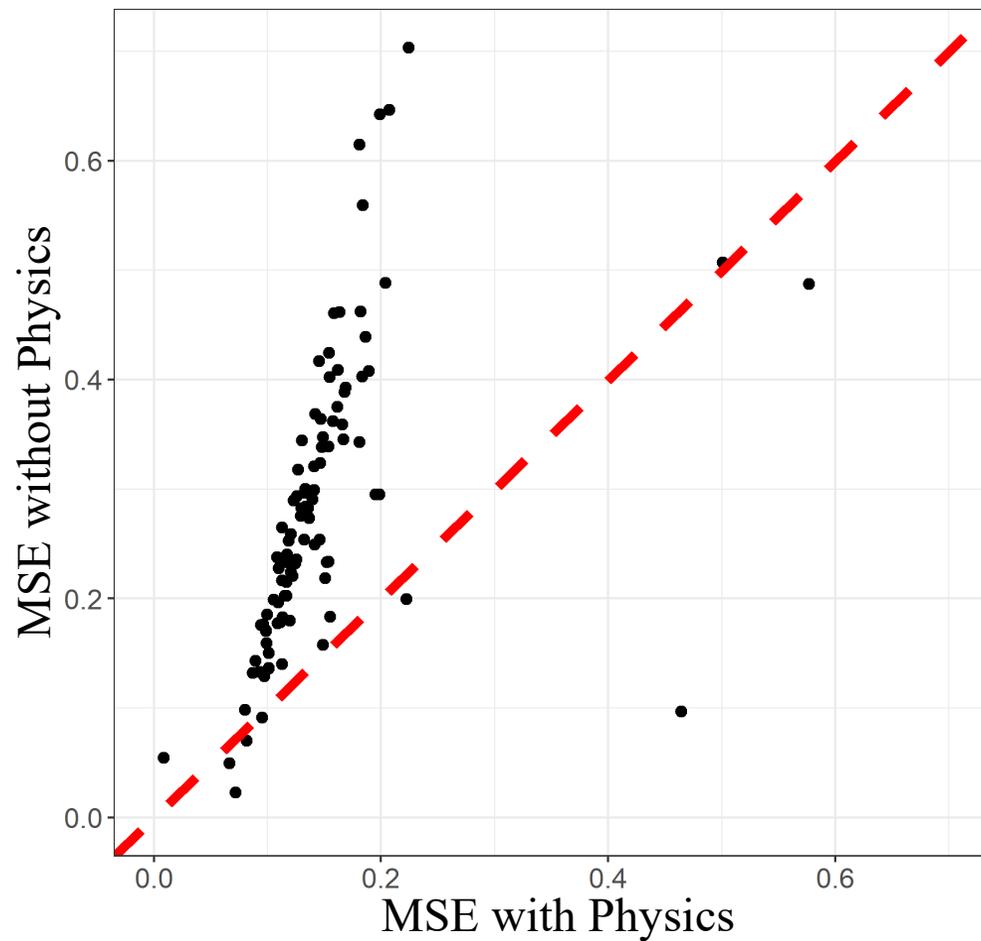
## PDE derivation

- Controlled experiment: 2D fluid at equilibrium
- We use time averaged (RANS) solution to Navier Stokes, which reduces to

$$\frac{u_\tau}{\delta} + \frac{u_\tau \delta}{\text{RE}_\tau} \frac{\partial^2 \bar{Y}}{\partial y^2} + \left\{ (\kappa y)^2 \left| \frac{\partial \bar{Y}}{\partial y} \right| \right\} \frac{\partial \bar{Y}}{\partial y} = 0$$

- $\bar{Y}(x, y)$ : average water velocity
- $\delta$ : boundary layer thickness (known)
- $\text{RE}_\tau$ : Reynolds number (known)
- $u_\tau$ : friction velocity (known)

Method	MSE
With Physics	0.13 (0.05)
No Physics	0.26 (0.17)



# Uncertainty quantification

Prediction Interval	Uncalibrated	Calibrated
95%	100.0 (0.0)	<b>95.0 (1.0)</b>
90%	100.0 (0.0)	<b>90.0 (2.0)</b>
80%	90.0 (3.0)	<b>80.0 (2.0)</b>

# Wrap up work 1

- **A PINN is just a penalized statistical model**
- It works really well in cases where the physics is well known
- No free lunch: you need the right approximation of the Navier Stokes equation
- We also performed uncertainty quantification and calibrated the forecast
- If the PDE penalty has some unknown physical parameters (e.g., **Reynolds number and viscosity**), they can be estimated as well!

