



PUC
RIO

MACH
2025



JOHNS HOPKINS
WHITING SCHOOL
of ENGINEERING

Scalable Multi-GPU Training of Neural Operators: Advancing Generalization in High- Dimensional Physical Systems

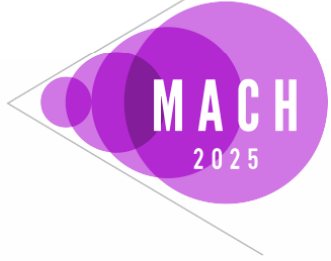
Luis Felipe dos Santos, PhD Candidate

Dibakar Roy Sarkar, PhD Student

Deane Roehl, PhD

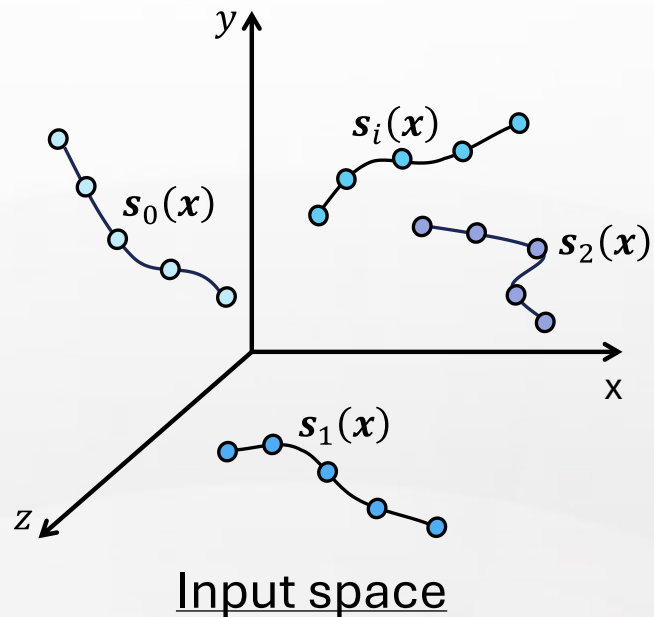
Somdatta Goswami, PhD

AGENDA



- ❑ INTRODUCTION
- ❑ MOTIVATION
- ❑ DATA PARALLEL OPERATOR LEARNING
- ❑ RESULTS
- ❑ FUTURE PERSPECTIVES

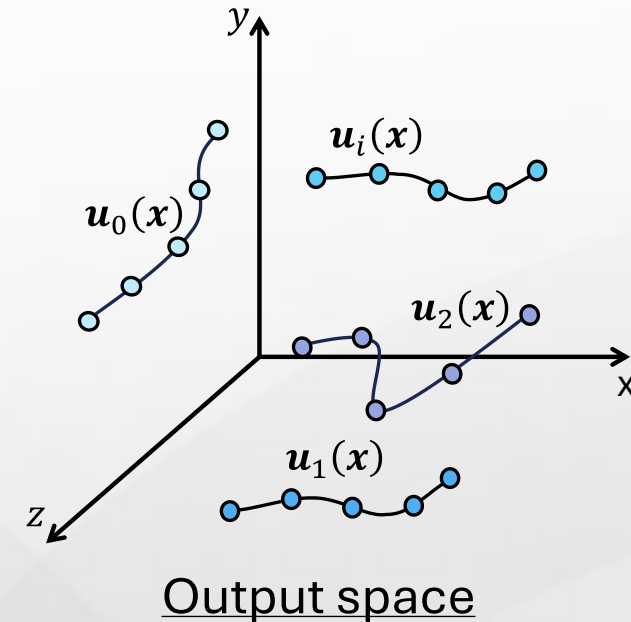
- Infinite - dimensional mapping between spaces



