# Progressive Fourier Neural Representation for Sequential Video Compilation
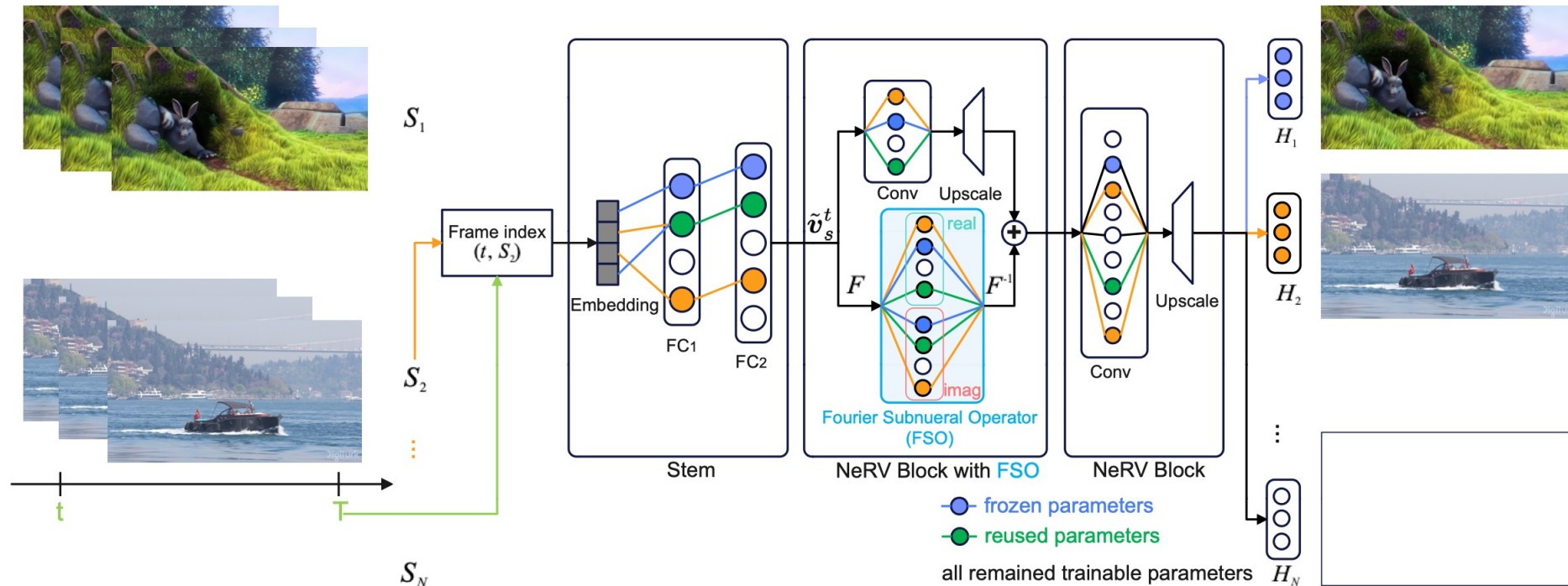
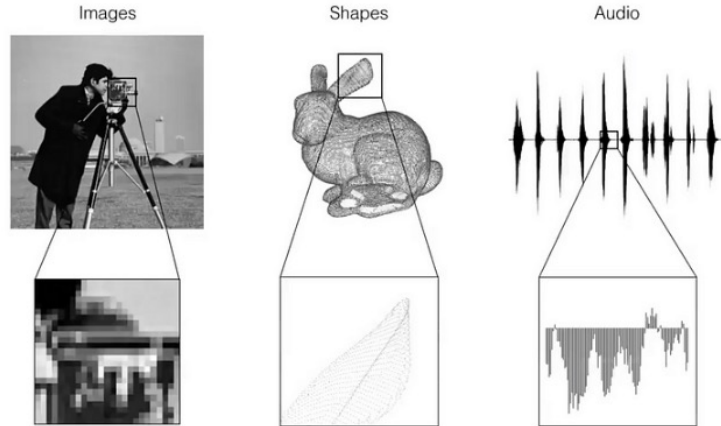**Haeyong Kang**, Jaehong Yoon, DaHyun Kim, Sungju Hwang, and Chang D. Yoo

**Motivation**:

(1) It has not been explored sequential video so far in Continual Learning.

(2) There is also no robust continual baseline in Video Continual Learning.



➜ We propose a new method to show the effectiveness of reused winning tickets (WSN, ICML2022) in Video Continual Representations.

## - Backgrounds of Neural Implicit Rep.



- The world around us is not discrete,
- Yet, we choose to represent real-world signals such as images or sound in a discrete manner.

(-) Discrete only contain a discrete amount of information regarding the signal.

(-) Given a 256x256 pixel grid for an image, we are not able to scale it up to a 512x512 image.
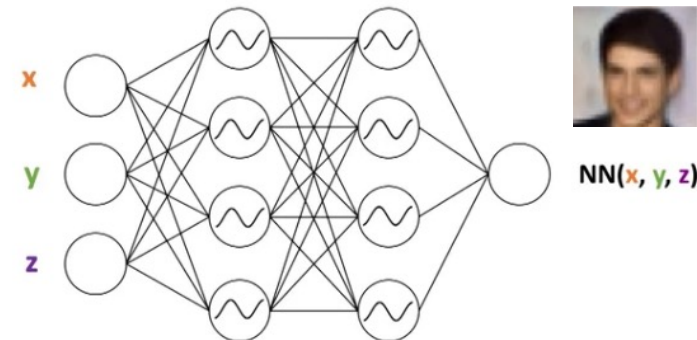
← Not enough information

- **Neural Implicit Representations (NIR)** are neural networks (e.g. MLPs) that estimate the function $f$ that represents a signal continuously, by training on discretely represented samples of the same signal.

- **NIP** learns how to estimate the underlying (continuous) function $f$ (denoted $F$ below):
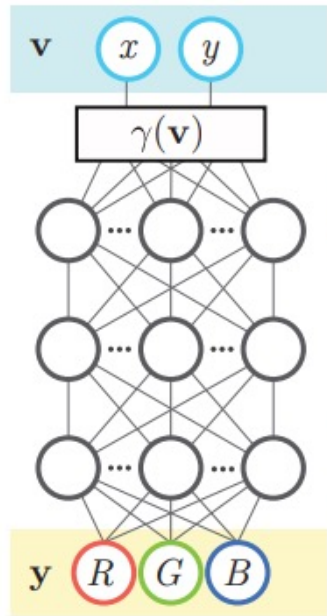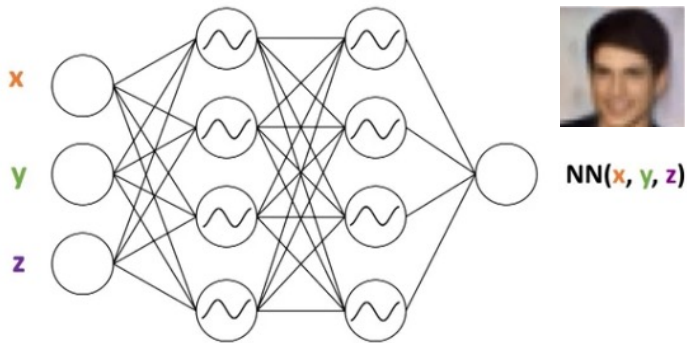
$$F(x, \Phi, \nabla_x \Phi, \nabla_x^2 \Phi, \cdots) = 0, \quad \Phi : x \mapsto \Phi(x)$$

- The network parameterizes $\Phi$. After training on the discretely represented samples, the estimated $f$ would be *implicitly* encoded in the network, hence the name "**Neural Implicit Representation**".



