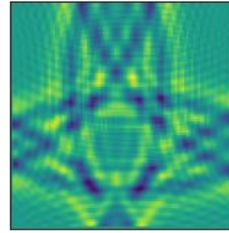
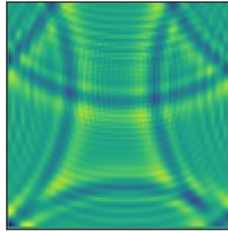
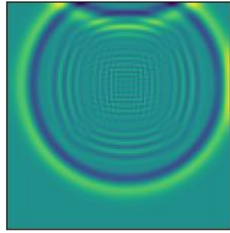
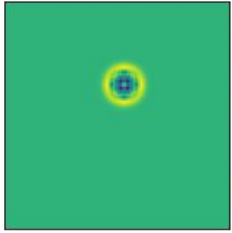


# Learning Earth system model dynamics with implicit schemes



**Marcel Nonnenmacher**

Model-driven Machine Learning  
Institute for Coastal Research  
Helmholtz-Centre Hereon

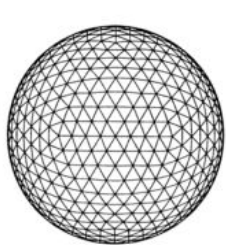


Model-driven  
Machine Learning

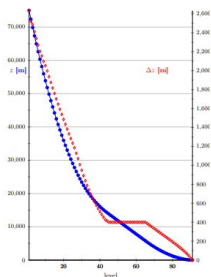
**HELMHOLTZ**AI | ARTIFICIAL INTELLIGENCE  
COOPERATION UNIT

# Earth Science models and their applications

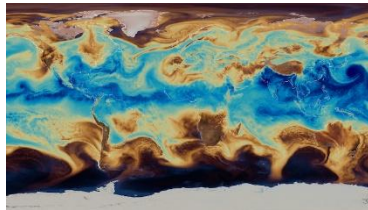
ICON model for numerical weather prediction (DWD, MPI-M)



2949120  
grid cells



90  
z-levels



10+  
variables / location

Total values / time point: **2.6 billion**

Or with a time step of 120 seconds,  
Total values / day: **1.9 trillion**

Images: Stevens et al., 2019, Reinert et al., 2021

Data assimilation application for weather forecasts (ECMWF)

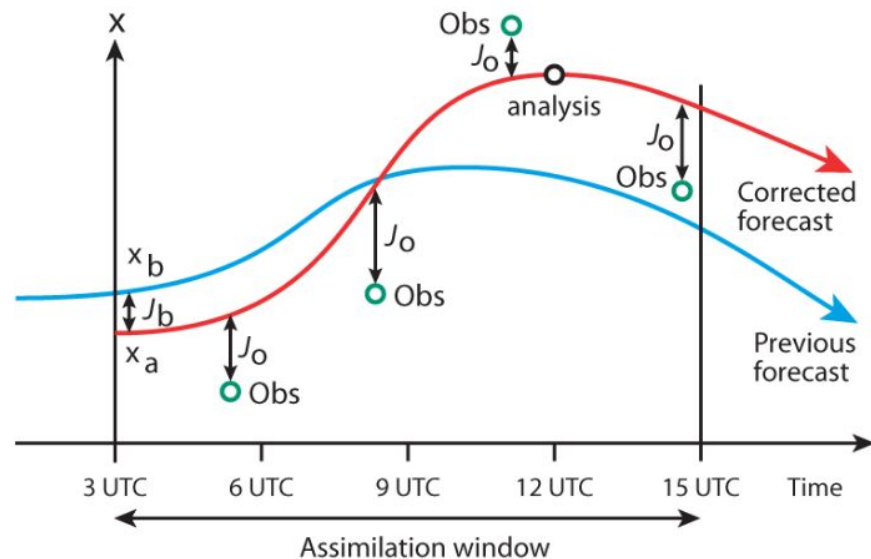


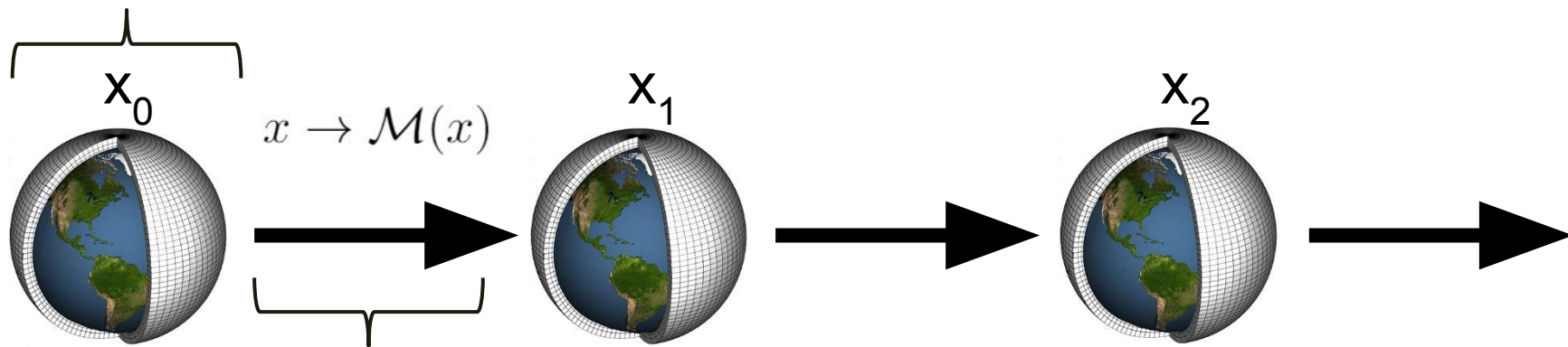
Image source: ecmwf.int

“20 years of 4D-Var: better forecasts through a better use of observations”

# Physics-based Simulation in Earth Science

## System state

- Physical variables (pressure, temperature, ...)
- Snapshot for one time point



## State update function

- Physical processes (radiation, turbulence, phase changes, ...)
- Mass transfer due to gravity, wind, precipitation, ...

Derivatives of state update function  $\frac{dX_1}{dX_0}$

- Data assimilation, model tuning, ...

# Emulation of numerical simulators

Simulation for Earth Science models is **expensive**.



Mistral, DKRZ ( $>10^5$  processors, 3.6 PetaFLOPS)

**Emulation:** ML model imitates the state update function of a numerical simulator.

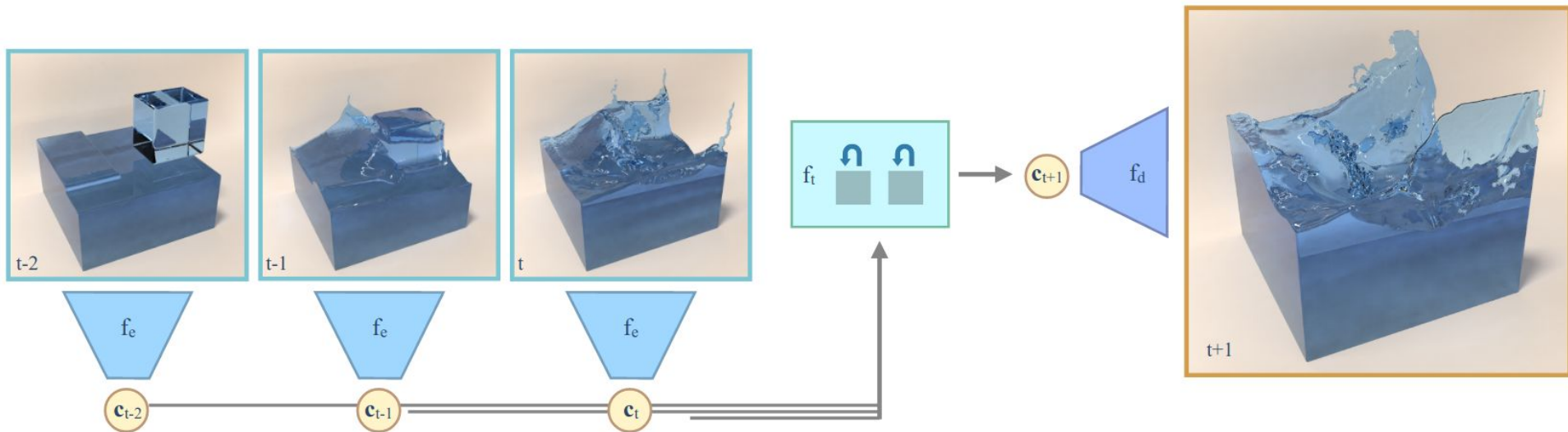
Possible advantages of ML models:

