

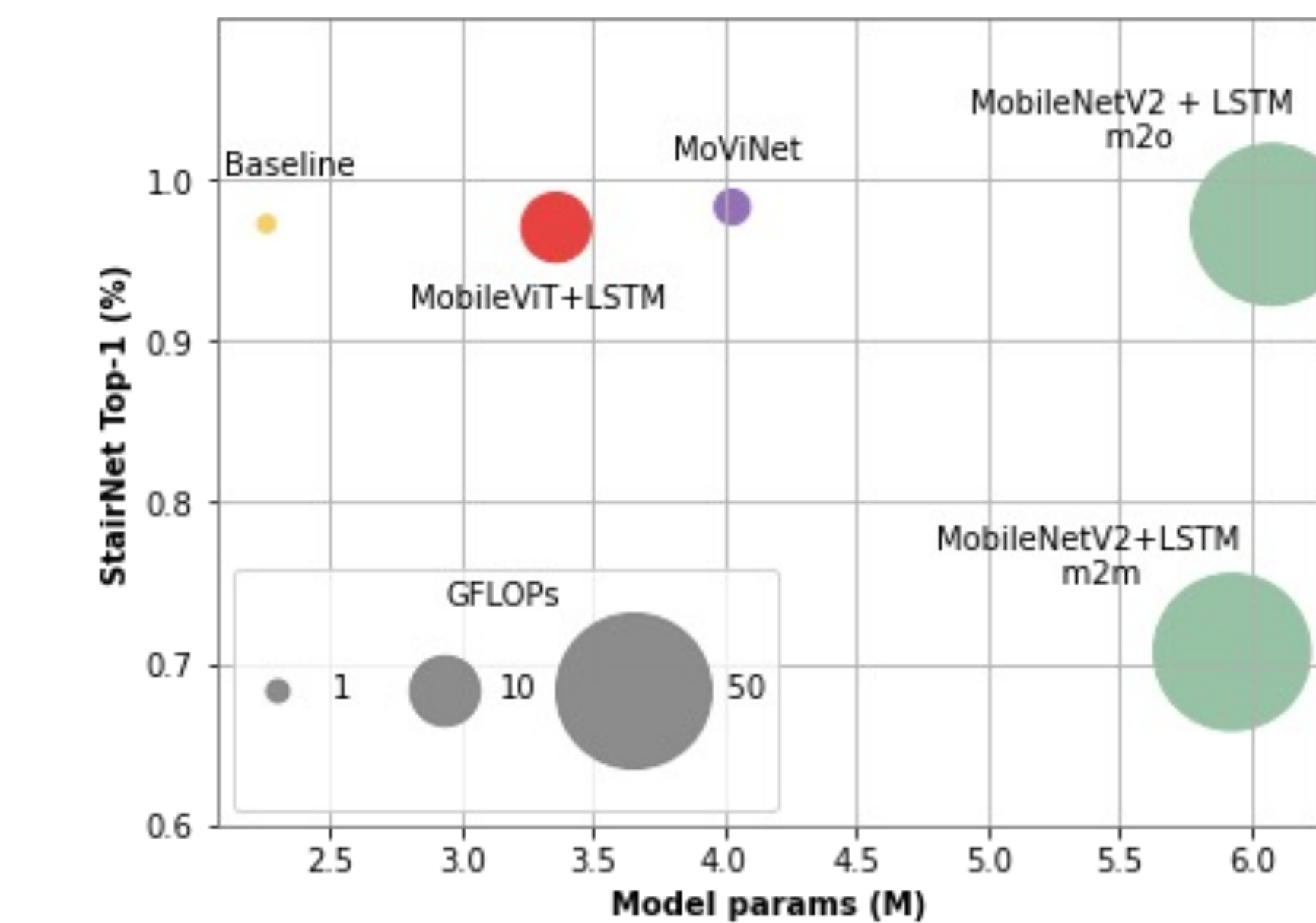
# Sequential Image Classification of Human-Robot Walking Environments using Temporal Neural Networks