Discrete representations of various signals (SIREN: Sinusoidal Representation Networks, NeuralPS2020)
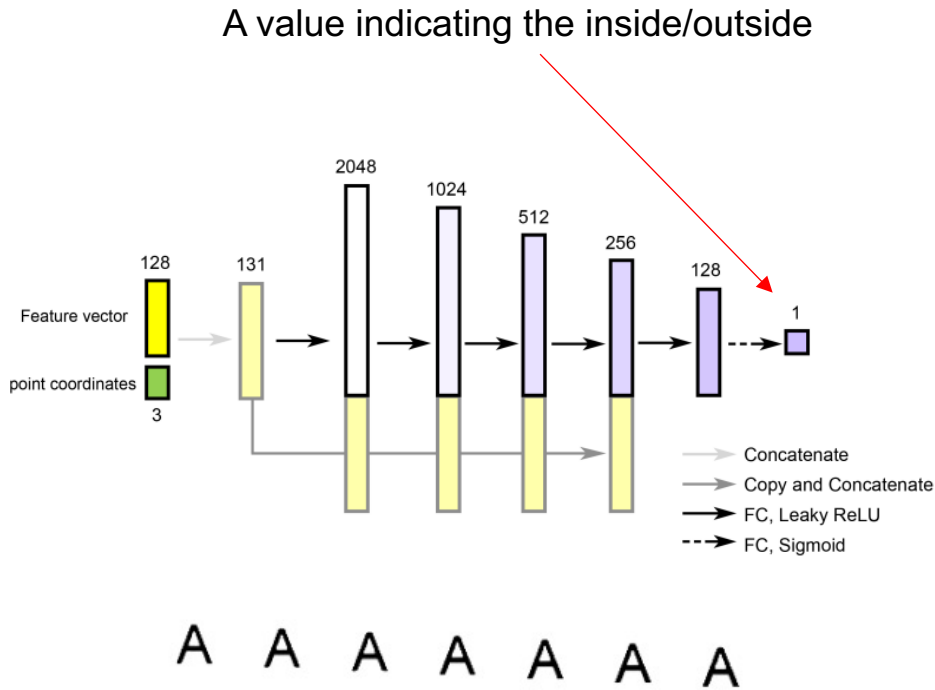
## Coordinate-based Methods



(a) Coordinate-based MLP

(b) Image regression
$(x,y) \rightarrow$ RGB

A value indicating the inside/outside

SIREN: Sinusoidal Representation Networks, NeuralPS2020
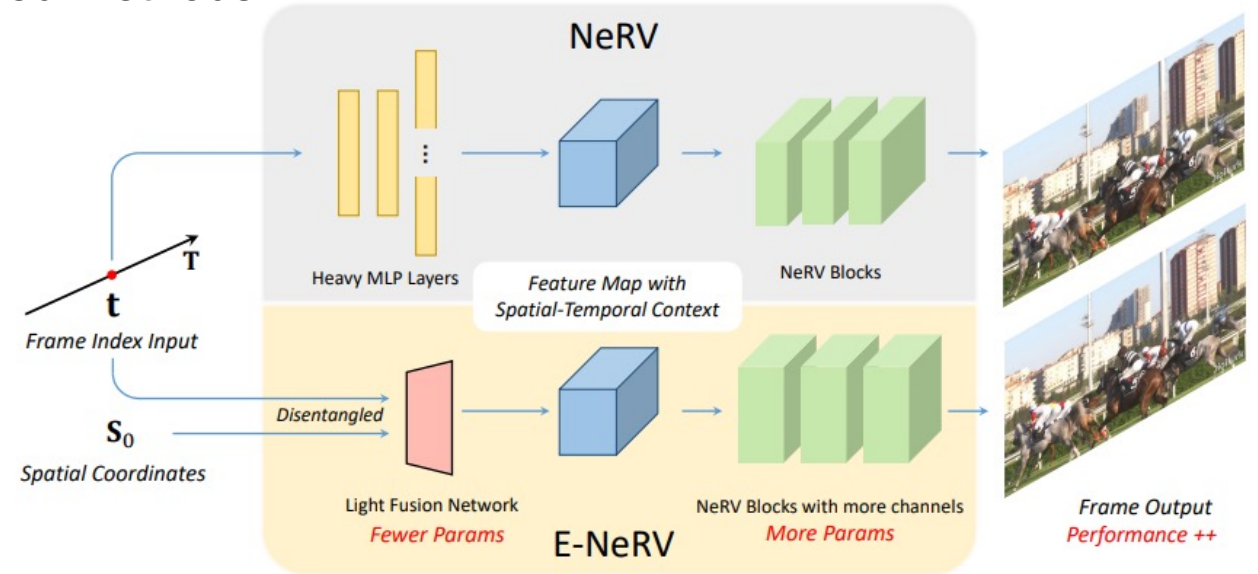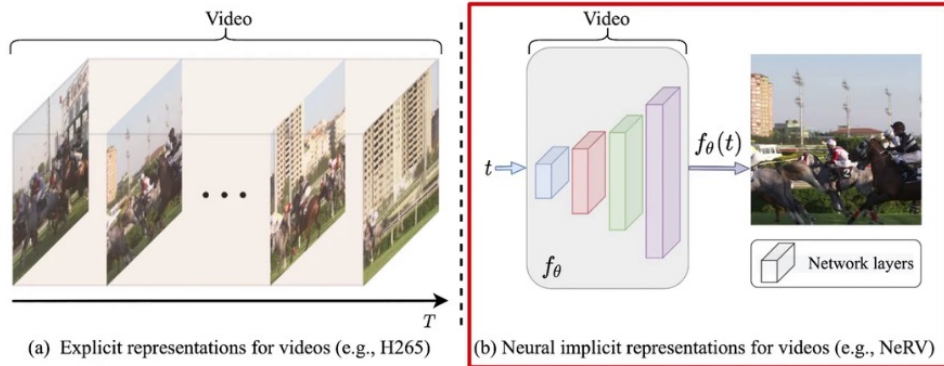
(+) map coordinates to images.