- parallelization on GPU clusters
- may find “shortcuts” from large datasets
- trained in “auto-diff” environments  
(no need to write and maintain derivative routines)

Simulation code: FORTRAN, often  $> 10^5$  lines

```
108 do i = 1,nx,1
109   im = max(1,i-1)
110   ip = min(nx,i+1)
111   |
112   if((d(i).gt.0.1).and.(d(ip).gt.0.1)) then
113     ce(i) = dt*dt*wimp*wimp*g*0.5*(h(i)+h(ip))/(dx*dx)
114   endif
115   |
116   if((d(i).gt.0.1).and.(d(im).gt.0.1)) then
117     cw(i) = dt*dt*wimp*wimp*g*0.5*(h(i)+h(im))/(dx*dx)
118   endif
119
120 -enddo
```

# Emulator examples



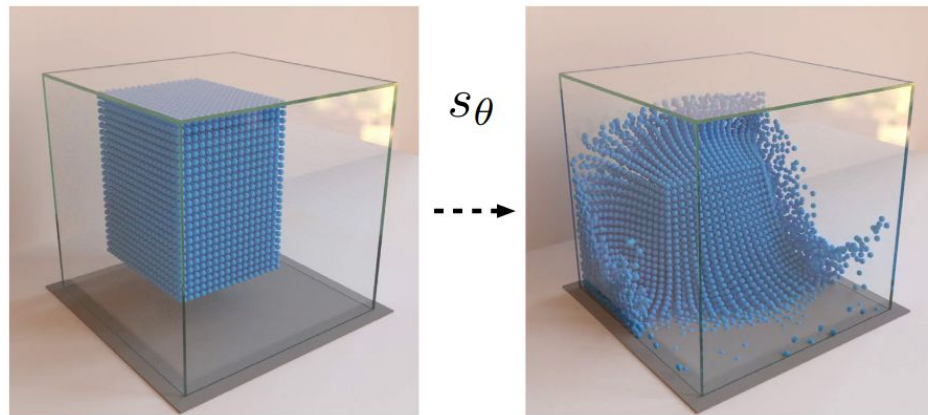
S. Wiewel<sup>1</sup>, M. Becher<sup>1</sup>, N. Thuerey<sup>1 †</sup>

<sup>1</sup>Technical University of Munich

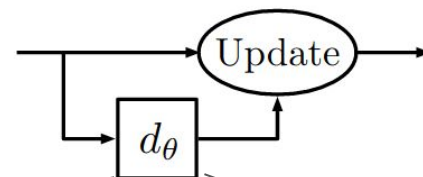
Computer Graphics Forum, Volume 38 (2019) – Issue 2

Presented at Eurographics Conference 2019, May 2019

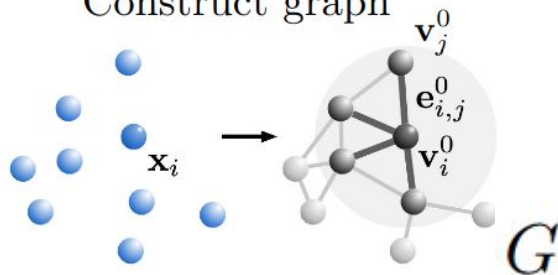
# Emulator examples



Learned simulator,  $s_\theta$



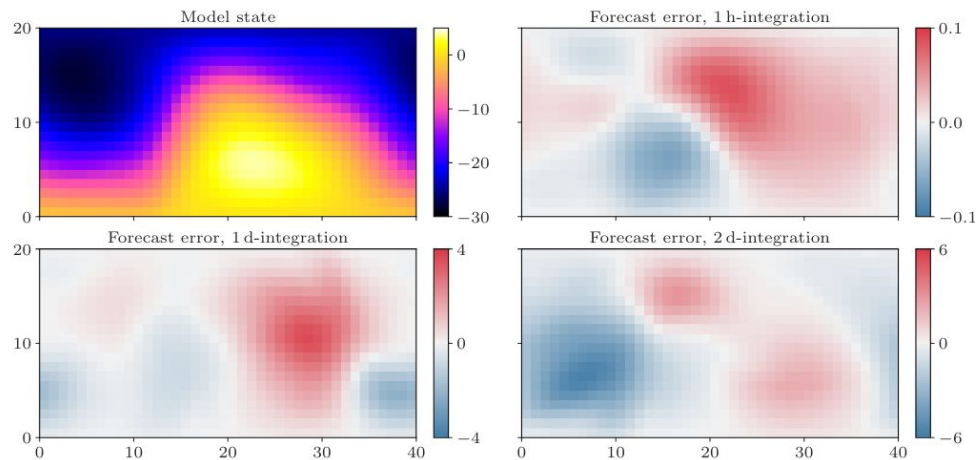
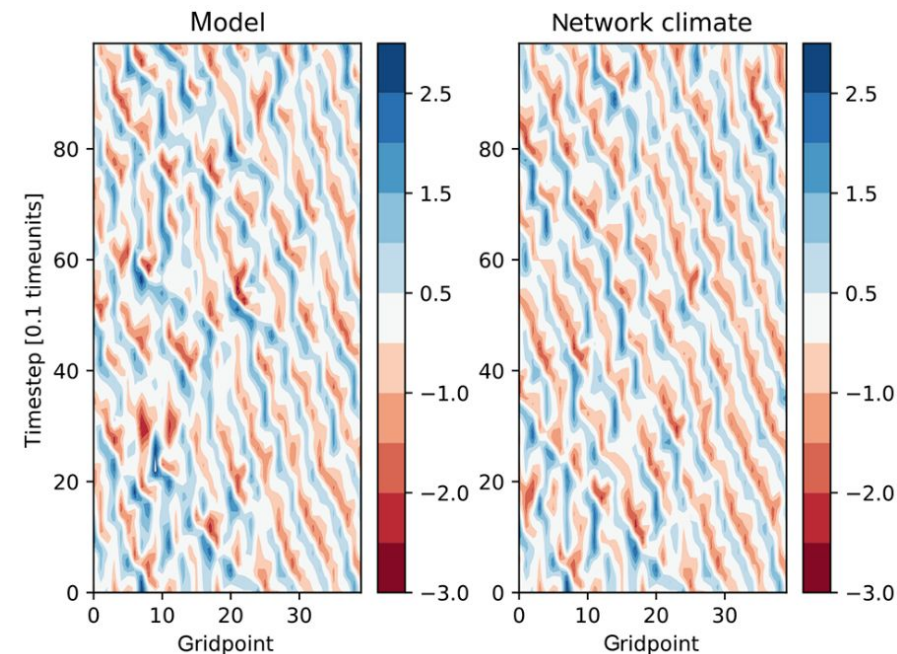
Construct graph



**Alvaro Sanchez-Gonzalez, Jonathan Godwin, Tobias Pfaff, Rex Ying, Jure Leskovec, Peter Battaglia** *Proceedings of the 37th International Conference on Machine Learning*, PMLR 119:8459-8468, 2020.