Bohdan Ivaniuk-Skulskyi<sup>1,4</sup>, Andrew Garrett Kurbis<sup>2,4</sup>, Alex Mihailidis<sup>3,4</sup>, and Brokoslaw Laschowski<sup>3,4</sup>

<sup>1</sup>National University of Kyiv-Mohyla Academy, Ukraine; <sup>2</sup>University of Waterloo, Canada; <sup>3</sup>University of Toronto, Canada; <sup>4</sup>KITE-Toronto Rehabilitation Institute, Canada

## Introduction

Robotic prosthetic legs and exoskeletons require real-time and accurate state estimation of the walking environment for smooth transitions between different locomotion mode controllers. However, previous studies have mainly been limited to static image classification, therein ignoring the temporal dynamics of human-robot walking.



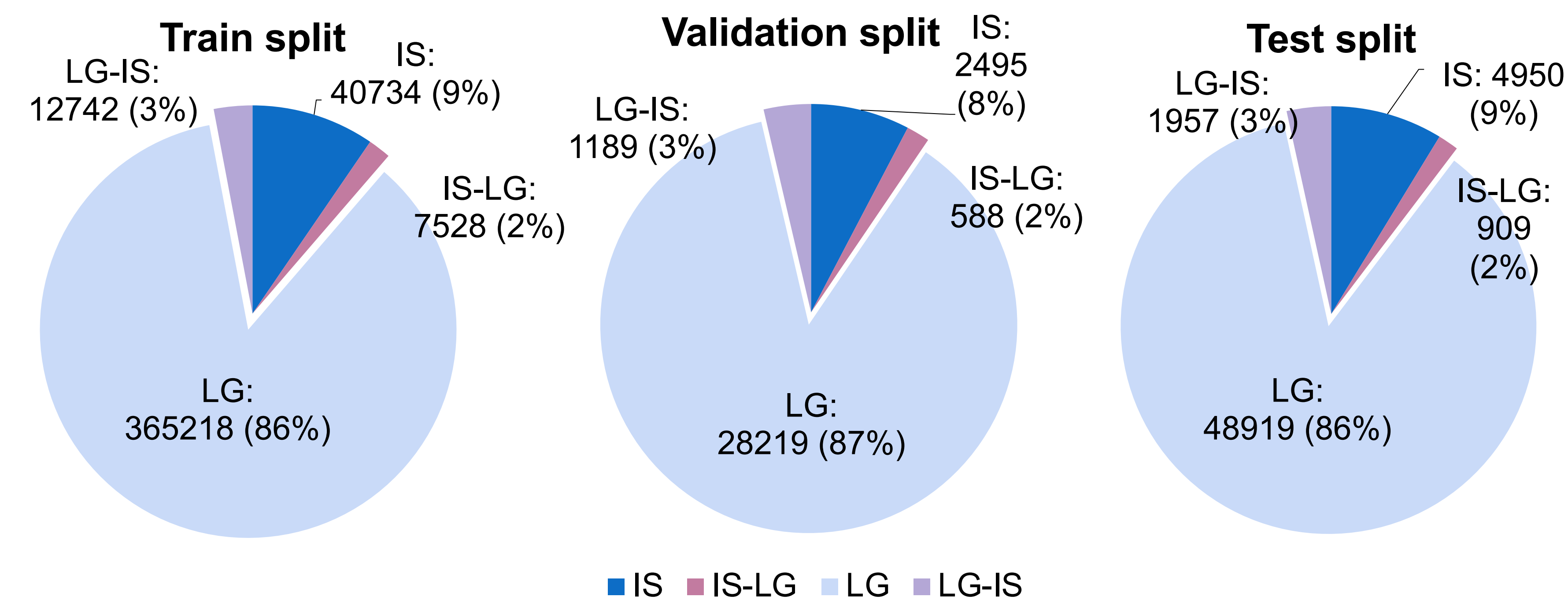
**Figure 1.** Comparison of model size, number of operations, and prediction accuracy. The 3D-CNN MoViNet model achieved better accuracy with slightly more FLOPs and parameters than the state-of-the-art 2D-CNN baseline.