Fourier features, NeuralPS2020

(+) to learn representation with high frequencies

IMNET, CVPR2019

(+) Feature + coordinate led to cleaner interpolation results.

## *Index-based Methods (1)*

(+) Faster and more accurate than Coordinate-based Methods



(a) Explicit representations for videos (e.g. H265)

(b) Neural implicit representations for videos (e.g., NeRV)

NeRV

Heavy MLP Layers — Feature Map with Spatial-Temporal Context — NeRV Blocks

Frame Index Input — t

$S_0$ — Spatial Coordinates — Disentangled

Light Fusion Network — *Fewer Params* — E-NeRV — NeRV Blocks with more channels — *More Params*

Frame Output — *Performance ++*

### NeRV, NeuralPS2021

### E-NeRV, ECCV2022

(+) NeRV represents videos or image datasets as neural neworks, **taking an image index as input**, and **outputting the whole image**.

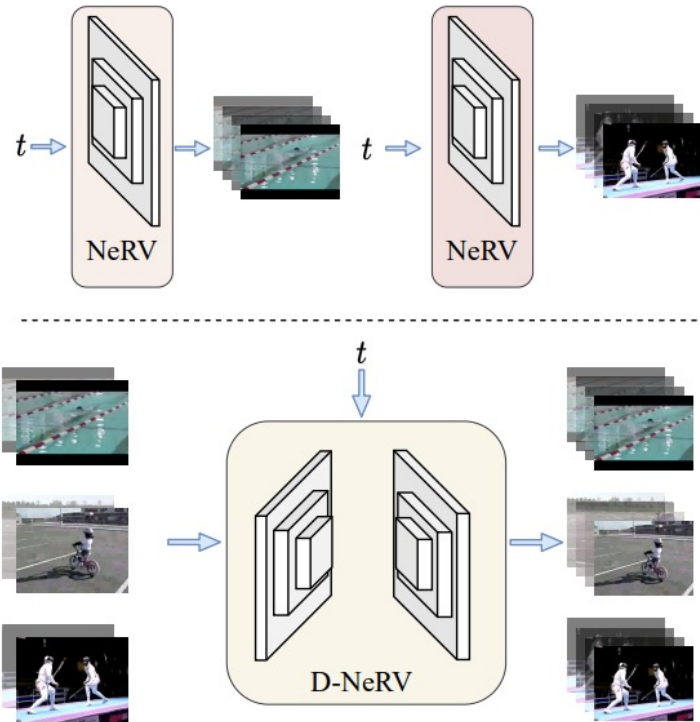(+) **Faster than coordinate-based methods**.

(+) The size of the parameters were reduced **by introducing disentangled spatial-temporal representations** with a light network, while maintaining the majority of performance.
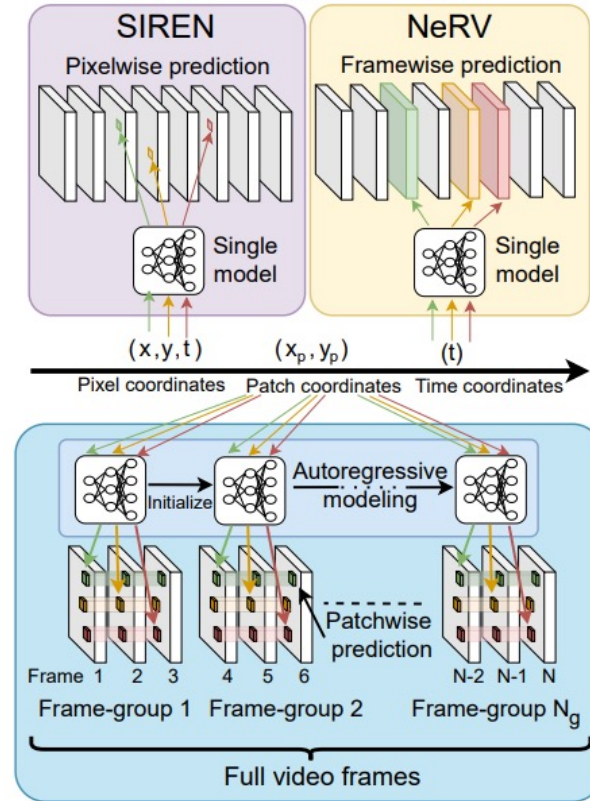
(+) They distributed the saved parameters to increase channel dimensions in convolution blocks, resulting in an E-NeRV model with similar or fewer parameters but much better performance.

## Index-based Methods (2)



D-NeRV, CVPR2023

(+) NeRV optimizes representation to every video independently while **D-NeRV encodes all videos by a shared model**.
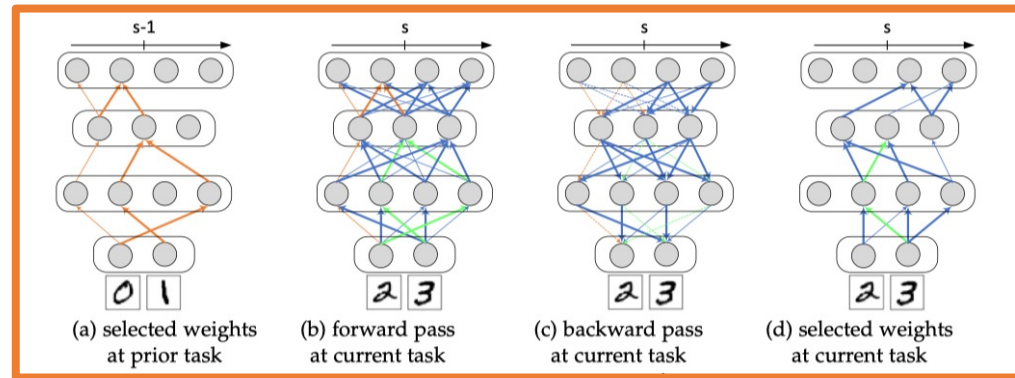