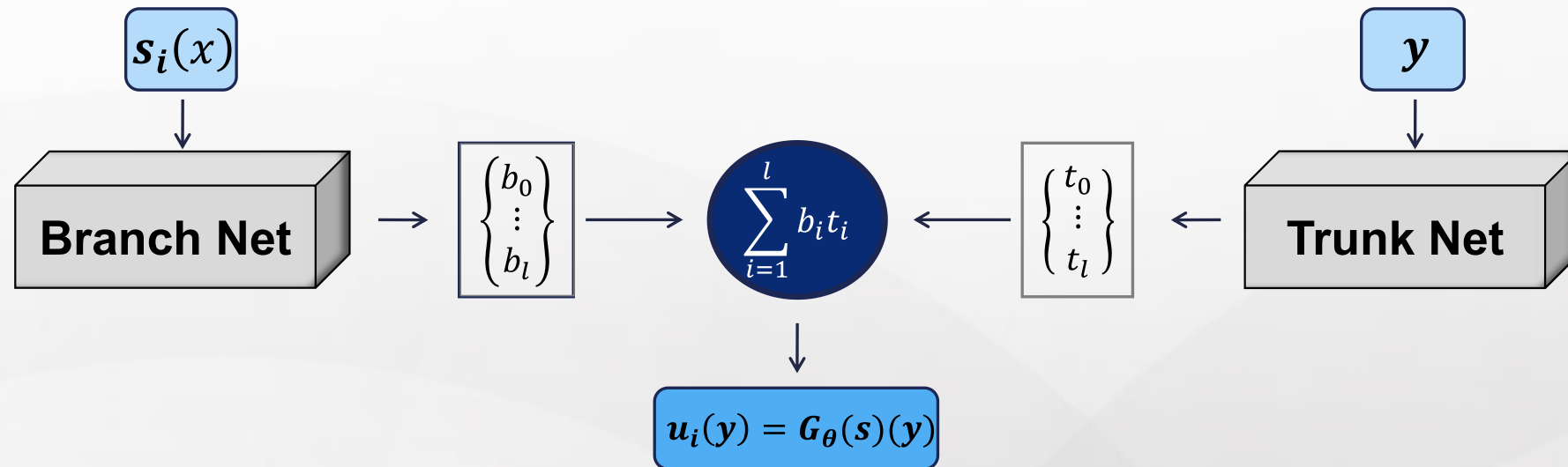
Operator Learning



$$G_{\theta} : s_i \rightarrow u_i$$



- ❑ Leverages the universal approximation capabilities of neural networks to enable efficient Deep Operator Network learning with a flexible architecture.



□ ADVANTAGES

- Real-time inference
- Reliable surrogate model

□ CHALLENGES

- Down-sampling input-output spaces miss fine details, reducing accuracy.
- Latent space learning is ineffective for non-dissipative systems.
- Small mini-batches limit the network's generalization ability.

❑ SOLUTION

- ❑ Developing a scalable framework for data-parallel operator learning to handle high-dimensional PDEs using JAX, optimized for multi-node and multi-GPU HPC

