



南京航空航天大学
NANJING UNIVERSITY OF AERONAUTICS AND ASTRONAUTICS

CVIT: CONTINUOUS VISION TRANSFORMER FOR OPERATOR LEARNING

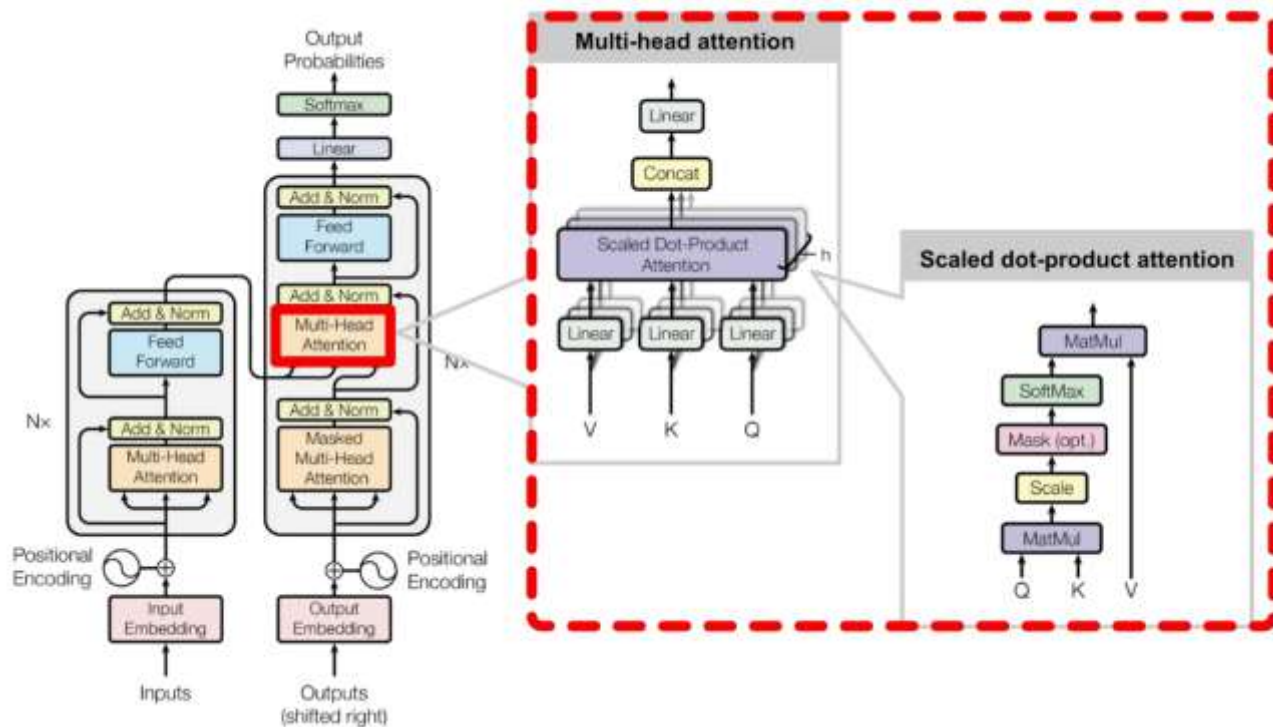
Sifan Wang¹, Jacob H. Seidman², Shyam Sankaran³, Hanwen Wang², George J. Pappas⁴ Paris
Perdikaris³

ICLR 2025

The emergence of neural operators has driven advancements in **solving partial differential equations** (PDEs) and **modeling physical systems**. Recent developments have focused on enhancing the **expressiveness** and **efficiency** of these models. For instance, extensions of the Fourier Neural Operator (FNO) (Li et al., 2021) have explored factorized representations (Tran et al., 2021) to reduce computational complexity while maintaining performance. Similarly, wavelet-based approaches (Gupta et al., 2021; Tripura & Chakraborty, 2022) have been proposed to capture multi-scale features more effectively.

Transformer-Based Operator Learning. The **self-attention mechanism** central to transformers offers a powerful tool for capturing long-range dependencies in spatial and temporal data, a crucial aspect in many PDE-governed systems. For example, **OFormer** (Li et al., 2022) introduced a novel way to embed continuous input functions and output queries into a transformer architecture. This approach demonstrated the potential of transformers in handling the inherent continuity of physical problems. **GNOT** (Hao et al., 2023) further extended this idea by incorporating graph structures, allowing for more flexible representation of input functions and improved handling of irregular domains.