NIRVANA, CVPR2023

(+) **NIRVANA performs spatio-temporal patch-wise prediction** and **fits individual neural networks to groups of frames (clips)** which are initialized using networks trained on the previous group.
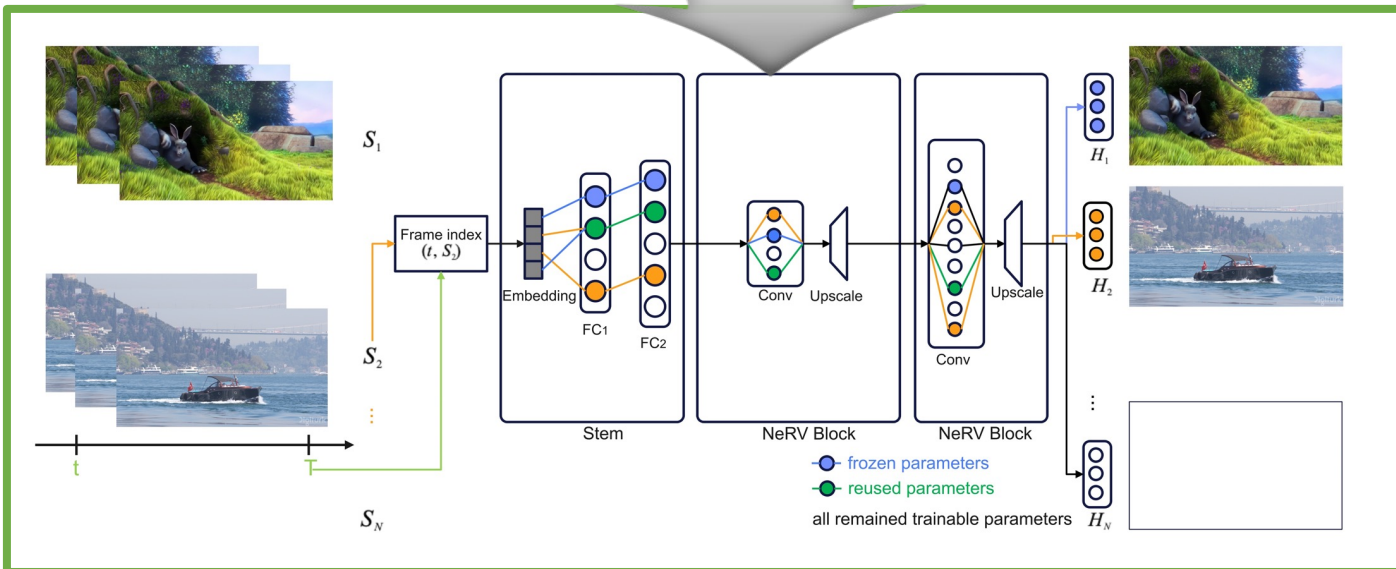
(+) Such an autoregressive patch-wise approach exploits both spatial and temporal redundancies present in videos while promoting scalability and adaptability to varying video content, resolution or duration.

## Index-based sequential Neural Implicit Representation (sequential NIR)

WSN, ICML2022



(a) selected weights at prior task
(b) forward pass at current task
(c) backward pass at current task
(d) selected weights at current task

(+) Providing forget-free continual learning solutions.
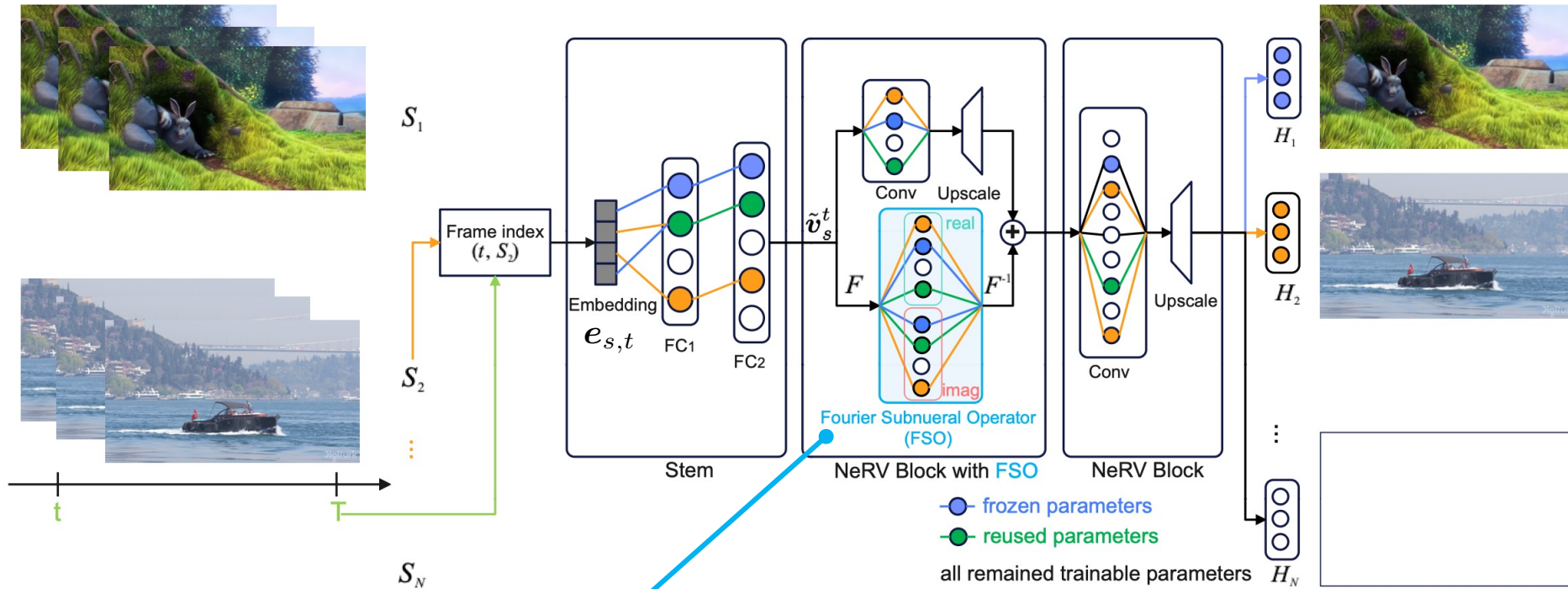(+) Reused Subnetworks Adaptably tailer to tasks.
(+) Fast Convergence .



(+) Providing forget-free video continual learning.
(+) Fast Convergence.

(-) WSN does not have enough parameters for complex video representations.

➔ Needs more parameters to deal with the issue.

## Index-based sequential Neural Implicit Representation (sequential NIR)

(+) Fourier Subneural Operator (FSO): Enough Reusable parameters in Frequency domain



$$\left(\mathcal{K}(\phi)\tilde{\boldsymbol{v}}_t^s\right)(\boldsymbol{e}_{s,t}) = \mathcal{F}^{-1}\left(R_\phi \cdot (\mathcal{F}\tilde{\boldsymbol{v}}_t^s)\right)(\boldsymbol{e}_{s,t}),$$

the Fourier transform of a periodic subnetwork function

$$R_\phi \text{ of } (\boldsymbol{\theta}^{real} \odot \boldsymbol{m}_s^{real}) \text{ and } (\boldsymbol{\theta}^{imag} \odot \boldsymbol{m}_s^{imag})$$