# Emulator examples



**Alban Farchi**

**Marc Bocquet**

CERA

joint laboratory École des Ponts ParisTech and EDF R&D  
Champs-sur-Marne, France

**Patrick Laloyaux**

**Massimo Bonavita**

ECMWF

Shinfield Park  
Reading, United Kingdom

Nonlin. Processes Geophys., 26, 381–399, 2019

<https://doi.org/10.5194/npg-26-381-2019>

Sebastian Scher<sup>1</sup> and Gabriele Messori<sup>1,2</sup>

<sup>1</sup>Department of Meteorology and Bolin Centre for Climate Research, Stockholm University,

<sup>2</sup>Department of Earth Sciences, Uppsala University, Uppsala, Sweden

arXiv:2010.12605v2 [stat.ML] 10 May 2021

The core principles (L96)

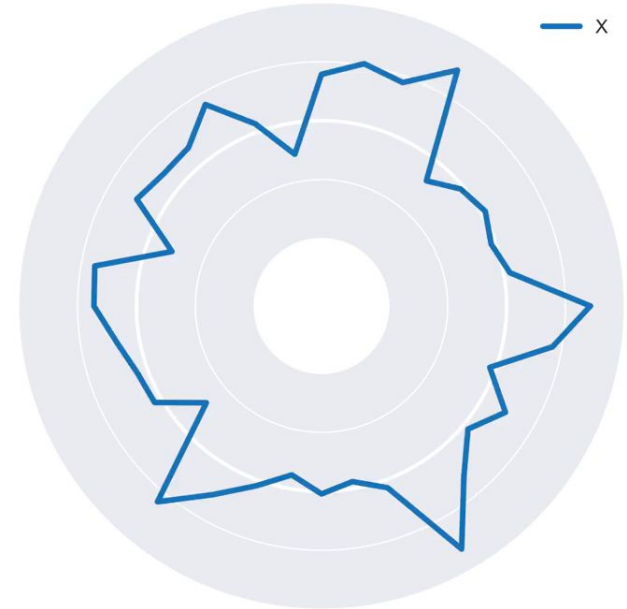
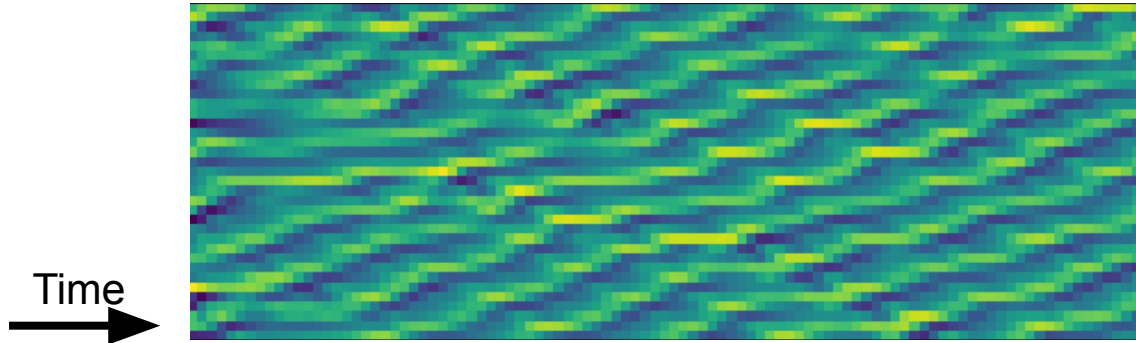


# Model System: Lorenz '96

$$\frac{dx_k}{dt} = -x_{k-1} (x_{k-2} - x_{k+1}) - x_k + F$$

Lorenz, 1996

- 40 coupled nonlinear differential equations
- Chaotic dynamics (Lyapunov time 1.67 for  $F = 8$ )



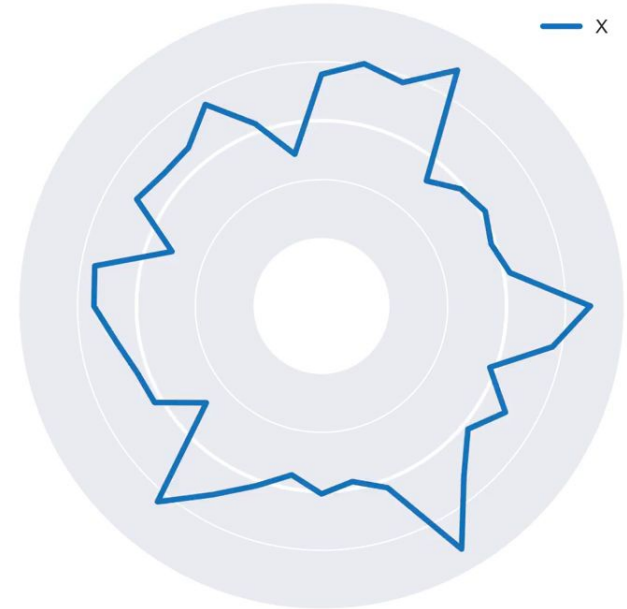
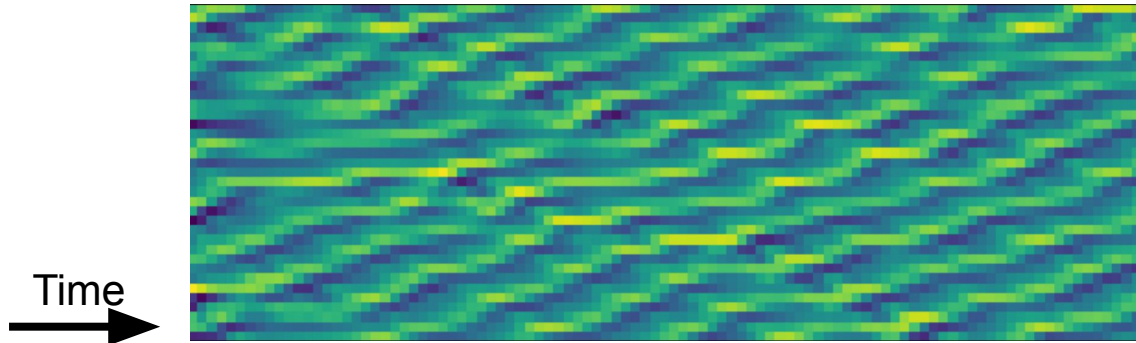
# Model System: Lorenz '96

$$\frac{dx_k}{dt} = -x_{k-1} (x_{k-2} - x_{k+1}) - x_k + F$$

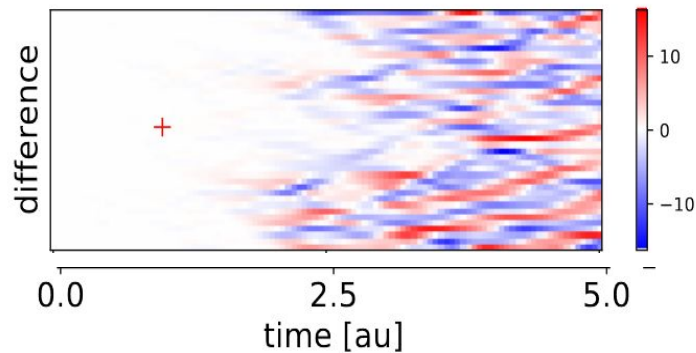
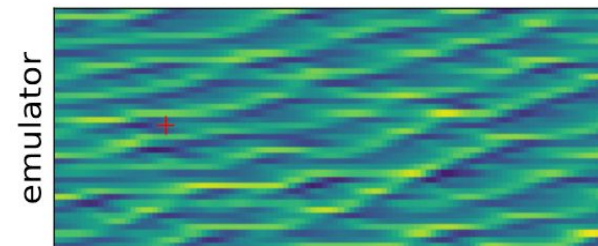
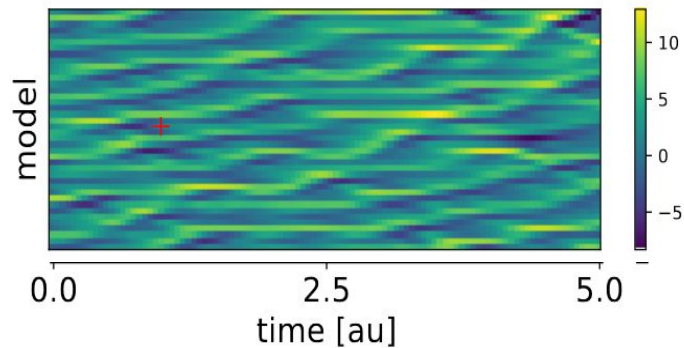
Lorenz, 1996