DATA PARALLEL OPERATOR LEARNING



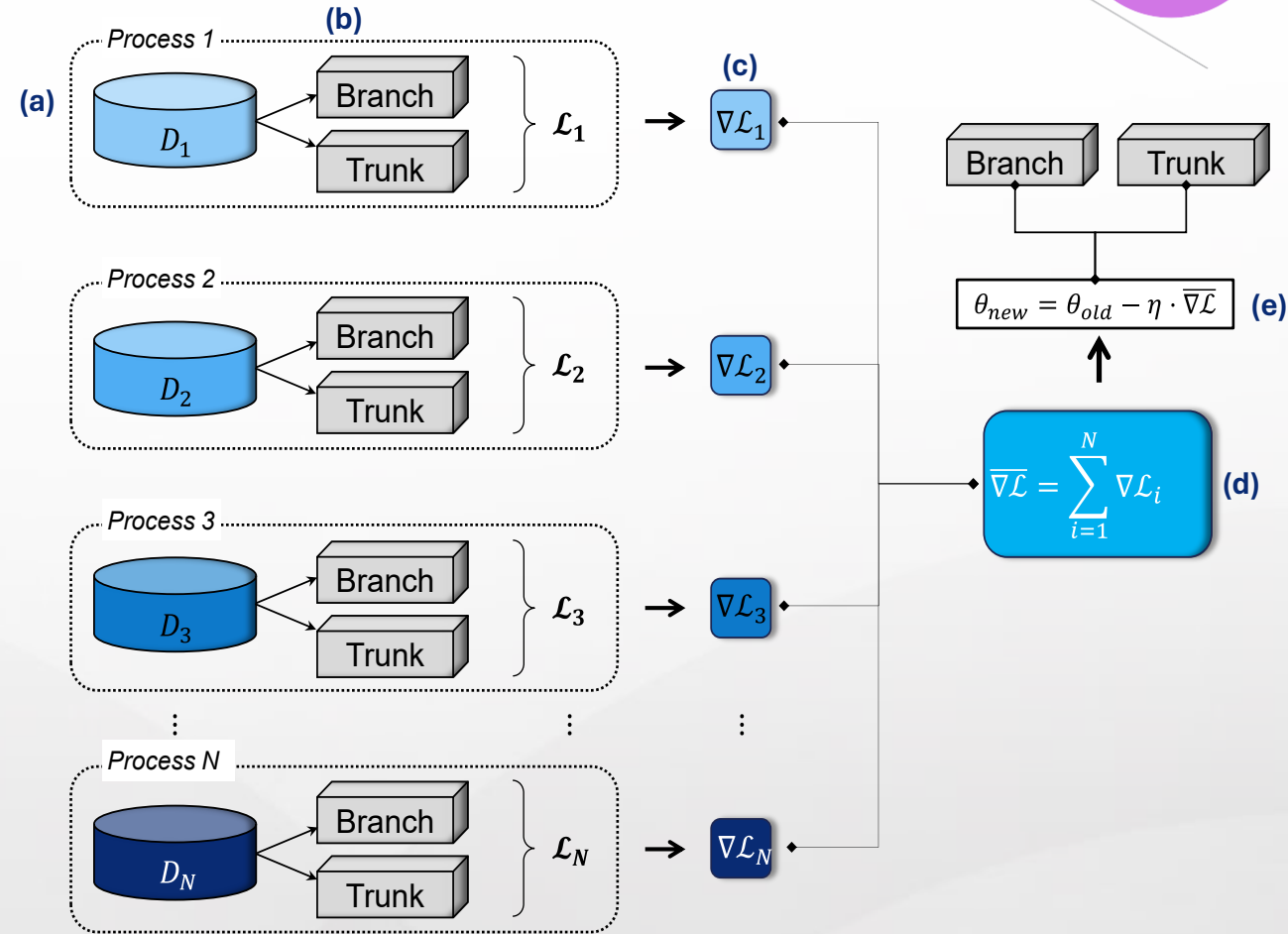
GENERAL FRAMEWORK:

Model:

- Create and replicate trunk and branch states

Train step:

- Data treatment:
 - Sharding inputs and outputs
- Get local predictions – losses – gradients
- Averaging gradients and sum losses
 - All-Reduce mean and Reduce sum
- Local parameters updating
- Go to step i. and repeat



□ SCALING METRICS:

- Speed up (S_{up})

$$S_{up} = \frac{T_1}{T_N}$$

- Efficiency (E_f)