## *Index-based sequential Neural Implicit Representation (sequential NIR)*

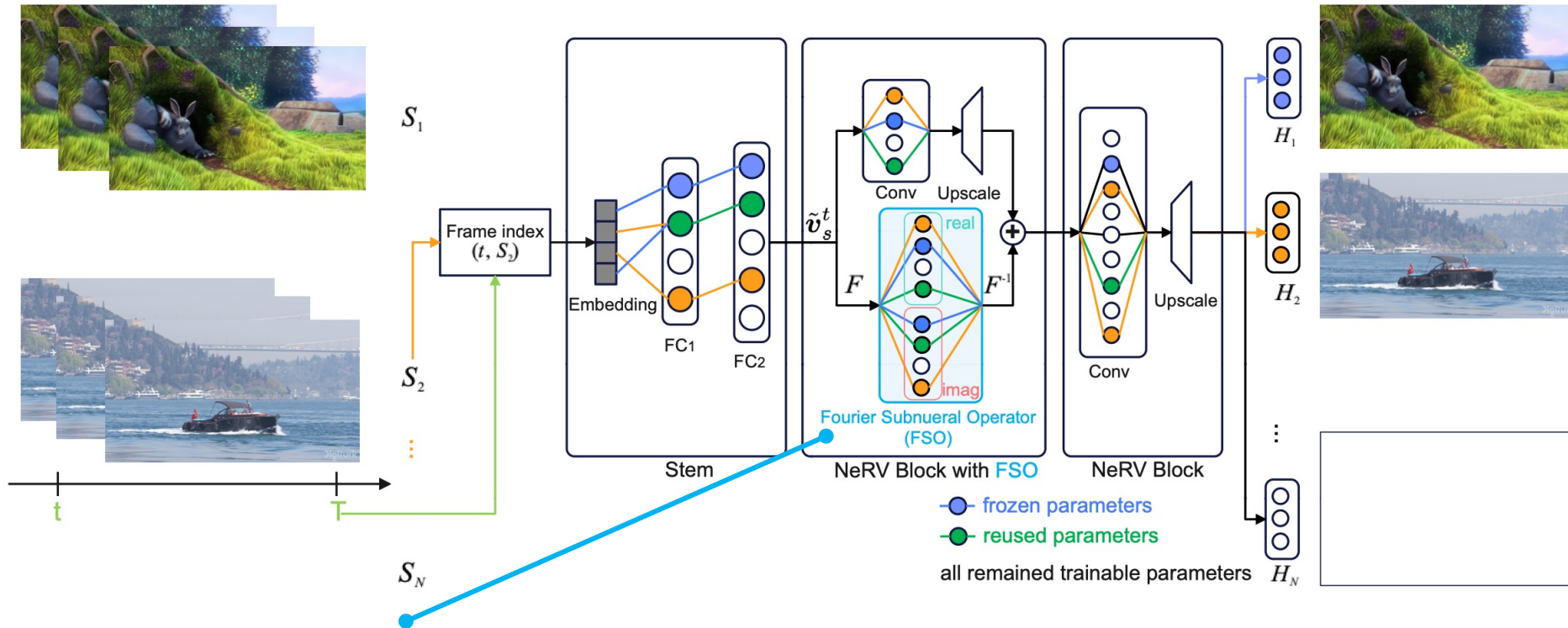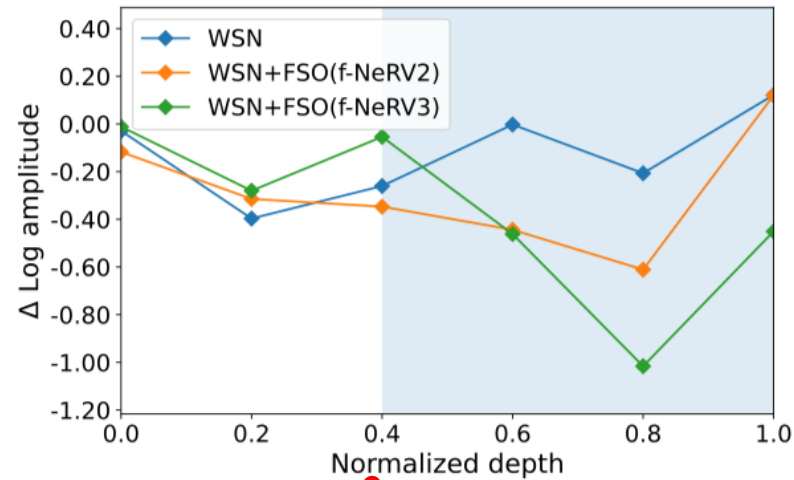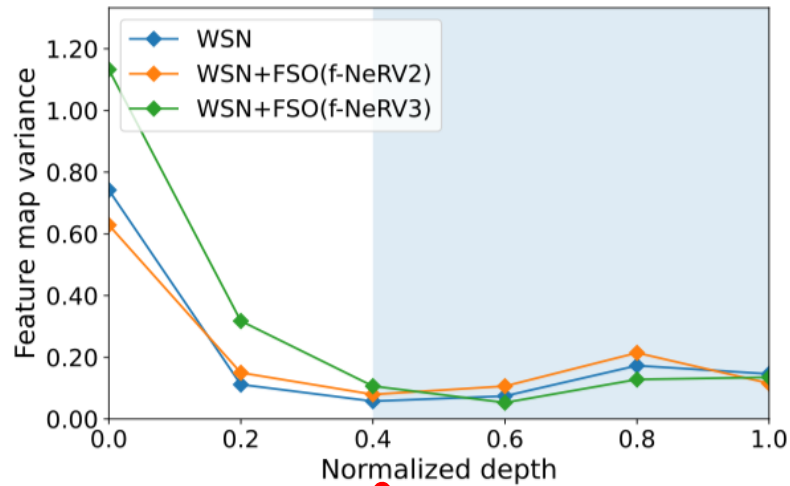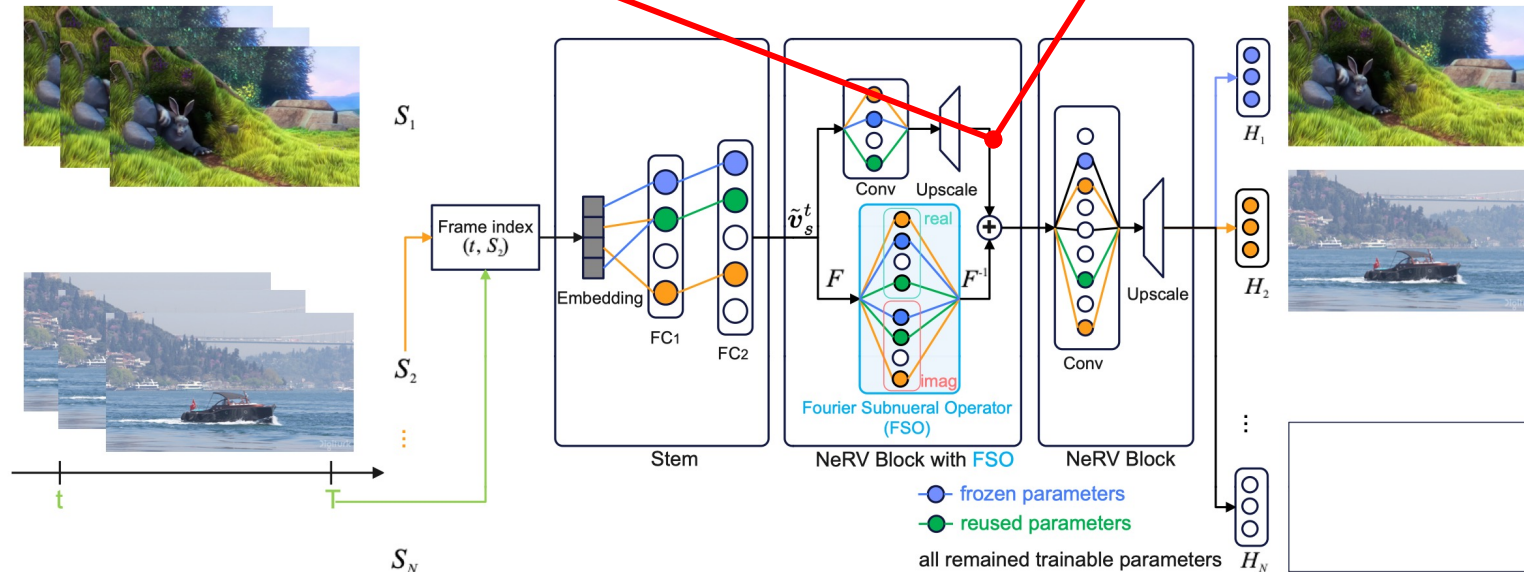(+) Fourier Subneural Operator (FSO): Enough Reusable parameters in Frequency domain