**local computations !**

- 40 coupled nonlinear differential equations
- Chaotic dynamics (Lyapunov time 1.67 for  $F = 8$ )

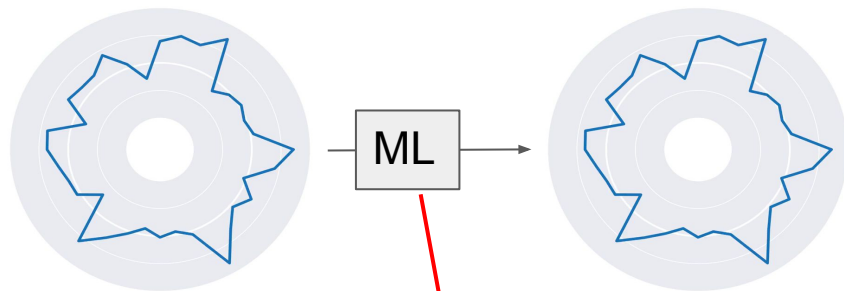


# Basic emulator training



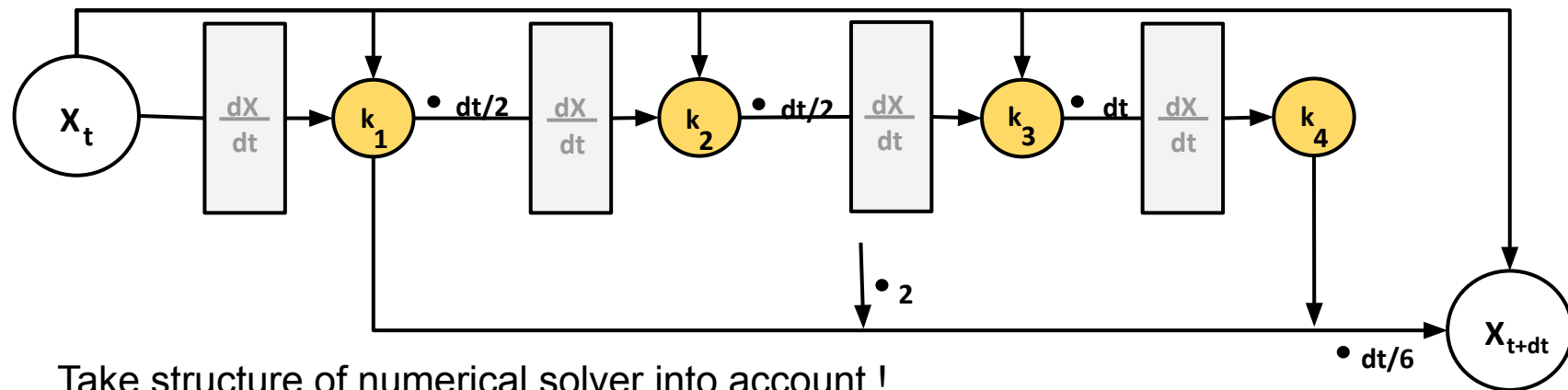
state at time  $t$

state at time  $t+1$



local function!

# Structured ML models for learning the update step



Take structure of numerical solver into account !

$$\text{RK4: } X_{t+dt} = X_t + \frac{dt}{6}(k_1 + 2(k_2 + k_3) + k_4)$$

$$k_1 = f(X_t)$$

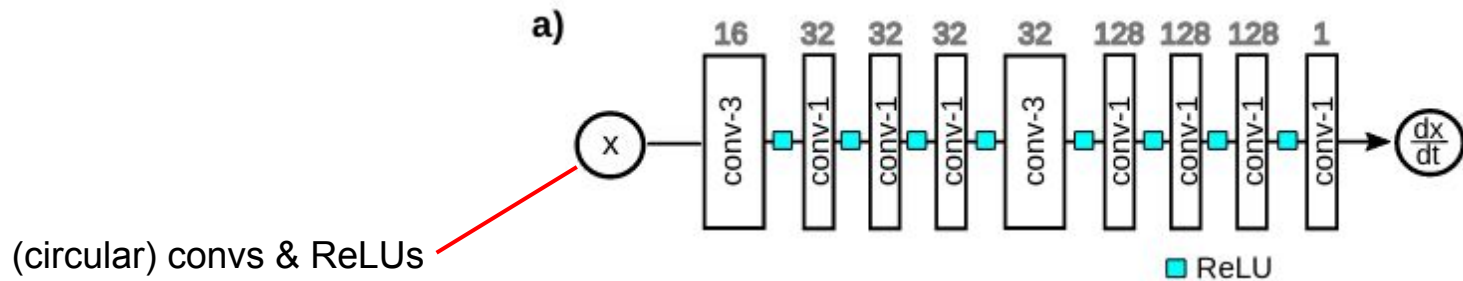
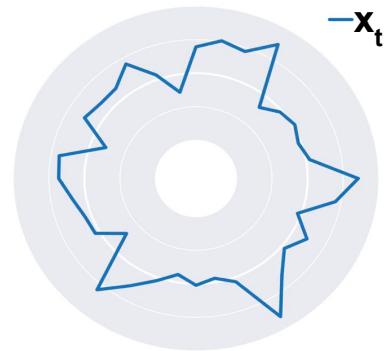
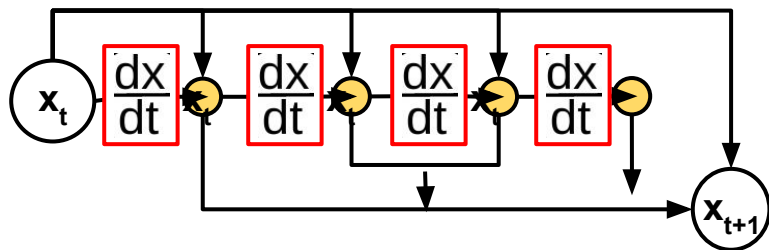
$$k_3 = f\left(X_t + k_2 \cdot \frac{dt}{2}\right)$$

$$k_2 = f\left(X_t + k_1 \cdot \frac{dt}{2}\right)$$

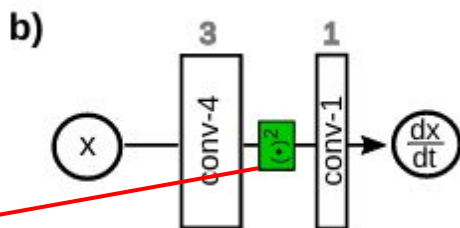
$$k_4 = f(X_t + k_3 \cdot dt)$$

If  $f$  is differentiable (wrt.  $\mathbf{x}_t$ ), so is  $\mathbf{x}_{t+1}$  !