$$E_f = N \cdot S_{up} = \frac{N \cdot T_1}{T_N}$$

T_1 : Time for processing considering one device

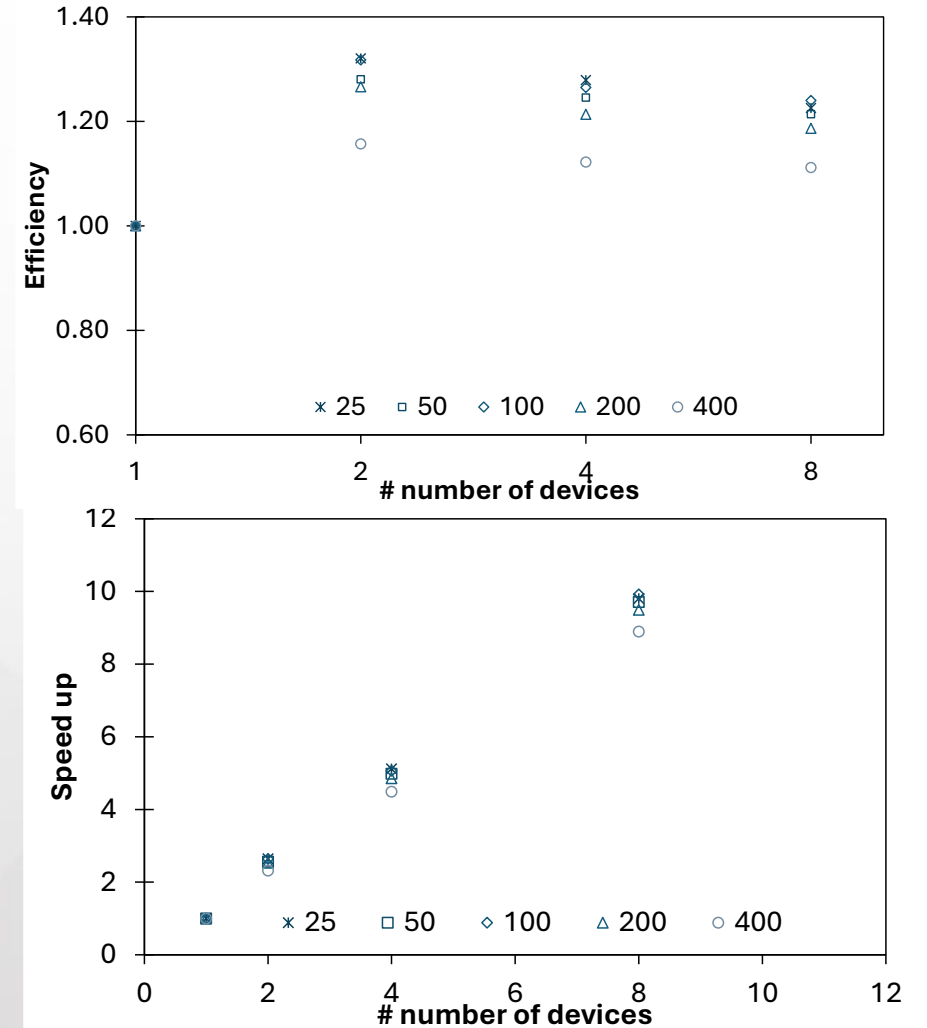
T_N : Time for processing considering N devices

Results for **Branch Parallel** training (**Physics-informed Burgers**):

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

- Input space $\rightarrow u(t = 0, x) \therefore (2000, 101)$
- Output space $\rightarrow u(t, x) \therefore (2000, 101, 101)$