## Objectives

Here we developed different temporal neural networks to compare the performances of static vs. sequential image classification of real-world walking environments in terms of prediction accuracy, and computational and memory storage requirements [1].

## Methods

Building of *ExoNet* and *StairNet*, the large-scale datasets of first-person videos of real-world walking environments, we organized the data into sequences of images for four classes, including level-ground terrain, incline stairs, and transitions to and from incline stairs.



Sequential data requires specialized deep learning models that account for temporal dynamics. We studied different encoders, including VGG, EfficientNet, MobileNetV2, ViT, and MobileViT, each coupled with a temporal long short-term memory (LSTM) backbone. We also built and studied MoViNet - a new video classification model designed for mobile and embedded devices with limited computational resources [2].

	IS	IS-LG	LG	LG-IS	IS	IS-LG	LG	LG-IS	IS	IS-LG	LG	LG-IS
IS	94.63	2.75	1.43	1.19	97.9	0.5	1.05	0.5	96.0	0.95	2.1	0.93
IS-LG	11.44	76.9	8.8	2.86	10.4	84.0	5.1	0.45	17.0	74.0	8.8	0.22
LG	0.07	0.32	99.08	0.53	0.035	0.041	99.6	0.34	0.15	0.2	99.07	0.58
LG-IS	6.44	2.4	23.97	67.19	2.7	0.3	24.0	73.0	6.6	0.15	38.1	55.15
	IS	IS-LG	LG	LG-IS	IS	IS-LG	LG	LG-IS	IS	IS-LG	LG	LG-IS
	a) Baseline: MobileNetV2				b) MoViNet				c) MobileViT-XXS + LSTM			
IS	96.0	0.7	2.1	1.2	61.0	16.0	6.9	16.1	58.0	26.6	15.0	0.38
IS-LG	25.0	67.0	7.3	0.66	26.0	47.2	17.0	9.8	19.0	69.0	12.0	0.0
LG	0.13	0.03	99.4	0.44	9.5	16.0	41.0	33.5	0.07	26.9	73.0	0.03
LG-IS	2.6	0.0	37.4	60.0	8.5	11.0	18.5	62.0	0.6	28.7	34.6	36.1
	IS	IS-LG	LG	LG-IS	IS	IS-LG	LG	LG-IS	IS	IS-LG	LG	LG-IS
	d) MobileNetV2 + LSTM				e) MobileNetV2 + LSTM M2M-M2M				f) MobileNetV2 + LSTM M2M-M2O			

**Table 1.** Confusion matrices on StairNet using different temporal neural networks.

## Results

The 3D-CNN MoViNet network outperformed the 2D-CNN encoders with LSTM backbones and the 2D-CNN baseline model in terms of prediction accuracy, suggesting that network architecture plays an important role in performance besides simply the incorporation of temporal data.

Model	Parameters	GFLOPs	Accuracy	F1	NetScore	FPS
MoViNet [2]	4.03	2.5	<b>0.983</b>	<b>0.982</b>	69.67	13.45
MobileViT - LSTM	3.36	9.84	0.970	0.968	64.28	17.55
MobileNet - LSTM	6.08	53.96	0.973	0.970	54.35	17.65
MobileNet - LSTM (seq2seq) - M2O	5.93	50.97	0.707	0.799	49.18	17.55
MobileNet - LSTM (seq2seq) - M2M	5.93	50.97	0.432	0.538	40.62	17.55
Baseline [1]	2.26	0.61	0.972	0.972	<b>78.12</b>	22.58

## Discussion

In this study, we compared static vs. sequential image classification of walking environments. Although 3D-CNN networks achieve higher prediction accuracy compared to 2D-CNN encoders with LSTM backbones, they also have higher computational and memory storage requirements, which can be disadvantageous for robotic prosthetic legs and exoskeletons using edge computing devices

## References

- Kurbis, et al. (2022). "Stair Recognition for Robotic Exoskeleton Control using Computer Vision and Deep Learning". bioRxiv.
- Kondratyuk, et al. (2021). "MoViNets: Mobile Video Networks for Efficient Video Recognition". arXiv.