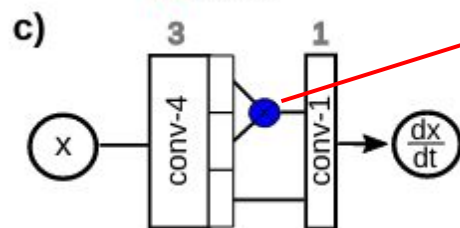
# Network architectures



(circular) convs & ReLUs  
convolution = local function

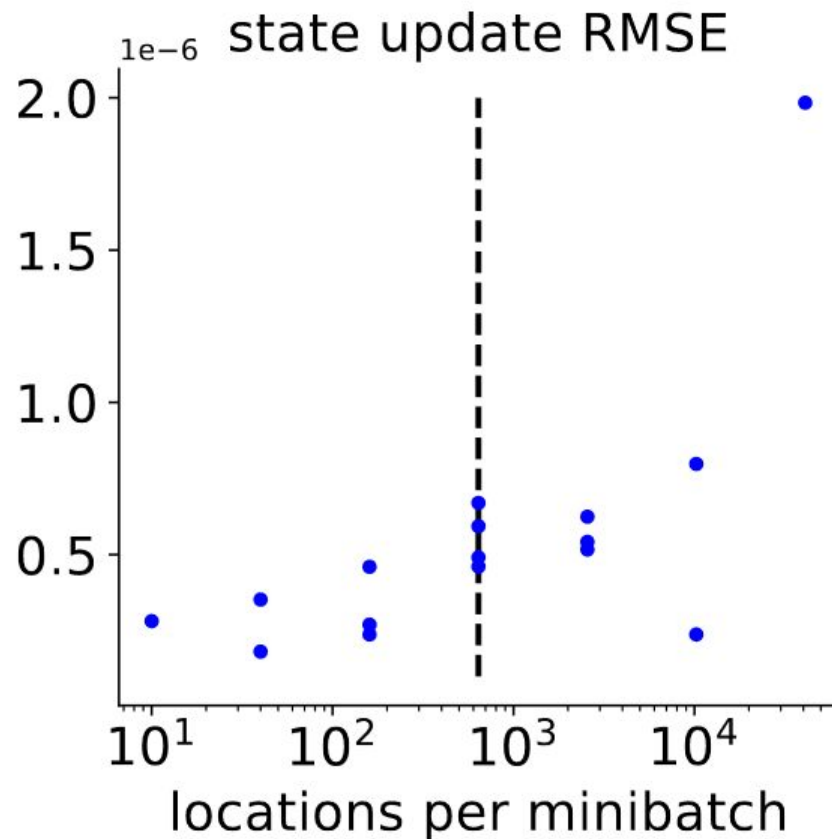
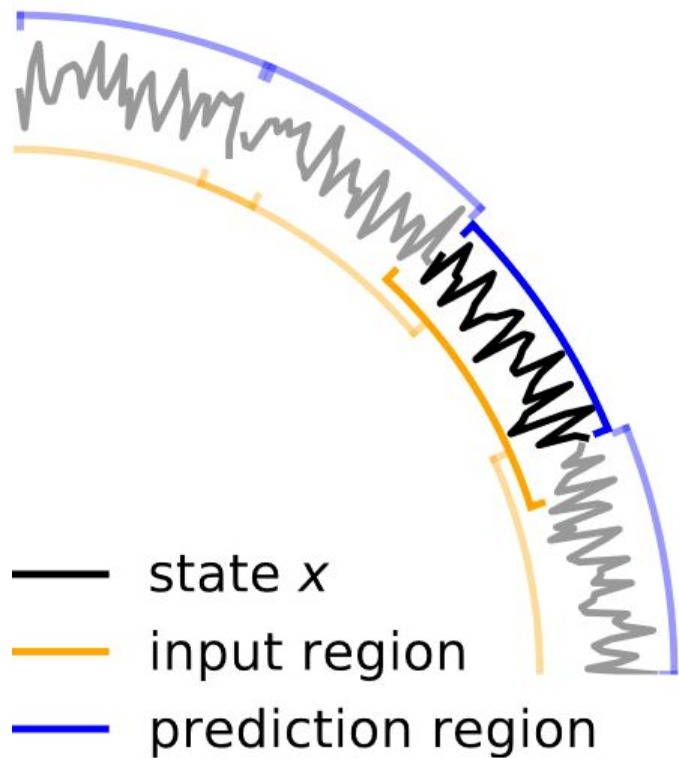


quadratic nonlinearity



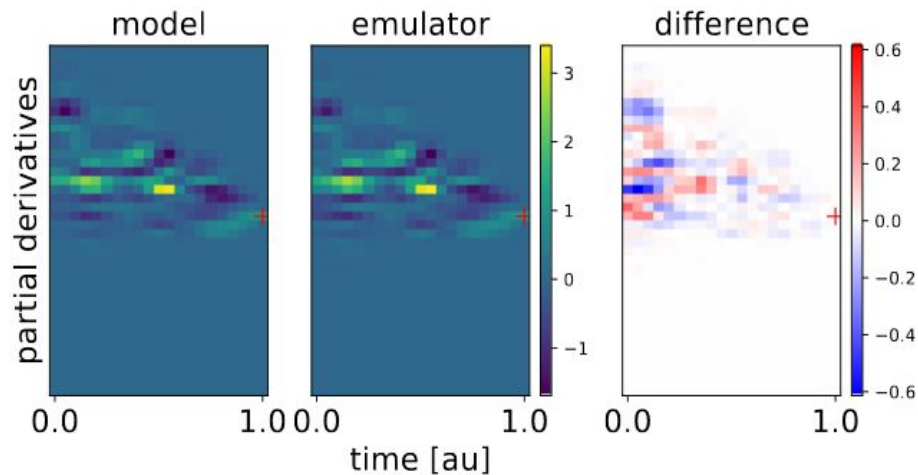
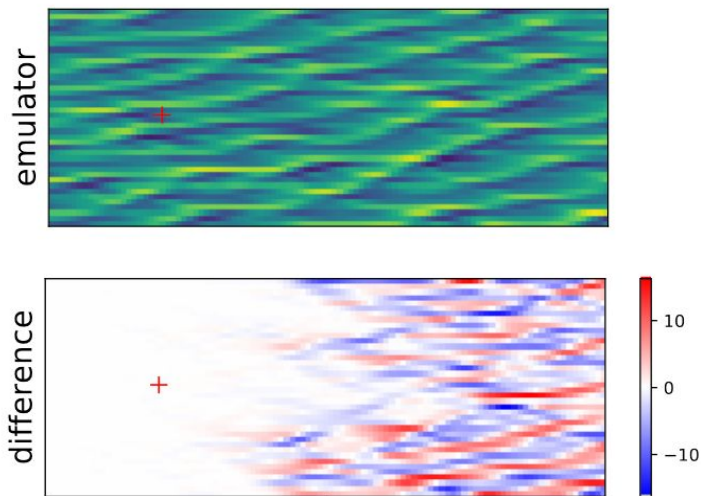
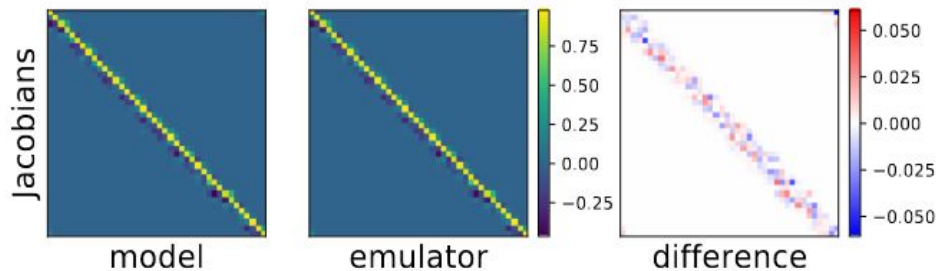
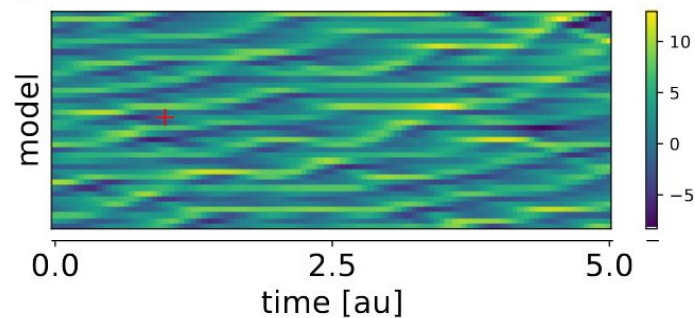
bilinear layer