$$G_{\theta}(u(t = 0, x))(t, x) \rightarrow u(t, x)$$

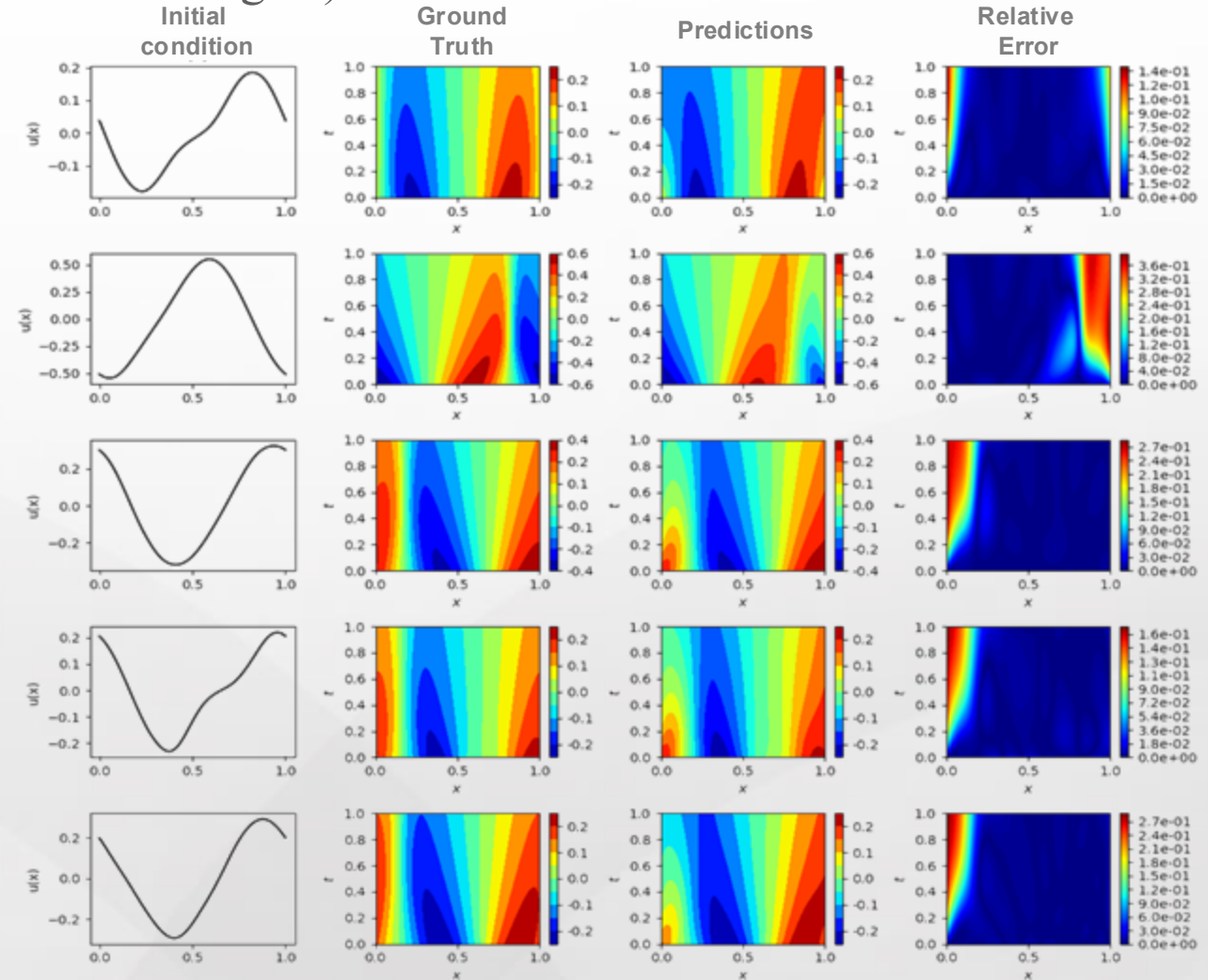


Results for Branch Parallel training (Physics-informed Burgers):

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

- Input space $\rightarrow u(t=0, x) \therefore (2000, 101)$
- Output space $\rightarrow u(t, x) \therefore (2000, 101, 101)$