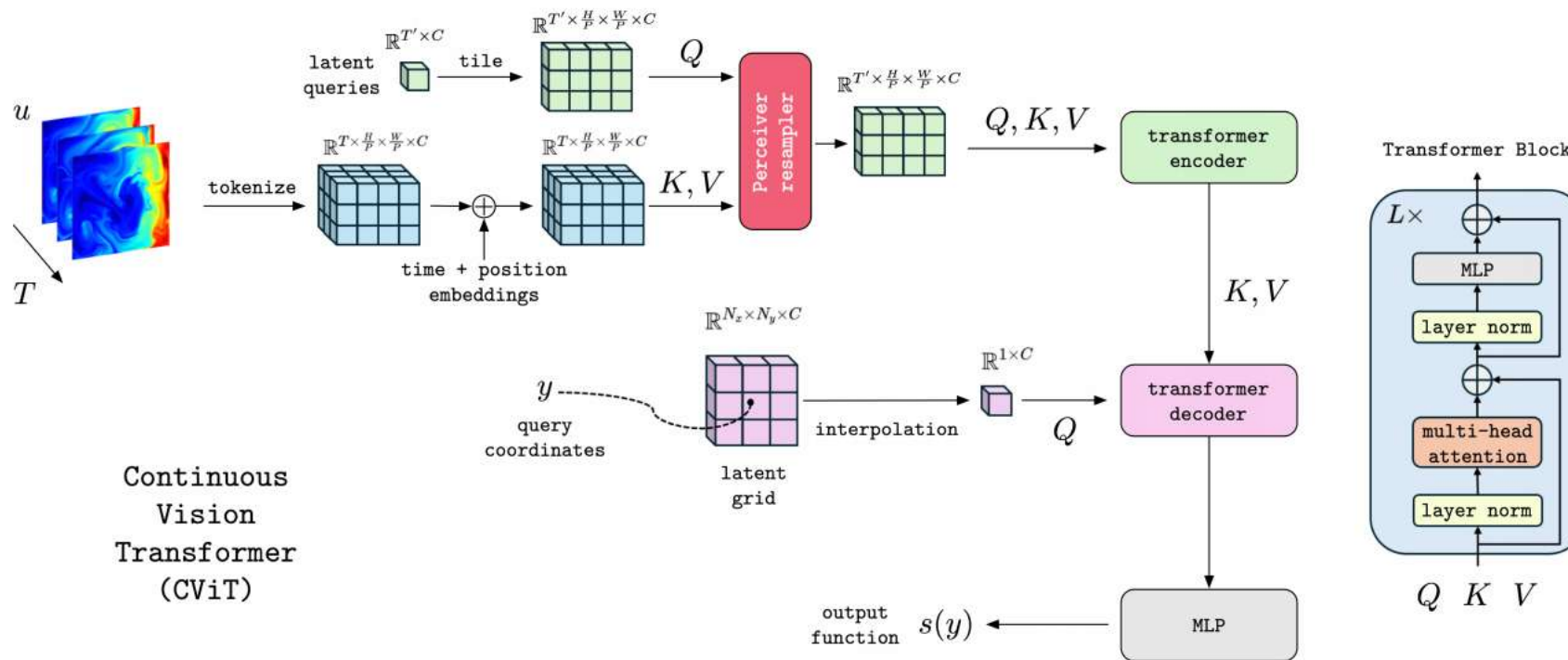
Transformer:



$$\text{Attention}(Q, K, V) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

- $Q \in \mathbb{R}^{N \times d_k}$: 查询矩阵 (Query), 由查询序列 X_q 经过投影得到。
- $K \in \mathbb{R}^{M \times d_k}$: 键矩阵 (Key), 由输入序列 X_k 经过投影得到。
- $V \in \mathbb{R}^{M \times d_v}$: 值矩阵 (Value), 由输入序列 X_v 经过投影得到。
- d_k : 键向量的维度。
- N : 查询序列的长度 (查询点个数)。
- M : 输入序列的长度 (键/值的个数)。