# Training on partial system states





# Trained emulators reproduce system states and derivatives



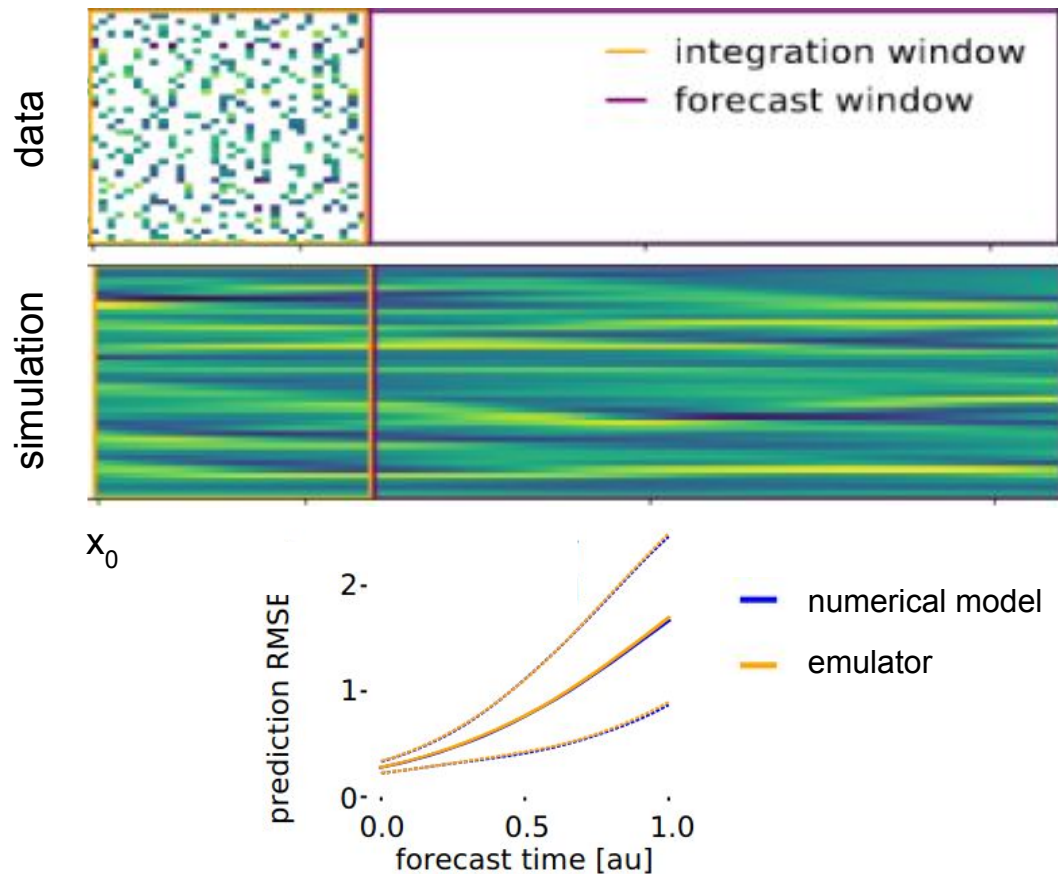
Data assimilation, parametrization tuning etc.

# Emulators for 4D-Var data assimilation

1) with **differentiable** dynamical model, **estimate initial state  $x_0$**  from noisy & incomplete data.

2) starting from  $x_0$ , simulate future states beyond final data point.

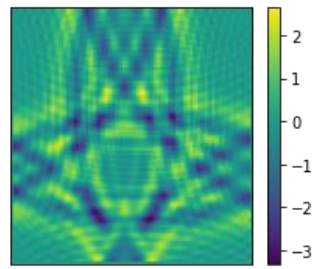
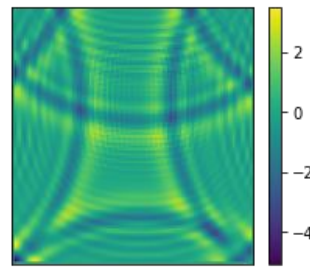
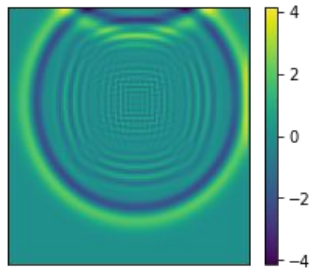
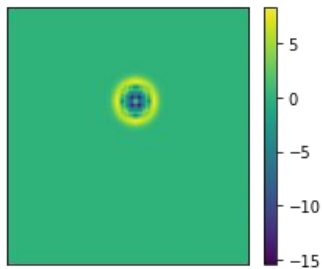
(optional) evaluate prediction error



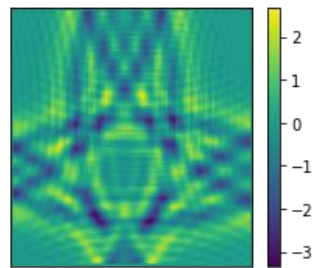
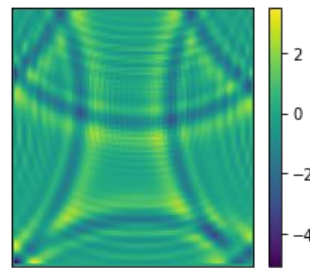
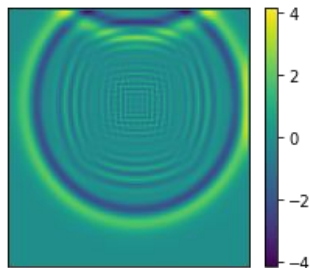
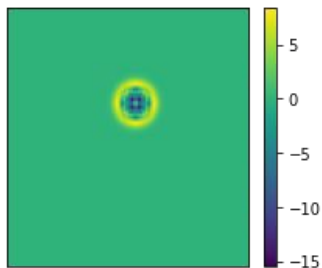
2D systems and beyond explicit solvers

# Shallow Water Equations

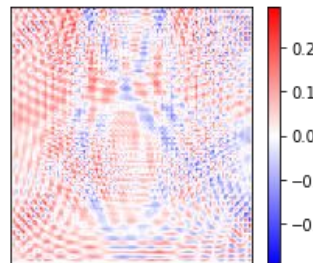
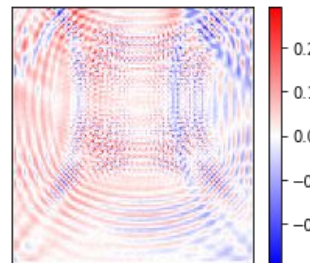
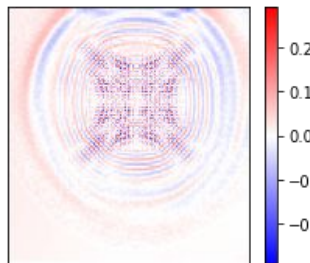
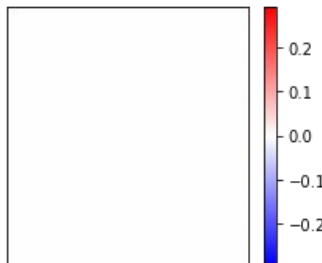
Shallow  
water  
equations