$$G_{\theta}(u(t=0, x))(t, x) \rightarrow u(t, x)$$



RESULTS

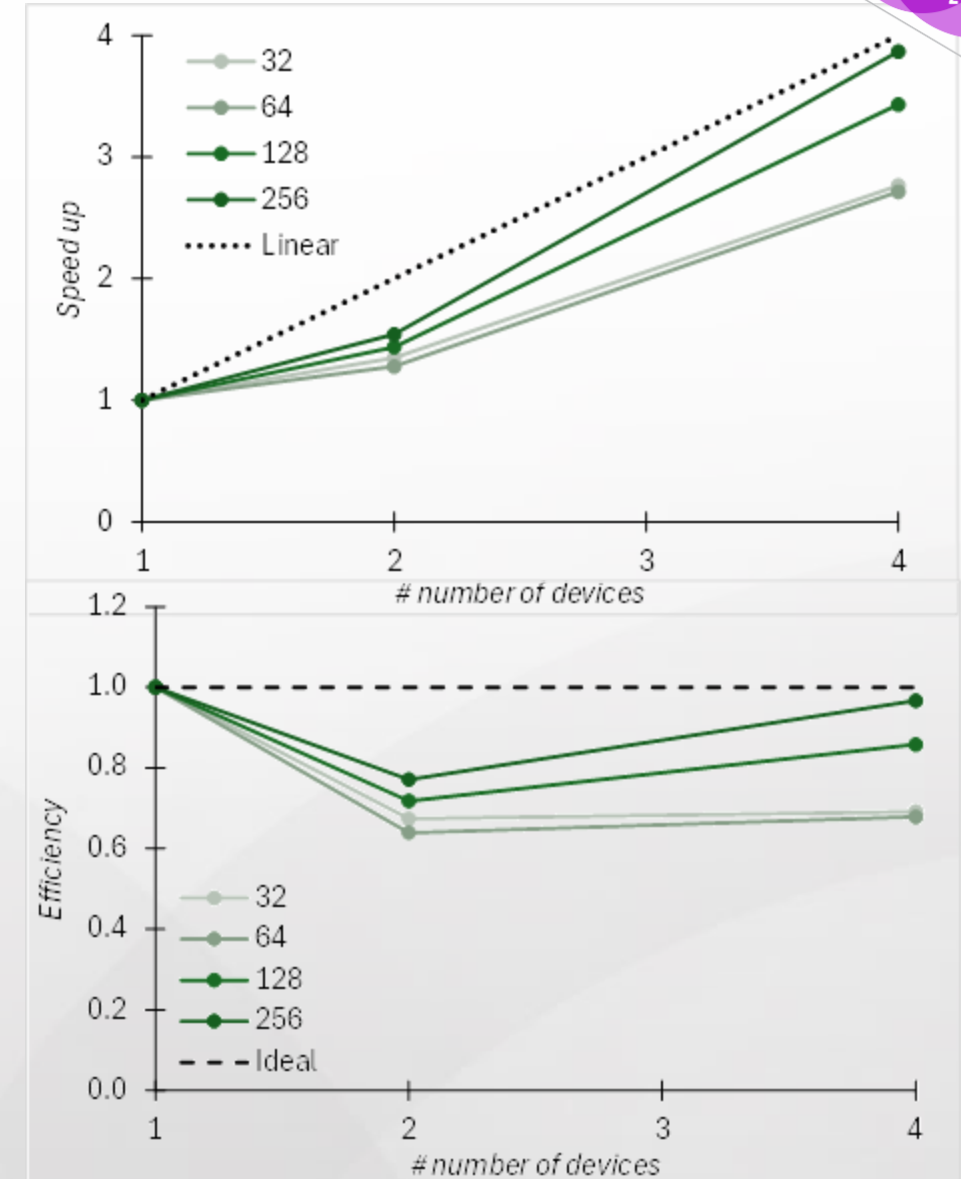


Results for **Trunk Parallel** training (Data-driven Darcy):

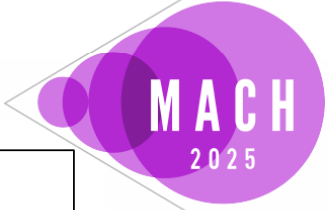
$$\frac{\partial \mathbf{P}}{\partial t} - K(x, y) * \nabla \mathbf{P} = \mathbf{q}$$

- Input space $\rightarrow K(x, y) \therefore (1000, 100, 100)$