## Work 2

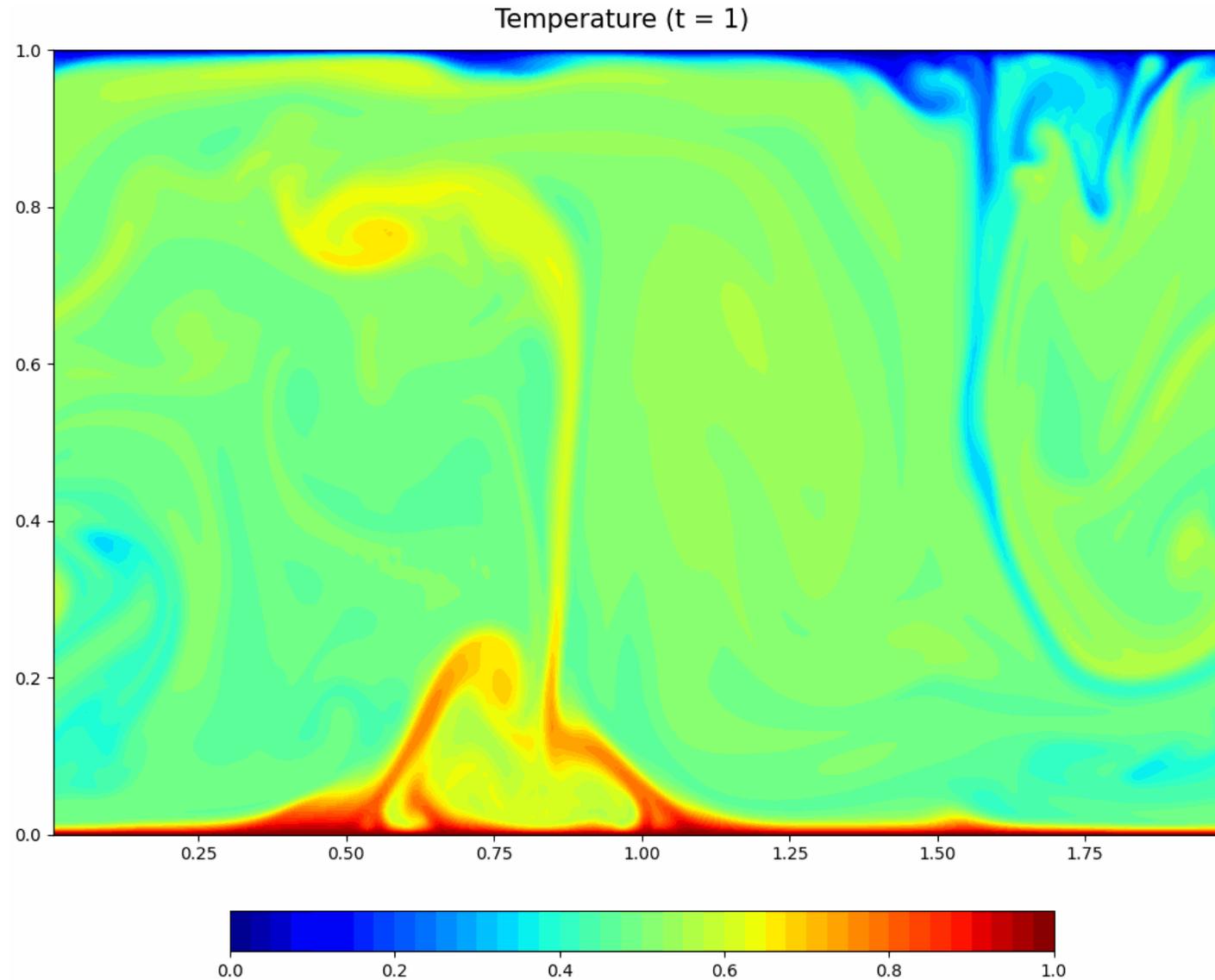
- We model the process

$$Z_t(\mathbf{s}) = \mathcal{H} \left( Y_t(\mathbf{s}), \boldsymbol{\theta}^{\text{data}}, \varepsilon_t(\mathbf{s}) \right)$$
$$Y_t(\mathbf{s}) = \mathcal{M}(\{Y_{t-1}(\mathbf{s}), \dots, Y_{t-k}(\mathbf{s})\}, \boldsymbol{\theta}^{\text{process}}, \eta_t(\mathbf{s}))$$

- $\mathcal{M}$  is some very large NN and is “nudged” towards the PDE solution

**Contribution 2 (Menicali et al., under review): A PINN to generate more realization from a numerical model (digital twin). Model: autoencoder for  $\mathcal{H}$ , an a recurrent NN for  $\mathcal{M}$ , informed with Navier-Stokes**