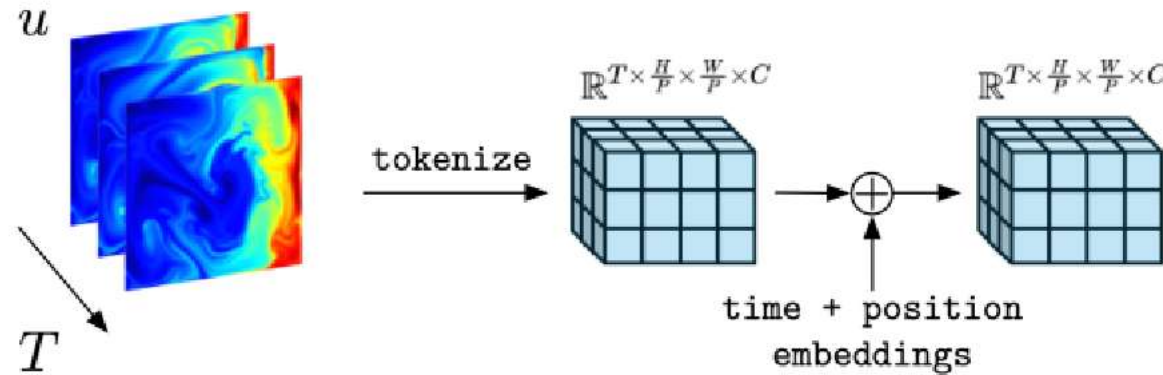
CViT:



Continuous Vision Transformer (CViT) Architecture: CViT consists of the following components:

- (1) **Spatio-temporal patch** embeddings to extract localized features.
- (2) **A temporal aggregation module** based on the Perceiver architecture, which captures temporal correlations to compresses tokens along the time axis.
- (3) **A Transformer encoder** that captures multi-scale spatial dependencies via self-attention layers.
- (4) **A novel grid-based positional encoding scheme for query coordinates**, allowing for flexible output representation and interpolation.
- (5) **A cross-attention decoder** that integrates information from the input function with query coordinates.

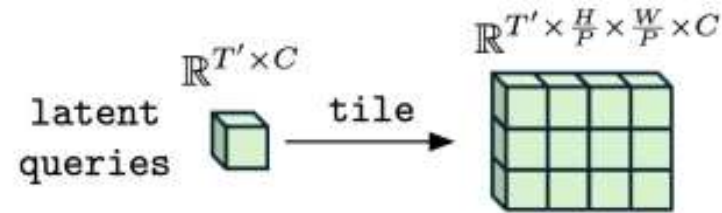
(1) Spatio-temporal patch embeddings to extract localized features.



The vision transformer encoder takes as input a gridded representation of the input function u , yielding a spatio-temporal data tensor $u \in \mathbb{R}^{T \times H \times W \times D}$ with D channels. We patchify our inputs into 3D tokens $u_p \in \mathbb{R}^{T \times \frac{H}{P} \times \frac{W}{P} \times C}$ by tokenizing each 2D spatial frame independently, following the process used in standard Vision Transformers. We then add trainable 1D temporal and 2D spatial positional embeddings to each token:

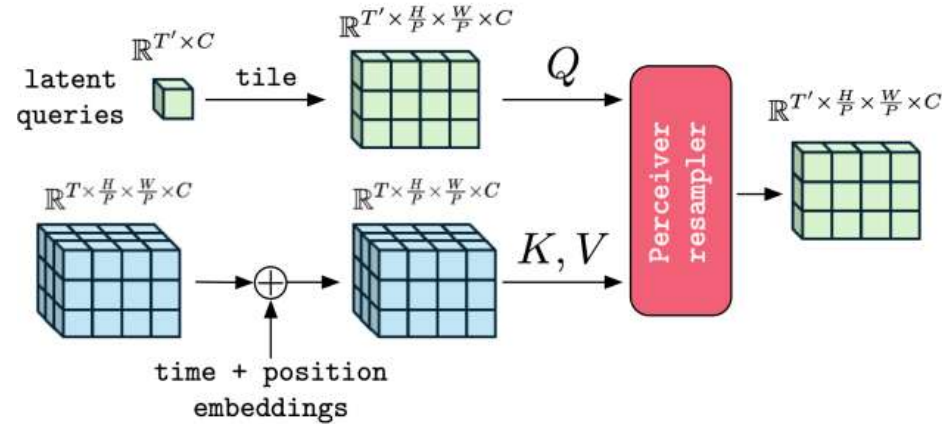
$$u_{pe} = u_p + PE_t + PE_s, \quad PE_t \in \mathbb{R}^{T \times 1 \times 1 \times C}, \quad PE_s \in \mathbb{R}^{1 \times \frac{H}{P} \times \frac{W}{P} \times C}$$

(2) A temporal aggregation module based on the Perceiver architecture.