(+) Providing forget-free video continual learning
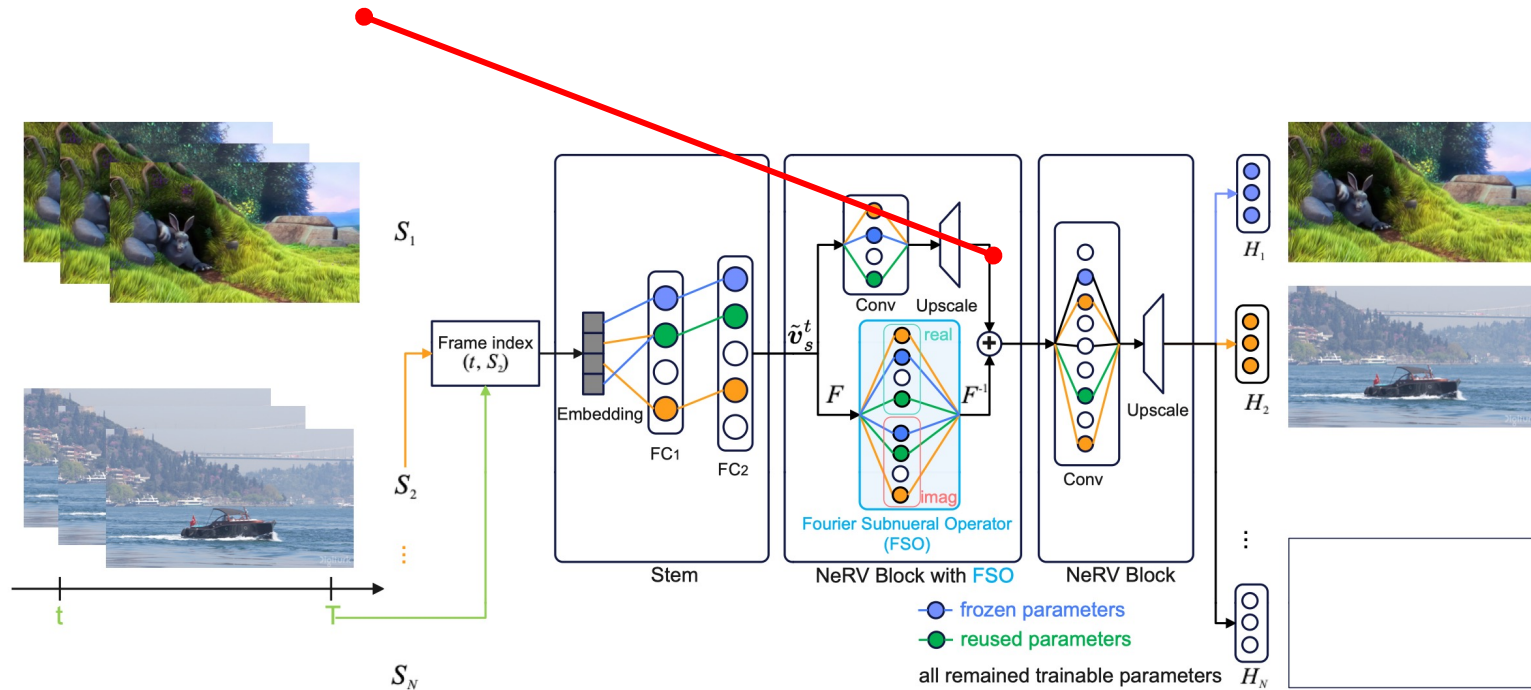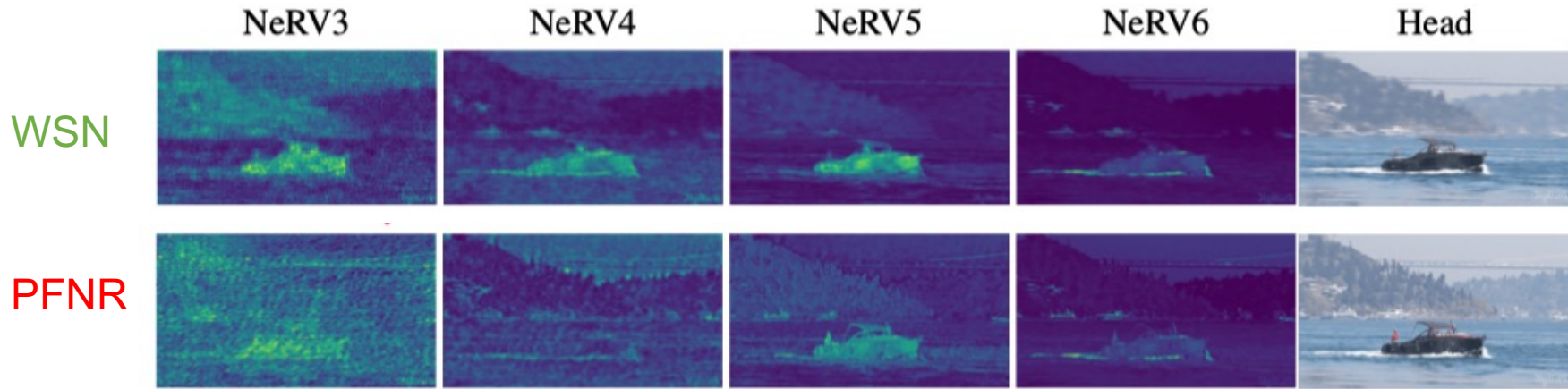(+) Fast Convergence
(+) Enough Reusable parameters in Frequency domain for complex video representation