# The simulation: Rayleigh–Bénard convection (RCB)

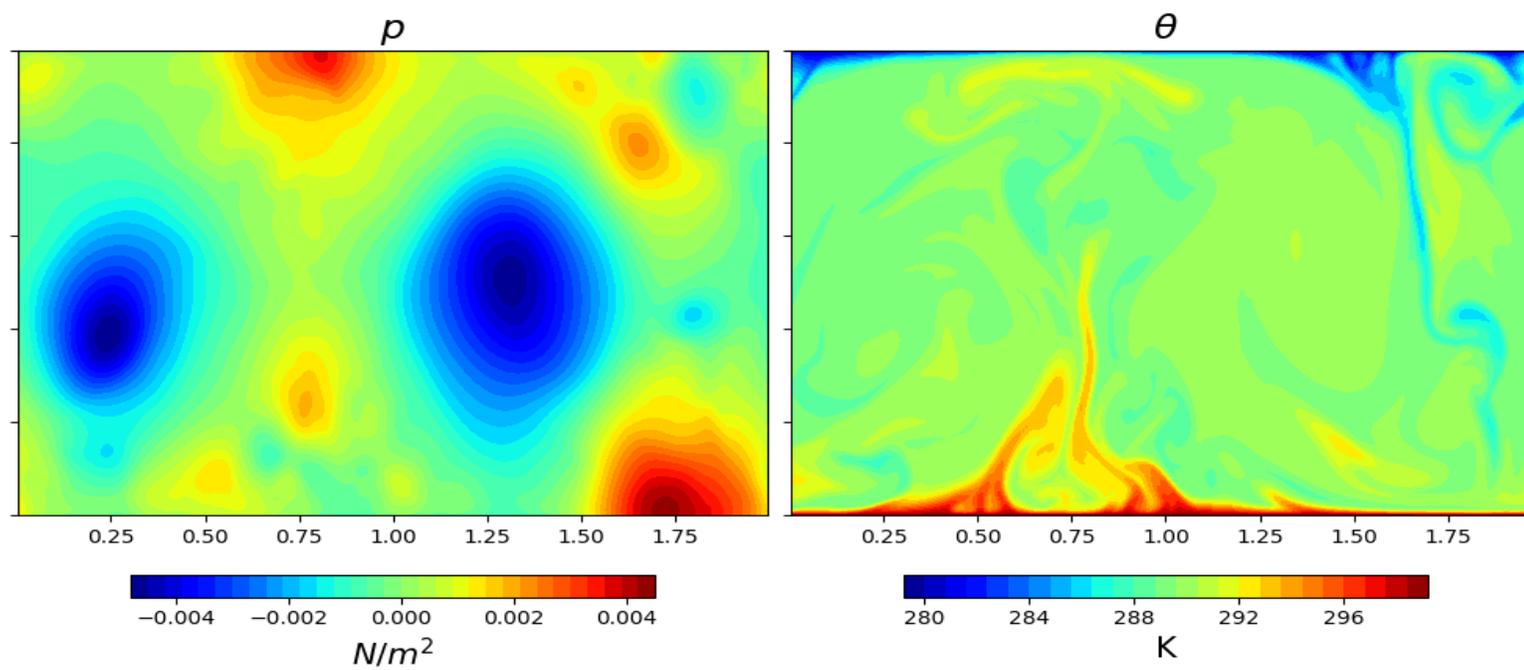
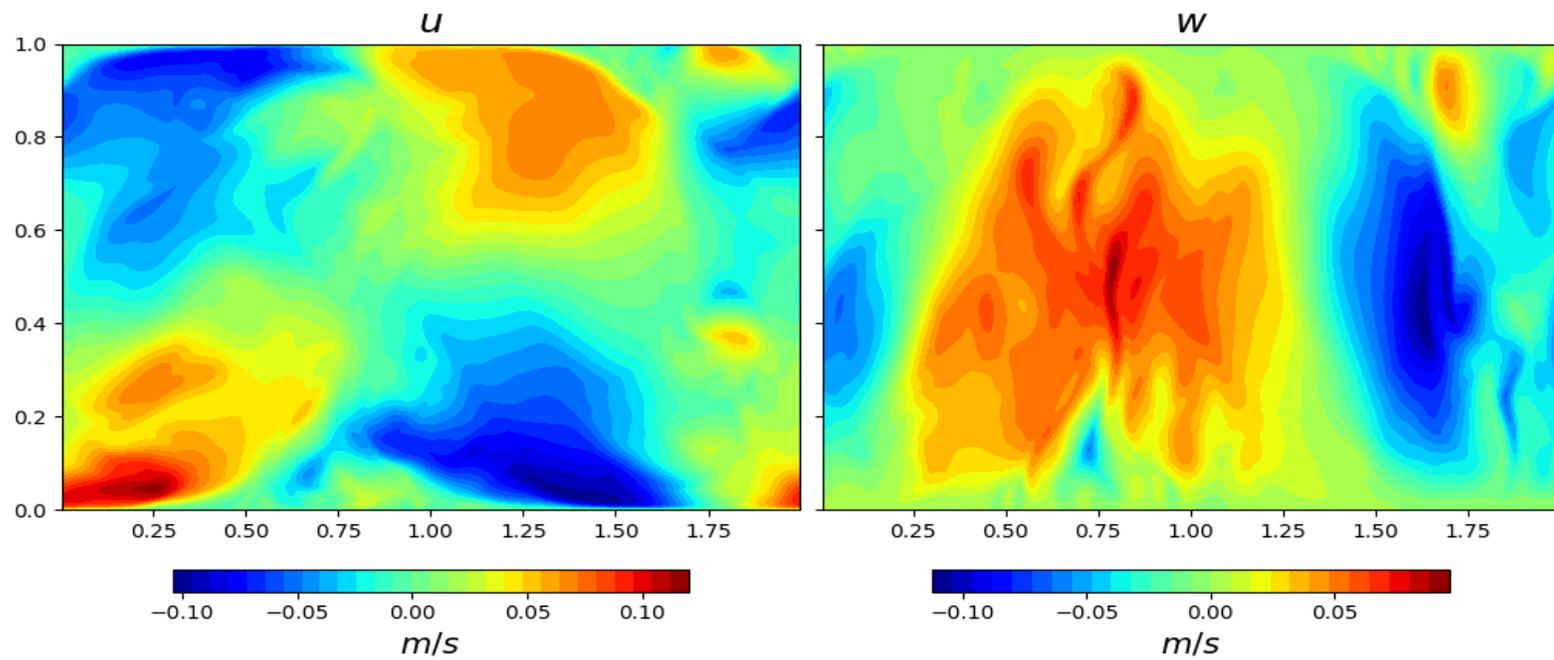