To reduce computational cost, we use a temporal aggregation layer based on the Perceiver architecture. Specifically, we learn a **pre-defined number of latent input queries** $z \in \mathbb{R}^{T' \times C}$. These queries serve as **learnable parameters** in a cross-attention module, which processes the visual features of our input.

(2) A temporal aggregation module based on the Perceiver architecture.



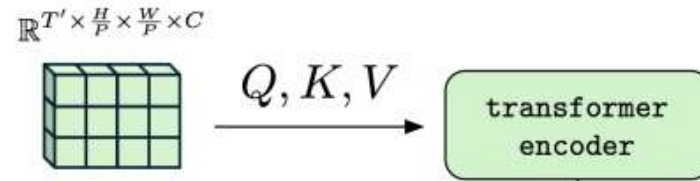
The Perceiver module operates on flattened visual features $u_f \in R^{(\frac{H}{P} \times \frac{W}{P}) \times T \times C}$ obtained by flattening the positionally encoded features u_{pe} . Through this cross-attention mechanism, we aggregate the temporal information into a compact latent representation $z_{agg} \in R^{(\frac{H}{P} \times \frac{W}{P}) \times T' \times C}$ as:

$$z' = \hat{z} + MHA(LN(\hat{z}), LN(u_f), LN(u_f)),$$

$$z_{agg} = z' + MLP(LN(z')).$$

Here, z is initialized by a **unit Gaussian distribution** and $\hat{z} \in R^{(\frac{H}{P} \times \frac{W}{P}) \times T' \times C}$ is obtained by tiling the latent query z ($\frac{H}{P} \times \frac{W}{P}$) times. Besides compressing tokens over the time axis, this module also enables the model to handle inputs with a variable number of time steps.

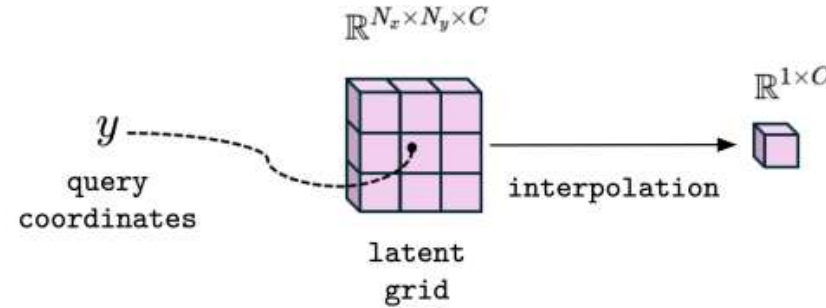
(3) A Transformer encoder that captures multi-scale spatial dependencies via self-attention layers.



We then process the aggregated tokens z_{agg} using a sequence of L pre-norm Transformer blocks.

$$\begin{aligned} z_0 &= LN(z_{agg}), \\ z'_l &= MSA(LN(z_{l-1})) + z_{l-1}, \quad l = 1, 2, \dots, L \\ z_l &= MLP(LN(z'_l)) + z'_l, \quad l = 1, 2, \dots, L \end{aligned}$$

(4) A novel grid-based positional encoding scheme for query coordinates.