- Output space $\rightarrow \mathbf{P}(t, x, y) \therefore (1000, 72, 100, 100)$

$$G_{\theta}(K(x, y))(t, x, y) \rightarrow \mathbf{P}(t, x, y)$$



RESULTS



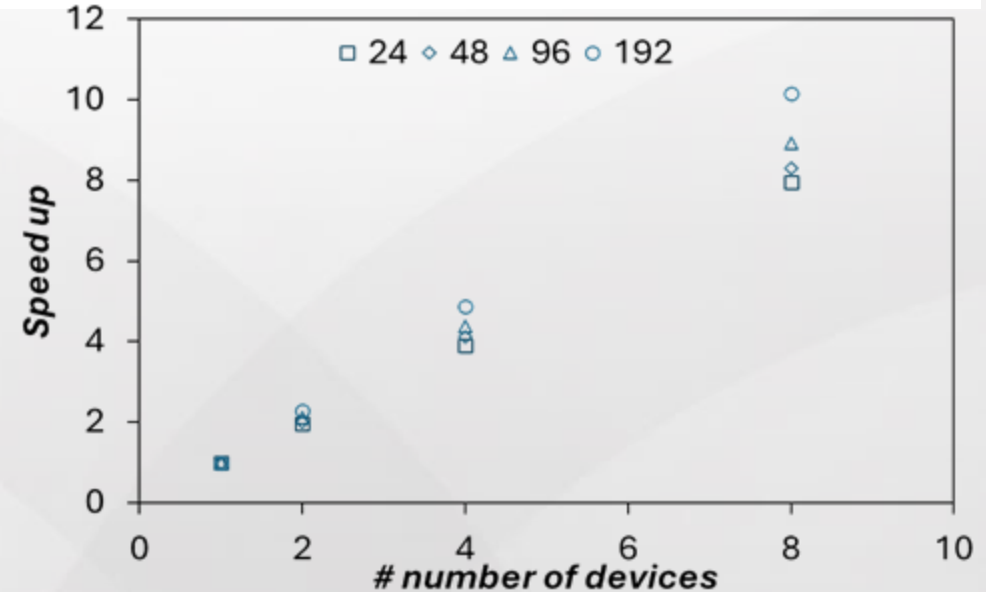
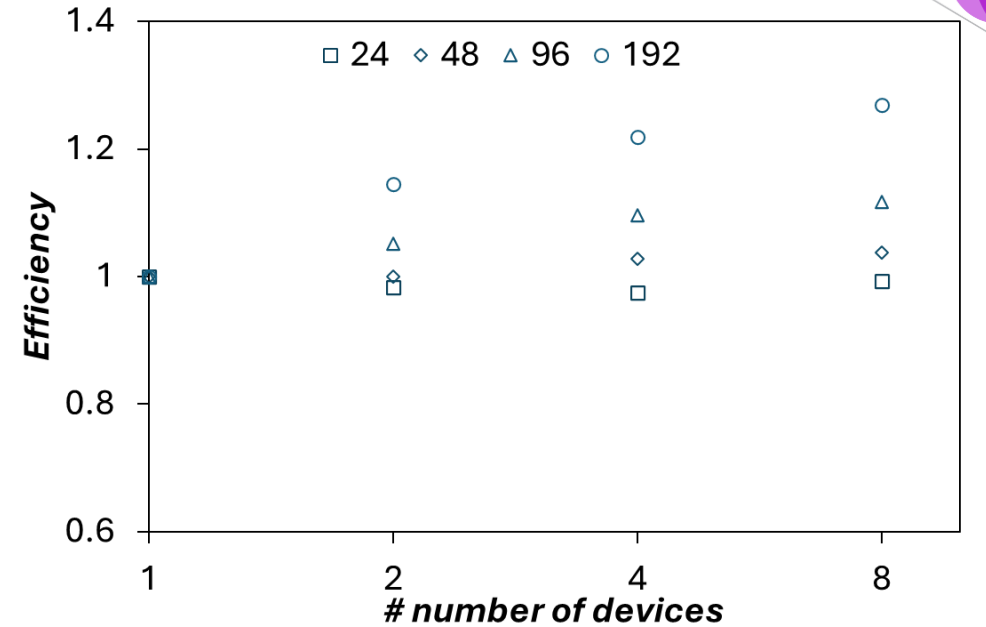
- Results for **Branch and Trunk Parallel** training
(Data-driven shallow water):