**Continual Learning: Forget-free Winning Subnetworks for Video Representations – Haeyong Kang, 2024.**

(+) High Feature Variance
(+) High Frequency Components

(+) High Feature Variance
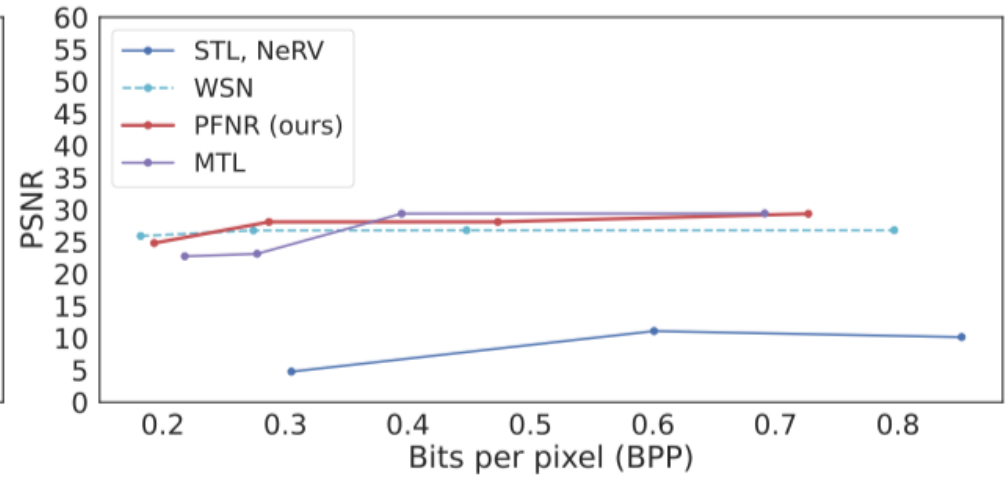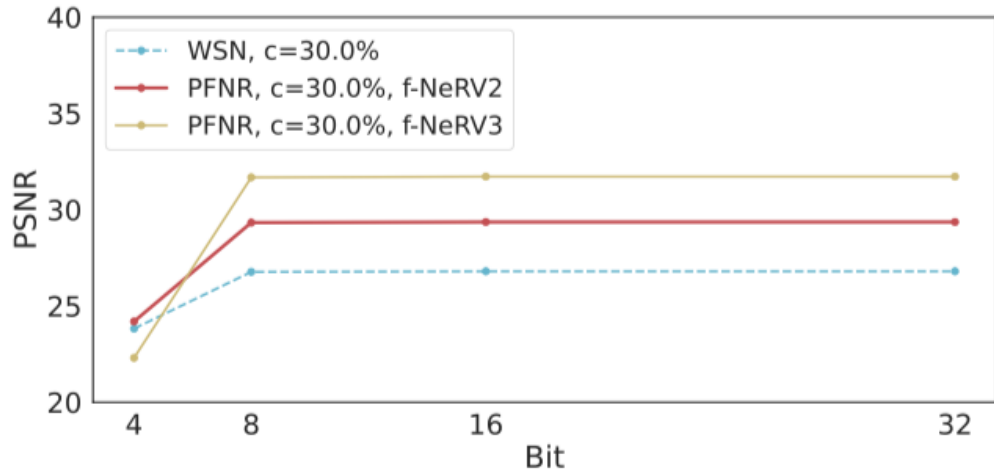(+) High Frequency Components

➔ High Quality Video Represen.

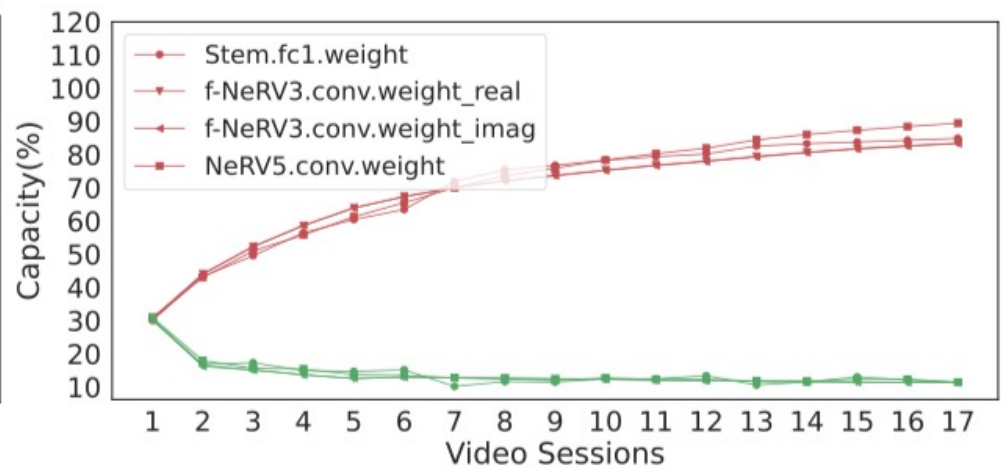**Algorithm 1** Progressive Fourier Neural Representation (PFNR) for VCL

**input:** $\{\mathcal{D}_s\}_{s=1}^N$, model weights of FSO $\boldsymbol{\theta}_* = \{\boldsymbol{\theta}, \boldsymbol{\phi}_{FSO}\}$, score weights of FSO $\boldsymbol{\rho}_* = \{\boldsymbol{\rho}, \boldsymbol{\rho}_{FSO}\}$, binary mask $\mathbf{M}_0 = \{\mathbf{0}^{|\boldsymbol{\theta}|}, \mathbf{0}^{|\boldsymbol{\theta}_{FSO}|}\}$, and layer-wise capacity $c\%$.