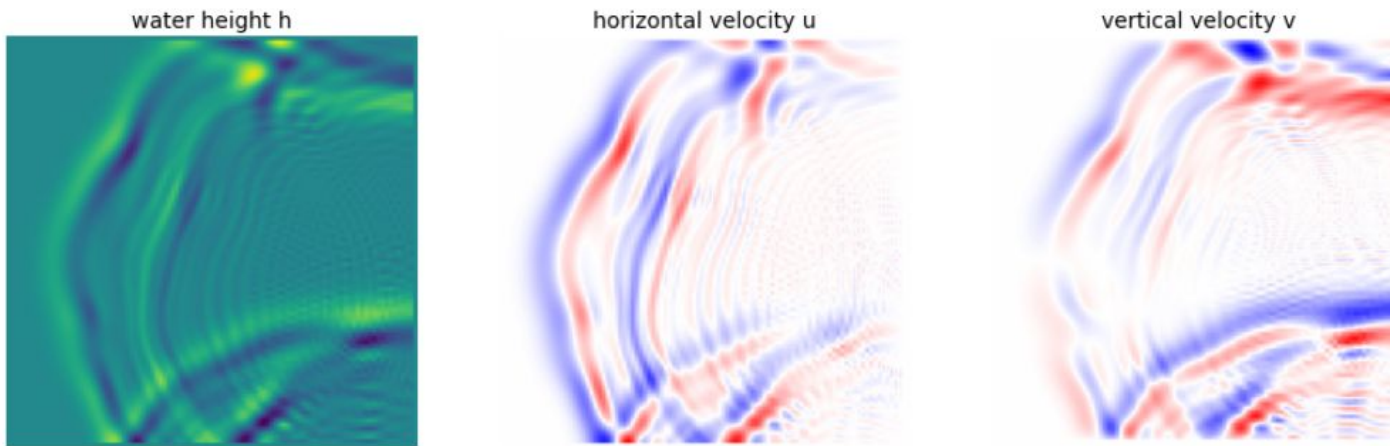
Emulator  
(convnet)



Difference



# Structure of SWE numerical solver



State update function of the numerical solver:

$$\left( h^{t+1}, u^{t+1}, v^{t+1} \right) = f \left( h^{t+1}, h^t, u^t, v^t \right) \quad \text{(semi-)implicit !}$$

(other fluid equations: similar for  $h$  = pressure)

More specifically here:

$$A h^{t+1} = b$$
$$A = A(h^t, u^t, v^t),$$
$$b = b(h^t, u^t, v^t)$$

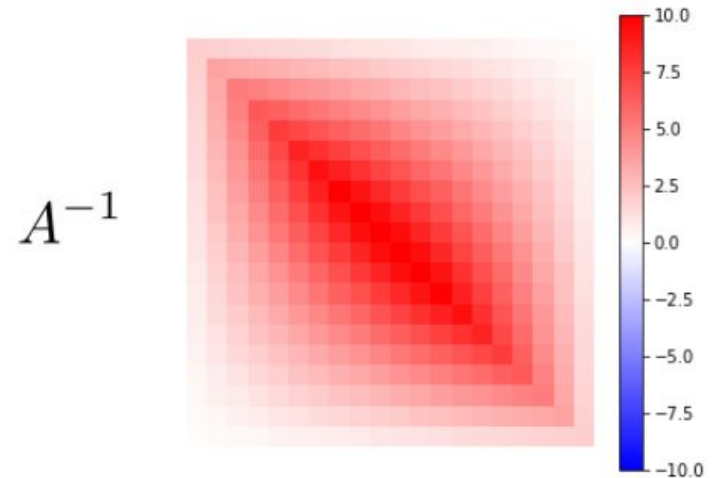
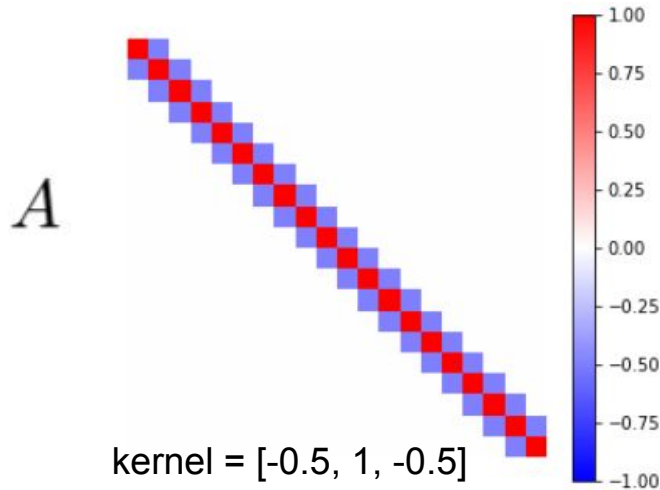
# Implicit schemes = non-local state update functions

$$A h^{t+1} = b$$

$$h^{t+1} = A^{-1}b$$

$A$  and  $b$  may generally be local functions (in  $h^t, u^t, v^t$ ),

but  $A^{-1}$  generally is not!





# Implicit schemes = non-local state update functions

$$A h^{t+1} = b$$

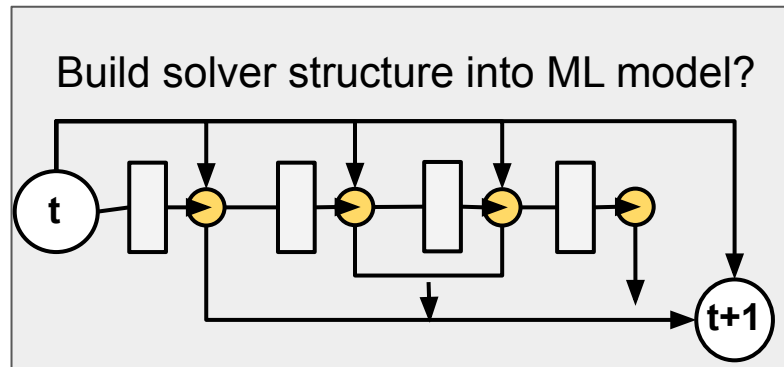
$$h^{t+1} = A^{-1}b$$

$A$  and  $b$  may generally be local functions (in  $h^t, u^t, v^t$ ),

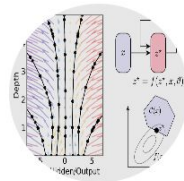
but  $A^{-1}$  generally is not!

We could try:

- non-local feed-forward networks  
(lots of network parameters)
- U-Nets  
(scaling with e.g. doubled resolution?)



# Linear implicit network layers



Deep Implicit Layers - Neural ODEs, Deep Equilibrium Models, and Beyond

<http://implicit-layers-tutorial.org/>

(linear) solve operations as part of a neural network model: for training need **backpropagation!**

1) forward pass:  $y = \text{layer}(x)$  , where  $A(x) y = b(x)$

2) backward pass:  $\frac{dy}{dx} = \frac{\partial y}{\partial \text{vec}(A)} \frac{d\text{vec}(A)}{dx} + \frac{\partial y}{\partial b} \frac{db}{dx}$

Implicit function theorem (IFT) gives:  $A \frac{\partial y}{\partial b} = I$   $A \frac{\partial y}{\partial a_j} = -y_j I$

$\frac{d\text{vec}(A)}{dx}$ ,  $\frac{db}{dx}$  depend on how we parametrize  $A(x), b(x)$

# Linear solvers

Any standard-library linear solver would do.

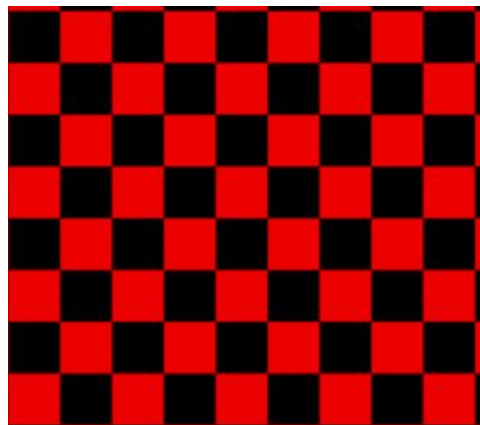
`numpy.linalg.solve`

`linalg.solve(a, b)`

For large system sizes ( $\gg 10^3$  grid points), use linear solvers that **scale** accordingly.