$$\frac{DV}{Dt} = -f\mathbf{k} \times \mathbf{V} - g\nabla h + \nu\nabla^2 \mathbf{V}$$

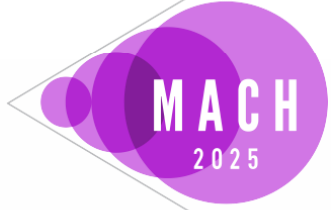
$$\frac{Dh}{Dt} = -\mathbf{h}\nabla \cdot \mathbf{V} + \nu\nabla^2 h$$

- Input space $\rightarrow \mathbf{h}(t=0, \mathbf{x}) \therefore (256, 256)$
- Output space $\rightarrow \mathbf{V}(t, \mathbf{x}) \therefore (72, 256, 256)$

$$G_{\theta}(\mathbf{h}(t=0, \mathbf{x}))(t, \mathbf{x}) \rightarrow \mathbf{V}(t, \mathbf{x})$$



RESULTS



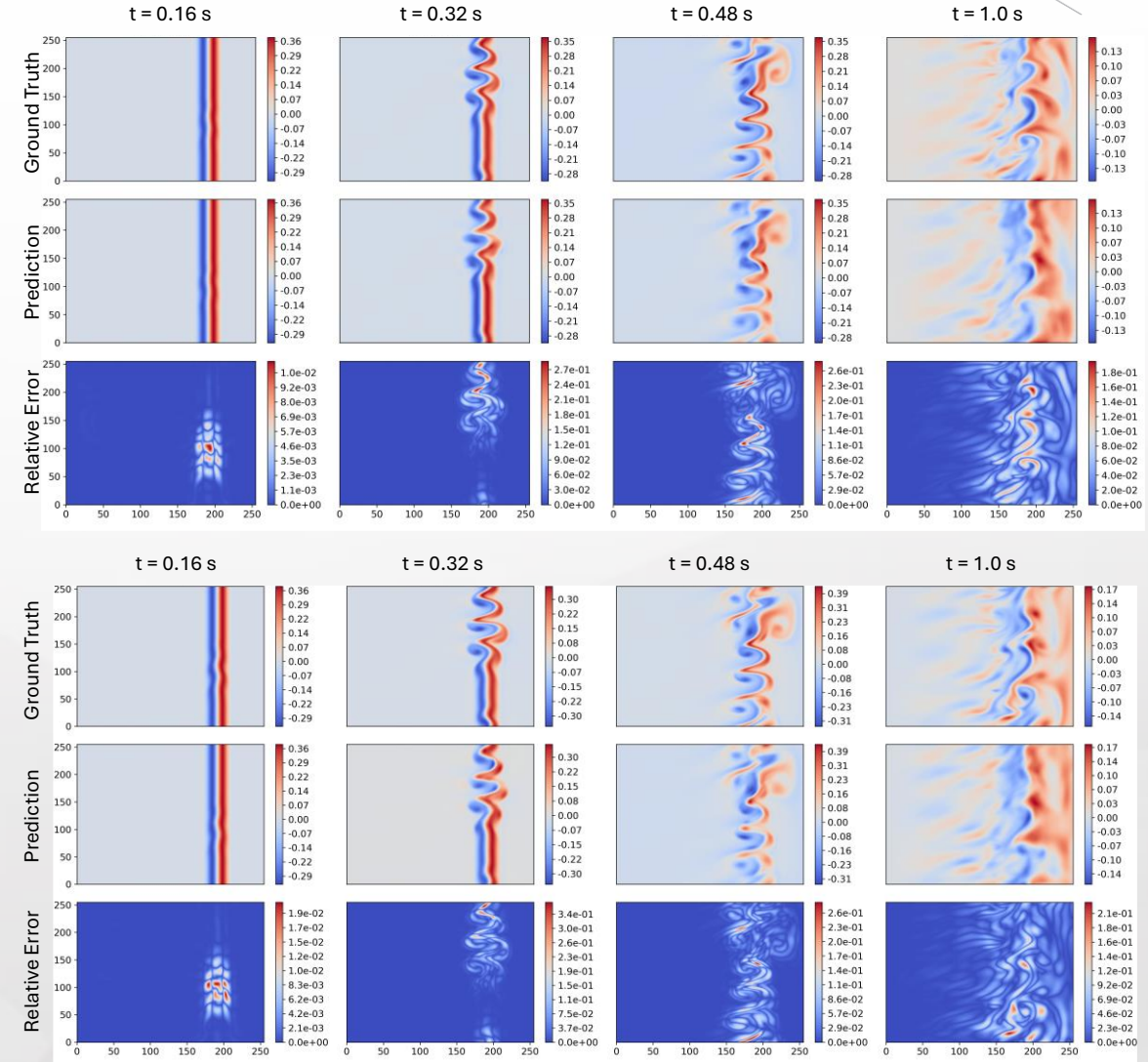
- Results for **Branch and Trunk Parallel** training
(Data-driven shallow water):

$$\frac{DV}{Dt} = -f\mathbf{k} \times \mathbf{V} - g\nabla h + \nu\nabla^2 \mathbf{V}$$

$$\frac{Dh}{Dt} = -\mathbf{h}\nabla \cdot \mathbf{V} + \nu\nabla^2 h$$

- Input space $\rightarrow \mathbf{h}(t = 0, \mathbf{x}) \therefore (256, 256)$
- Output space $\rightarrow \mathbf{V}(t, \mathbf{x}) \therefore (72, 256, 256)$

$$G_{\theta}(\mathbf{h}(t = 0, \mathbf{x}))(t, \mathbf{x}) \rightarrow \mathbf{V}(t, \mathbf{x})$$



- ❑ The **efficiency** of the proposed DP training is evident in its superior performance, which **exceeds the minimum admissible (0.8)** for all cases and is particularly pronounced (superior to 1) in the majority of the evaluated scenarios.
- ❑ The **speed-up** consistently **surpasses** the **ideal (linear) scaling** and exhibits a marked enhancement as the number of devices increases, thereby demonstrating the effective management of high-dimensional problems.
- ❑ **Designed** to be **generic**, allowing **easy extension** to various types of **neural networks** and **operator learning** being **adaptable to** different numbers of **nodes and GPUs**, ensuring **scalability** and **flexibility**. Available at <https://github.com/Centrum-IntelliPhysics>
- ❑ **Future directions** include the integration of **domain decomposition** with the proposed data parallel operator learning.

ACKNOWLEDGMENTS



Thank You!

FUNDING



NAIRR Pilot