1: randomly initialize $\boldsymbol{\theta}_*$ and $\boldsymbol{\rho}_*$.
2: **for** session $s = 1, \cdots, |\mathcal{S}|$ **do**
3:    **if** $s > 1$ **then**
4:       randomly re-initialize $\boldsymbol{\rho}_*$.
5:    **end if**
6:    **for** batch $\mathbf{b}_t \sim \mathcal{D}_s$ **do**
7:       obtain mask $\mathbf{m}_s$ of the top-$c\%$ scores $\boldsymbol{\rho}_*$ at each layer
8:       compute $\mathcal{L}\left(f(\mathbf{e}_{s,t}; \boldsymbol{\theta}_* \odot \mathbf{m}_s), \mathbf{b}_t\right)$, where input embedding, $\mathbf{e}_{s,t} = [\mathbf{e}_s; \mathbf{e}_t]$.
9:       $\boldsymbol{\theta}_* \leftarrow \boldsymbol{\theta}_* - \eta\left(\frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}_*} \odot (\mathbf{1} - \mathbf{M}_{s-1})\right)$           $\triangleright$ trainable weight update
10:      $\boldsymbol{\rho}_* \leftarrow \boldsymbol{\rho}_* - \eta(\frac{\partial \mathcal{L}}{\partial \boldsymbol{\rho}_*})$                       $\triangleright$ weight score update
11:    **end for**
12:    $\hat{\boldsymbol{\theta}}_s = \boldsymbol{\theta}_* \odot \mathbf{m}_s$
13:    $\mathbf{M}_s \leftarrow \mathbf{M}_{s-1} \vee \mathbf{m}_s$                          $\triangleright$ accumulate binary mask
14: **end for**
   **output:** $\{\hat{\boldsymbol{\theta}}_s\}_{s=1}^N$
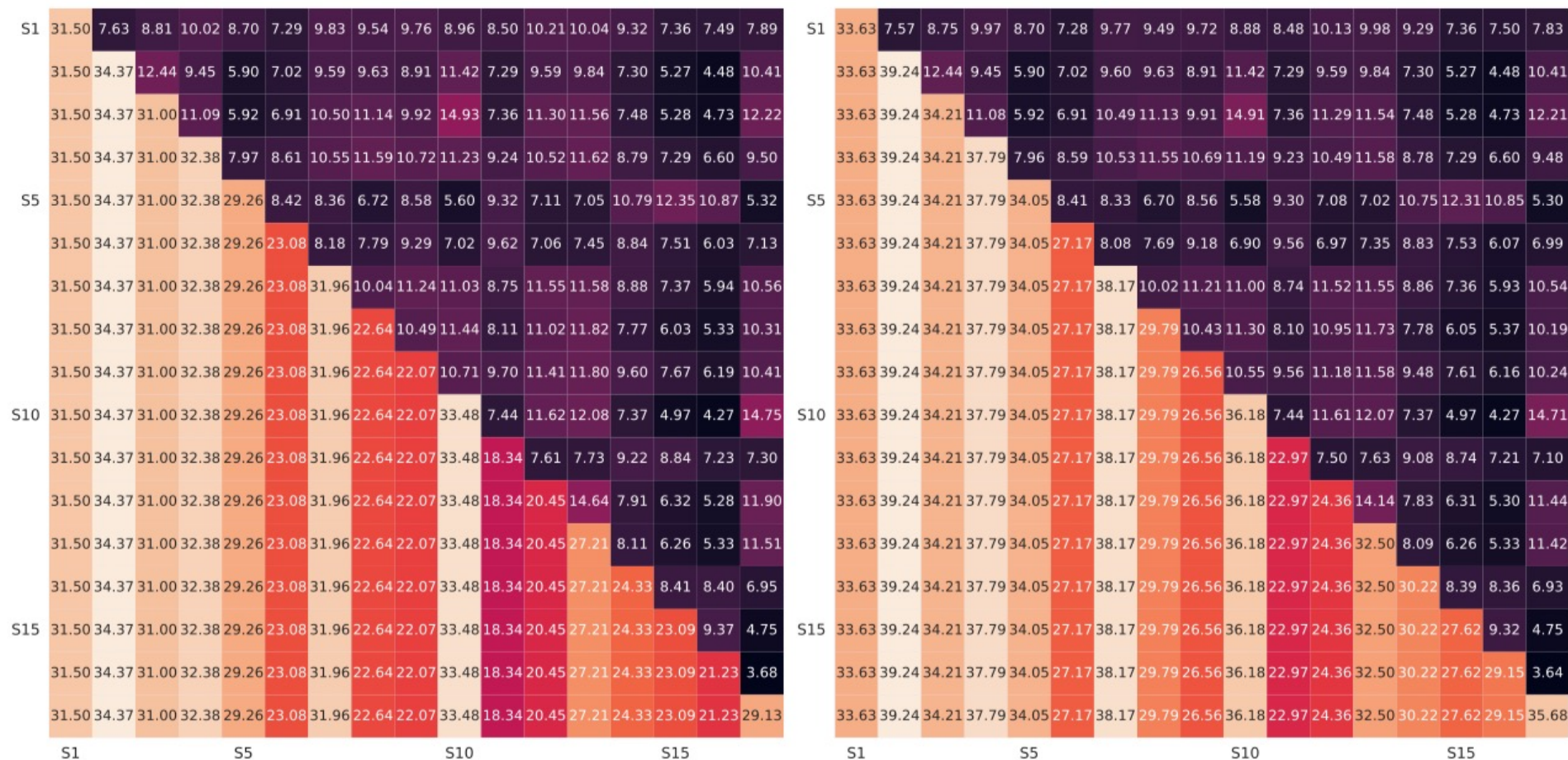
(+) Quantization of PFNR & Bits per pixel (BPP)



(+) PFNR's Performances and Progressive Model Capacity

(a) WSN, $c = 30.0\%$

(b) PFNR, $c = 30.0\%$, $f$-NeRV3

(+) PFNR's Forget-free Video Continual Learning

WSN (29.26, PSNR)

PFNR, *f*-NeRV2 (31.24, PSNR)

PFNR, *f*-NeRV3 (34.05, PSNR)

(+) PFNR's Video Generation

(+) 8-bit PFNR's Video Generation

# Conclusions

- Neural Implicit Representation (NIR) has recently gained significant attention due to its remarkable ability to encode complex and high-dimensional data into representation space and easily reconstruct it through a trainable mapping function.

- However, NIR methods assume a one-to-one mapping between the target data and representation models regardless of data relevancy or similarity. This results in poor generalization over multiple complex data and limits their efficiency and scalability.

- Motivated by continual learning, this work investigates how to accumulate and transfer neural implicit representations for multiple complex video data over sequential encoding sessions. To overcome the limitation of NIR, we propose a novel method, <u>Progressive Fourier Neural Representation (PFNR), that aims to find an adaptive and compact sub-module in Fourier space</u> to encode videos in each training session.

- This sparsified neural encoding allows the neural network <u>to hold free weights, enabling an improved adaptation for future videos</u>. In addition, when learning a representation for a new video, PFNR transfers the representation of previous videos with frozen weights.

- This design allows the model to continuously accumulate high-quality neural representations for multiple videos while ensuring lossless decoding that perfectly preserves the learned representations for previous videos. We <u>validate our PFNR method on the UVG8/17 and DAVIS50 video sequence benchmarks</u> and achieve impressive performance gains over strong continual learning baselines.