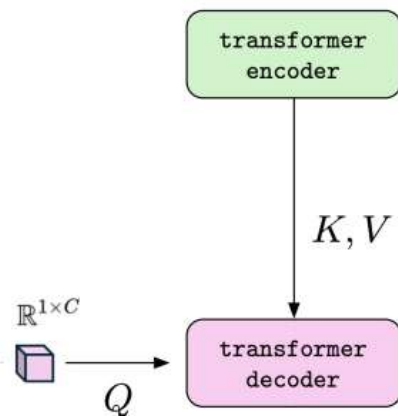


To enable continuous evaluation of outputs, we design a novel and efficient coordinate embedding to capture the fine-scale features of the target functions. Specifically, we create a **uniform grid** $\{y_{ij}\} \subset [0,1]^2$, for $i = 1, \dots, N_x$ and $j = 1, \dots, N_y$, along with **associated trainable latent grid features** $x \in \mathbb{R}^{N_x \times N_y \times C}$. For a query point $y \in \mathbb{R}^2$, we compute a Nadaraya-Watson interpolant over grid latent features:

$$\mathbf{x}' = \sum_{i=1}^{N_x} \sum_{j=1}^{N_y} w_{ij} \mathbf{x}_{ij}, \quad w_{ij} = \frac{\exp(-\beta \|y - \mathbf{y}_{ij}\|^2)}{\sum_{ij} \exp(-\beta \|y - \mathbf{y}_{ij}\|^2)}.$$

Here $\beta > 0$ is a **hyperparameter** that determines the locality of the interpolated features. Specifically, larger values of β yield more localized weight distributions w_{ij} , resulting in a higher-frequency interpolant that captures finer-scale variations. Conversely, smaller β values produce a smoother interpolant by averaging over a broader neighborhood of points.

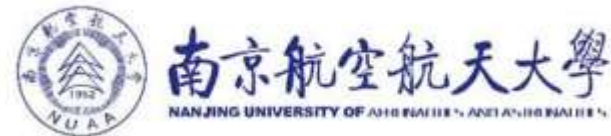
(5) A cross-attention decoder that integrates information from the input function with query coordinates.



The interpolated grid feature $x_0 = \mathbf{x}' \in \mathbb{R}^{1 \times C}$ is used as query input to a transformer decoder. This decoder uses the output of the vision transformer encoder z_L as keys and values in a cross-attention mechanism:

$$\begin{aligned} \mathbf{x}'_k &= \mathbf{x}_{k-1} + \text{MHA}(\text{LN}(\mathbf{x}_{k-1}), \text{LN}(\mathbf{z}_L), \text{LN}(\mathbf{z}_L)), & k = 1 \dots K, \\ \mathbf{x}_k &= \mathbf{x}'_k + \text{MLP}(\text{LN}(\mathbf{x}'_k)), & k = 1 \dots K. \end{aligned}$$

Experiment



For all experiments, unless otherwise stated, we use a patch size of 8×8 for tokenizing inputs.

Table 1: Details of Continuous Vision Transformer model variants.

Model	Encoder layers	Embedding dim	MLP width	Heads	# Params
CViT-S	5	384	384	6	13 M
CViT-B	10	512	512	8	30 M
CViT-L	15	768	1536	12	92 M

Model	# Params	NS	CNS	DR
FNO	0.5 M	9.12 %	9.60 %	12.00 %
FFNO	1.3 M	8.39 %	5.20 %	5.71 %
GK-T	1.6 M	9.52 %	3.77 %	3.59 %
GNOT	1.8 M	17.20 %	4.20 %	3.11 %
Oformer	1.9 M	13.50 %	6.25 %	1.92 %
DPOT-Ti	7 M	12.50 %	3.97 %	3.21 %
MPP-S	30M	-	-	1.12 %
DPOT-S	30 M	9.91 %	3.37 %	3.79 %
MPP-L (Pre-trained)	400 M	-	-	0.98 %
DPOT-L (Pre-trained)	500 M	7.98 %	2.16 %	2.32 %
DPOT-L (Fine-tuned)	500 M	2.78 %	1.31 %	0.73 %
DPOT-H (Pre-trained)	1.03 B	3.79 %	1.80 %	1.91 %
CViT-S	13 M	3.75 %	2.71 %	1.13 %
CViT-B	30 M	3.18 %	1.99 %	1.11 %
CViT-L	92 M	2.35 %	1.29 %	0.68 %