# Data and problem

- Two dimensional fluid
- Variables: velocity  $\mathbf{u}_t = (u_t, w_t)$ , temperature  $\theta_t$  and pressure  $p_t$
- Navier Stokes becomes:

$$\begin{aligned}\nabla \cdot \mathbf{u}_t &= 0, \\ \frac{\partial \mathbf{u}_t}{\partial t} + (\mathbf{u}_t \cdot \nabla) \mathbf{u}_t &= -\nabla p_t + \sqrt{\frac{\text{Pr}}{\text{Ra}}} \nabla^2 \mathbf{u}_t + \theta_t \mathbf{e}_z, \\ \frac{\partial \theta_t}{\partial t} + (\mathbf{u}_t \cdot \nabla) \theta_t &= \sqrt{\frac{1}{\text{PrRa}}} \nabla^2 \theta_t,\end{aligned}$$

- $\mathbf{e}_z$ : (0,1) vector
- Pr and Ra: Prandtl and Rayleigh constants
- The data are simulated on a  $256 \times 256$  grid and 1,000 time points



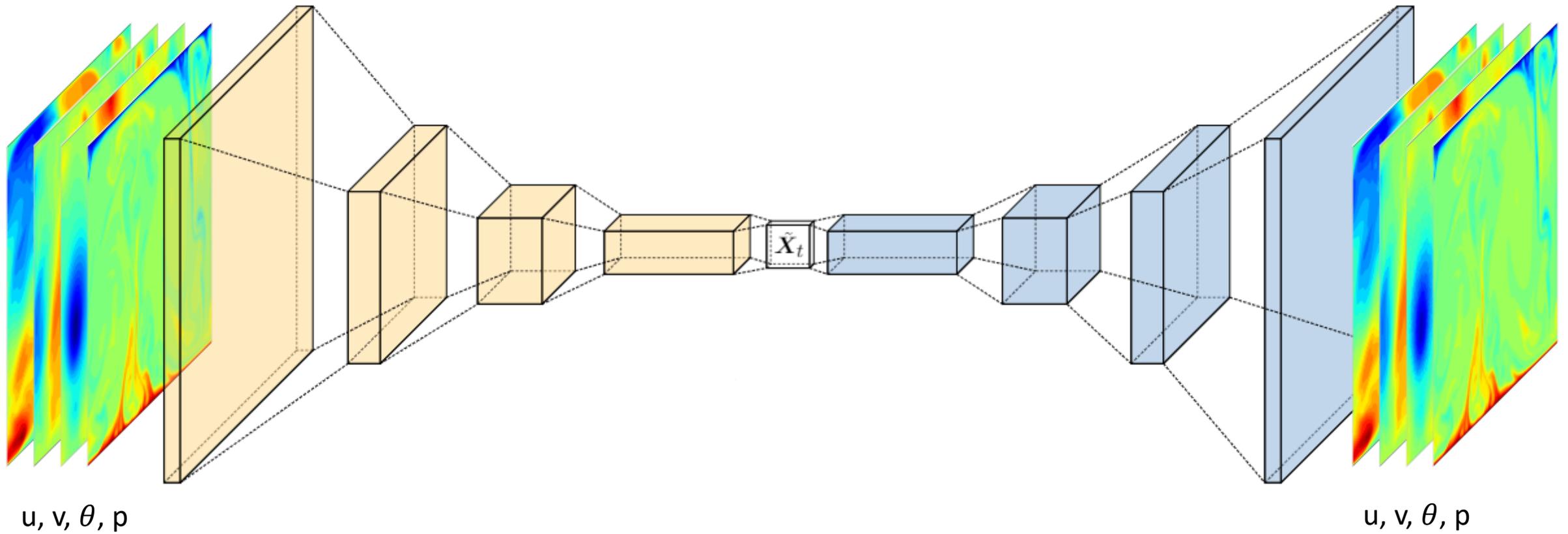
# Convolutional Autoencoder

- We model the spatio-temporal process

$$\mathbf{Y}_t(\mathbf{s}) = (u_t(\mathbf{s}), w_t(\mathbf{s}), p_t(\mathbf{s}), \theta_t(\mathbf{s})),$$
$$\mathbf{Y}_t = \{\mathbf{Y}_t(\mathbf{s})\}$$

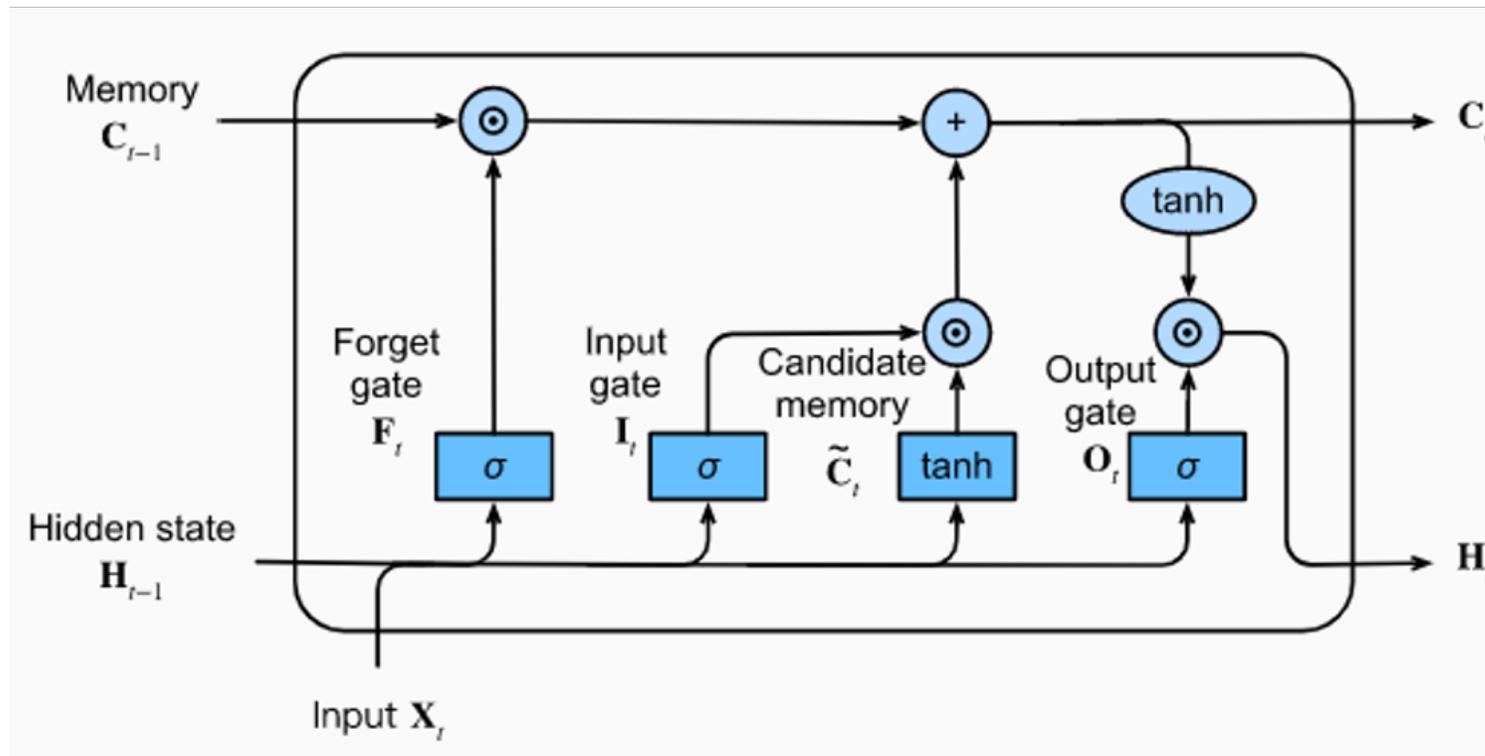
- We model the spatial structure of  $\mathbf{Y}_t$  with a Convolutional Autoencoder (CAE)
- The input is projected onto a latent space: the **spatial encoder**
- Then projected back onto the original space: **spatial decoder**
- It's a **dimension reduction technique** in space, like fixed rank Kriging, predictive processes, etc.