One possible choice: implement **multi-color Gauss-Seidel** in pytorch

- valid for an important class of **banded matrices**  $A$ .
- **iterative method**, initialize with estimate of solution  $\hat{y}$
- divide grid points into colors (e.g. red vs black)
- iterations only require dot products  $A \hat{y} \rightarrow \text{GPU !}$



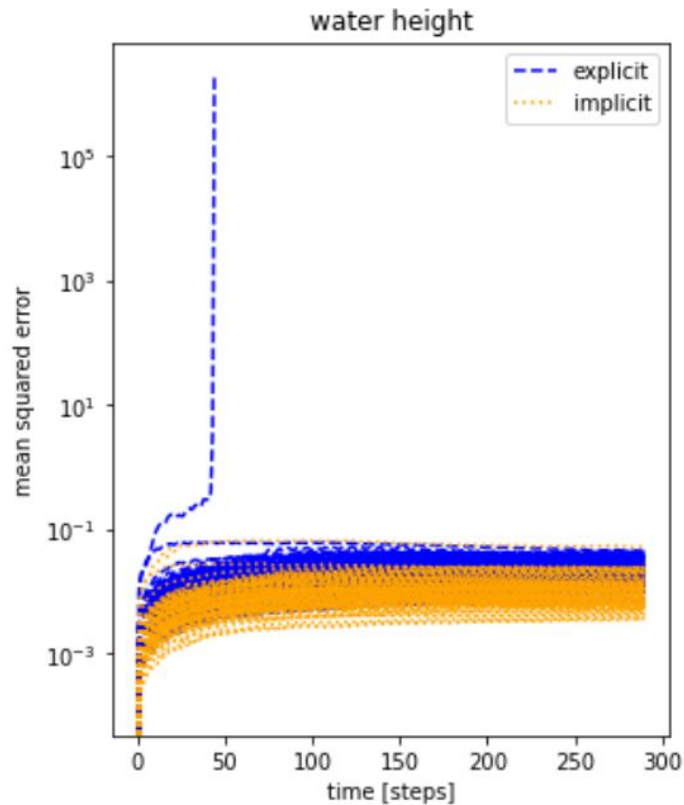
# Linear implicit layer in simple pytorch

[ this is essentially a re-write of torch.linalg.solve()  
in pure pytorch (no c++) and for custom solvers ]

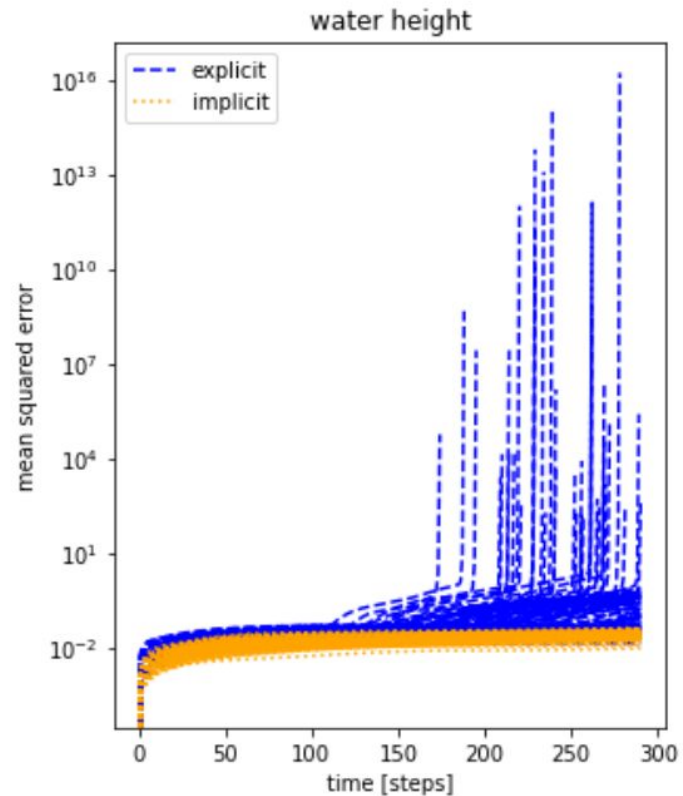
```
1 import torch
2 from utils import linear_solve
3
4 class LinearSolve(torch.autograd.Function):
5
6     @staticmethod
7     def forward(ctx, b, A):
8         y = linear_solve(A, b) solve #1
9
10        ctx.save_for_backward(A, y) # torch.autograd logistics
11
12        return y
13
14    @staticmethod
15    def backward(ctx, grad_output):
16
17        A, y = ctx.saved_tensors # torch.autograd logistics
18
19        z = linear_solve(A.T, grad_output) solve #2
20        # tranpose of A
21        # backward gradient from reverse-mode differentiation : dL / dy
22        )
23
24        grad_input = (z,
25                      -torch.bmm(z, y.transpose(-2, -1)) # dL/dy * dy / db
26                      )
27
28        return grad_input
29
30
31 class LinearLayer(torch.nn.Module):
32
33     def __init__(self, comp_bA):
34         self.comp_bA = comp_bA # function that returns (b, A) as function of x
35
36     def forward(self, x):
37         b, A = self.comp_bA(x) b, A = f(x) and (db/dx, dA/dx)
38         y = LinearSolve.apply(b, A)
39         return y
```

# Results on SWE

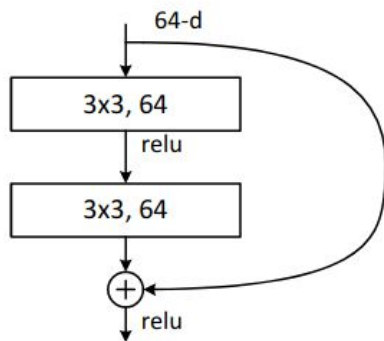
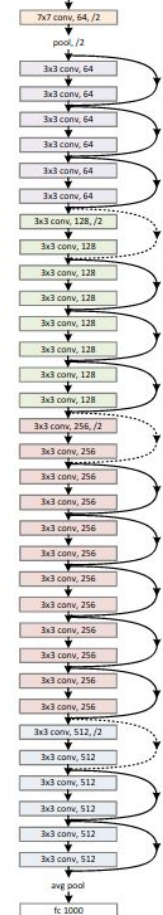
Solver step length:  $dt = 300s$



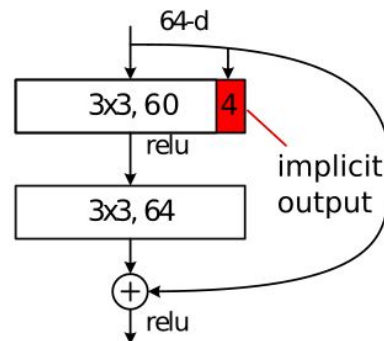
$dt = 600s$



# Composability: implicit output in neural architectures



Residual block: conv & relu



Residual block: stack(conv, implicit) & relu

## Deep Residual Learning for Image Recognition

Kaiming He Xiangyu Zhang Shaoqing Ren Jian Sun

Microsoft Research

{kahe, v-xiangz, v-shren, jiansun}@microsoft.com

arXiv:1512.03385v1 [cs.CV] 10 Dec 2015



# Outlook

- (linear) implicit layers as flexible building blocks of neural networks
- scaling to large systems (Gauss-Seidel only a start: multigrid methods)
- expressivity of linear implicit layer & parametrization of  $A(x)$

# Thank you!

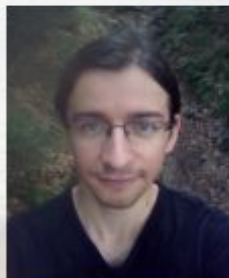
**HELMHOLTZ**AI | ARTIFICIAL INTELLIGENCE  
COOPERATION UNIT

Model-driven Machine Learning Group

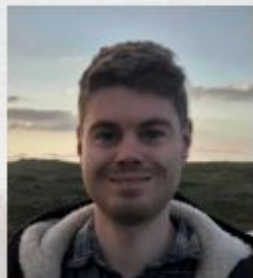
<http://m-dml.org>



David Greenberg



Marcel Nonnenmacher



Tobias Machnitzki



Shivani Sharma



Vadim Zinchenko



Kubilay Demir