Experiment

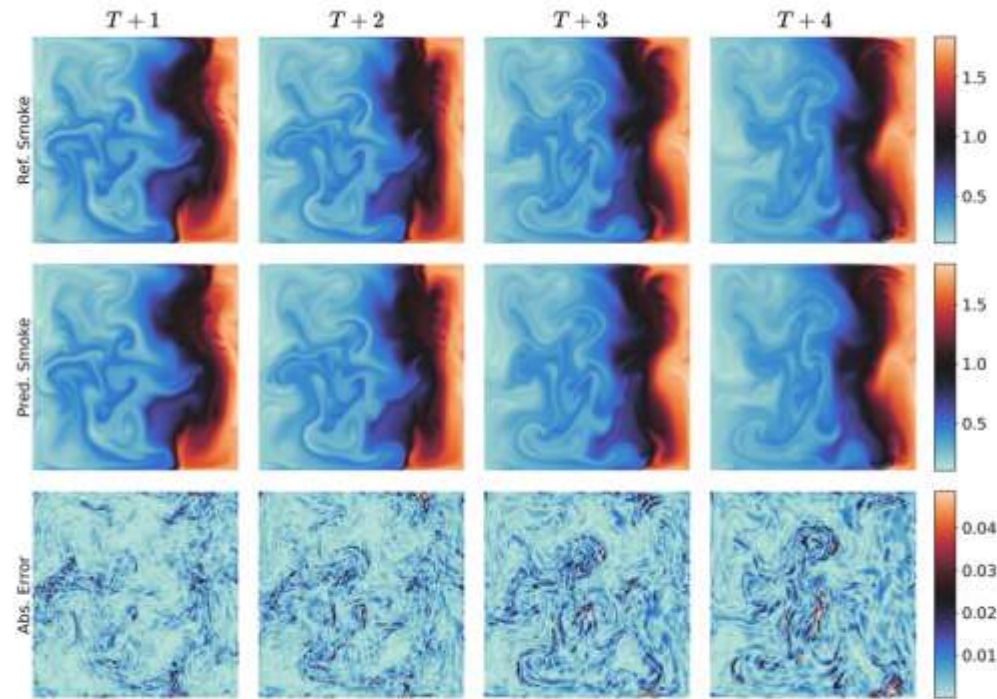


Figure 5: *Incompressible Navier-Stokes benchmark (NS)*. Representative CViT rollout predictions of the passive scalar field, and point-wise error against the ground truth.

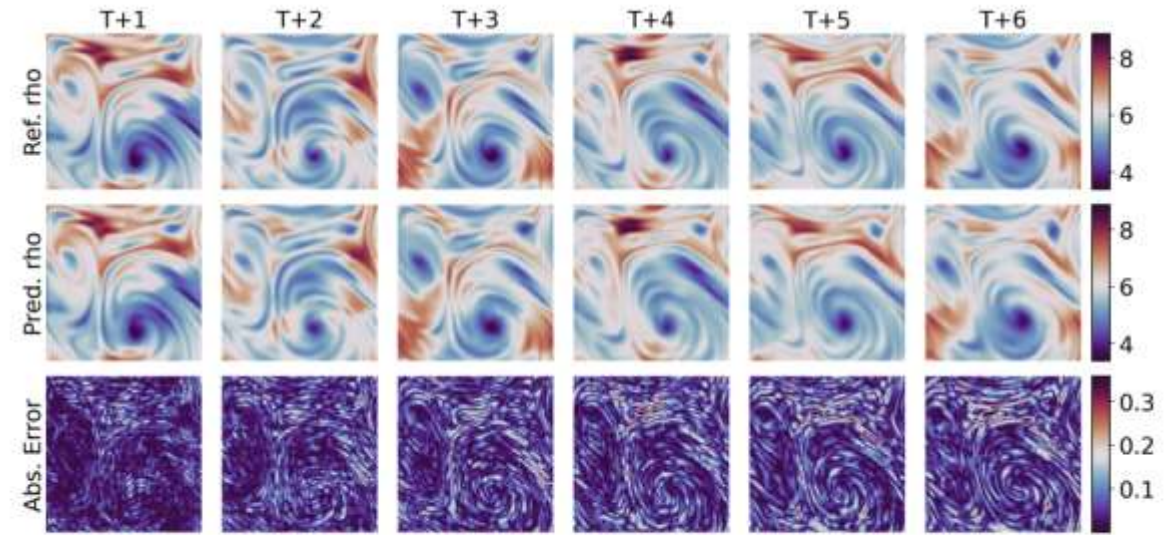


Figure 6: *Compressible Navier-Stokes Benchmark (CNS)*. Representative CViT rollout prediction of the density field ρ , and point-wise error against the ground truth.

Thanks