# A CAE in a picture

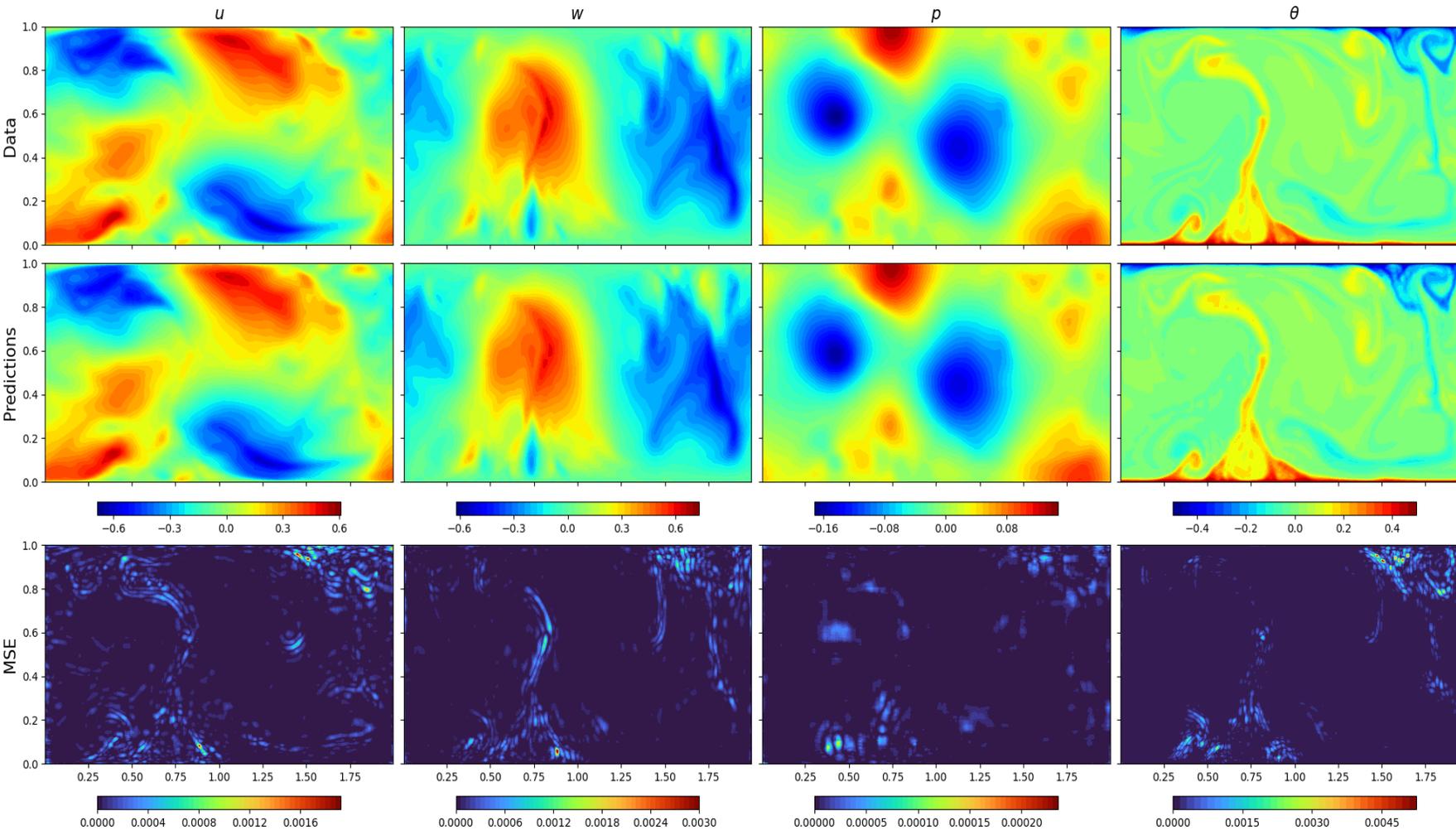


# ConvLSTM

- We assume an LSTM model for the temporal structure
- This is called a Convolutional Long-Short Term Memory (ConvLSTM) model



# Results: Spatial part



MSE

Fixed Kriging:  $5.0 \times 10^5$

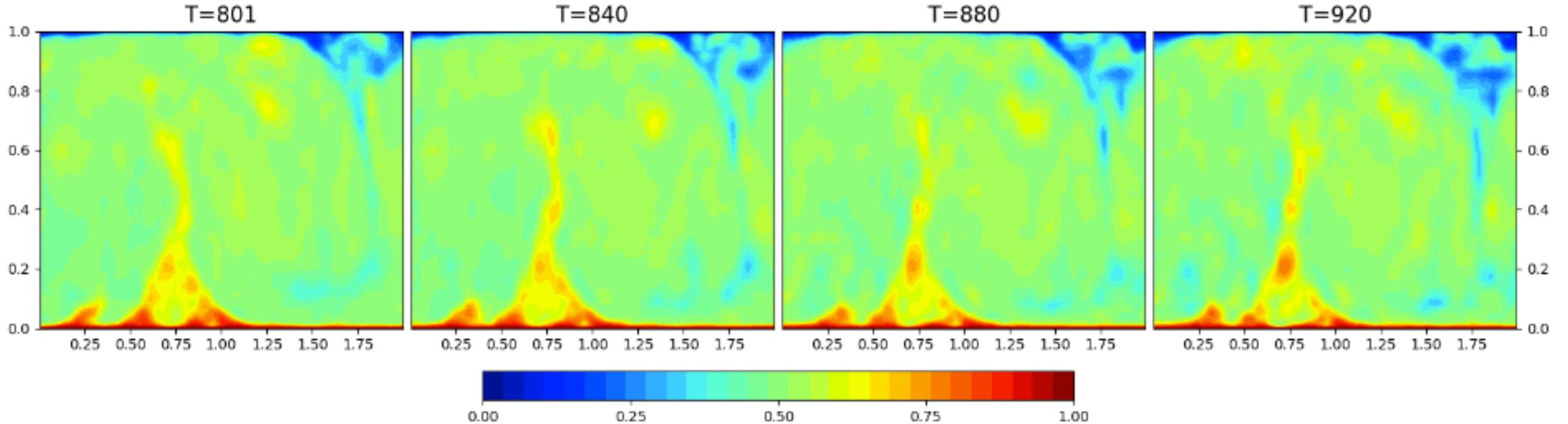
PCA:  $1.5 \times 10^{-3}$

ICA:  $3.6 \times 10^{-3}$

**CAE:  $3.1 \times 10^{-5}$**

# Results: Forecast

## Temperature Forecast



### MSE

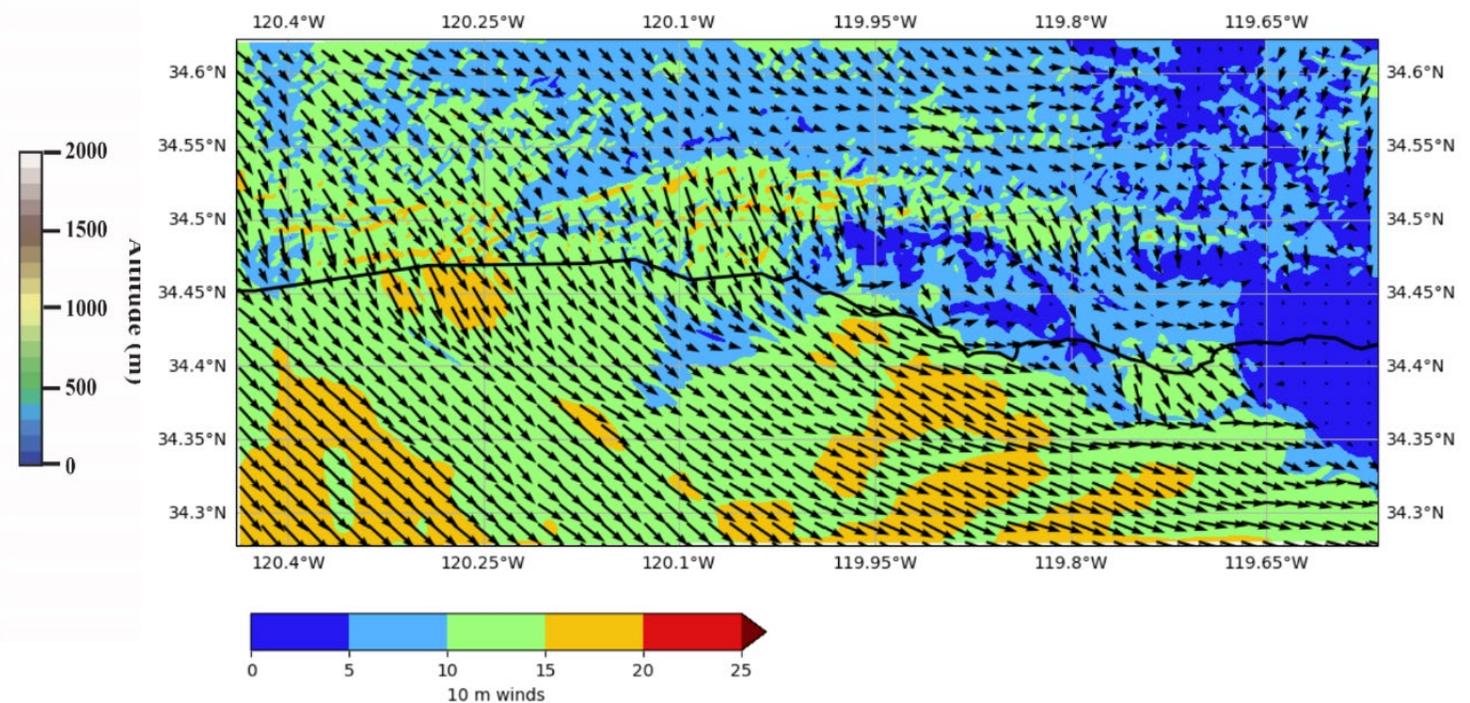
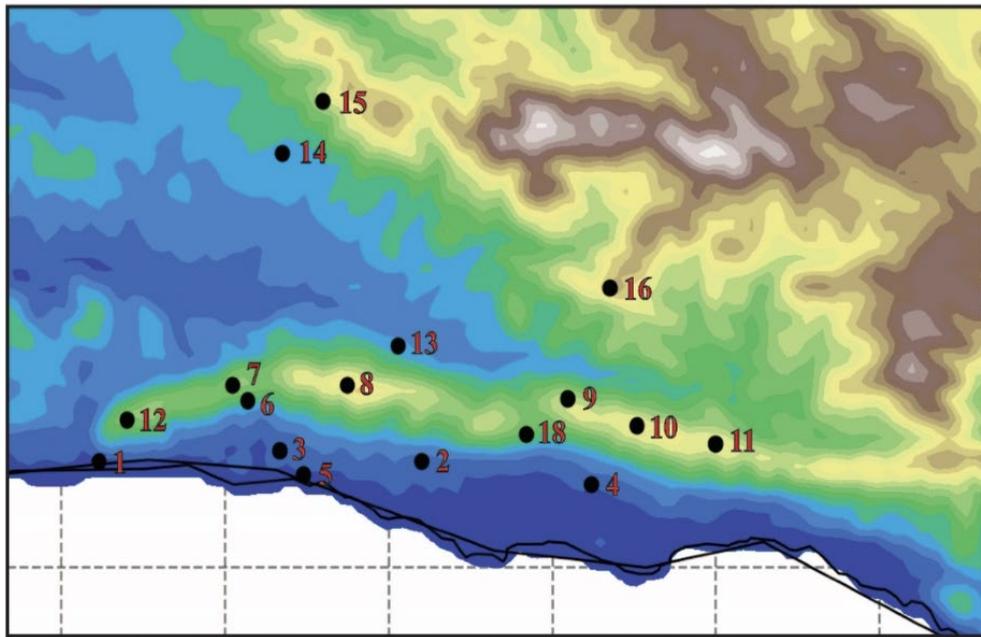
ARMA: 163.25

PI-ESN: 21.48

CRNN: 0.36

**Our Model: 0.22**

# PINNs and WRF-LES: the SWEX experiment



- A PINN for WRF-LES (mesoscale-microscale): no more simplifying Navier Stokes
- Tasks: integrate a statistical model with the WRF Fortran (!) and build a PINN

# Other interesting projects

- **PINNs and physics-informed emulators** for satellite retrievals
- **Physics-informed priors**
- **Land-vegetation dynamics** in the Midwest Holocene
- **Pushing emulations to exascale** for Petabytes of data (2024 Gordon Bell!)
- Statistical models for climate output **compression**
- Visualization in VR and other 3D environments
- High energy particle physics, optics, etc.



Thanks very much for your attention!

