

Correlation-Aware Joint-Pruning-Quantization using Graph Neural Networks

Muhammad Nor Azzafri Nor-Azman, Usman Ullah Sheikh, Mohammed Sultan Mohammed, Jeevan Sirkunan, Muhammad Nadzir Marsono



UTM
UNIVERSITI TEKNOLOGI MALAYSIA

Abstract

Deep learning excels in image classification but is constrained by its complexity. While joint pruning-quantization offers improvements, it can be further enhanced by considering layer correlations. This **1** exposes redundant computations across layers, **2** facilitates faster convergence in finding optimal pruning-quantization configurations, and **3** achieves better or comparable complexity reduction compared to other works. This paper introduces Graph Neural Networks (GNNs) to aggregate these inter-layer relationships.

Methodology

Algorithm 1 Joint pruning-quantization using GNN and RL

Input: $E_{max}, R_{solve}, N_{max}, C_{target}, I_{update}, K, Q$

Output: $model_{pruned-quantized}$

```

1:  $e \leftarrow 0$ 
2: while  $e \leq E_{max}$  and  $R \leq R_{solve}$  do
3:   Reset  $model_{base}$  and state
4:   Set state through graph representation of  $model_{base}$ 
5:    $n \leftarrow 0$ 
6:   while  $n \leq N_{max}$  and  $C_{current} \leq C_{target}$  do
7:     Set  $C_{current}$  based on agent's action
8:     Prune  $model_{base}$  based on agent's action
9:     Quantize  $model_{pruned}$  based on agent's action
10:    Finetune  $model_{pruned-quantized}$  for  $F$  epochs
11:    Update state through graph representation
    of  $model_{pruned-quantized}$ 
12:    Get reward from  $model_{pruned-quantized}$ 
13:    Increment  $n$ 
14:  end while
15:  if  $e \bmod I_{update} = 0$  then
16:    Train agent for  $K$  epochs
17:  end if
18:  Increment  $e$ 
19: end while
    
```

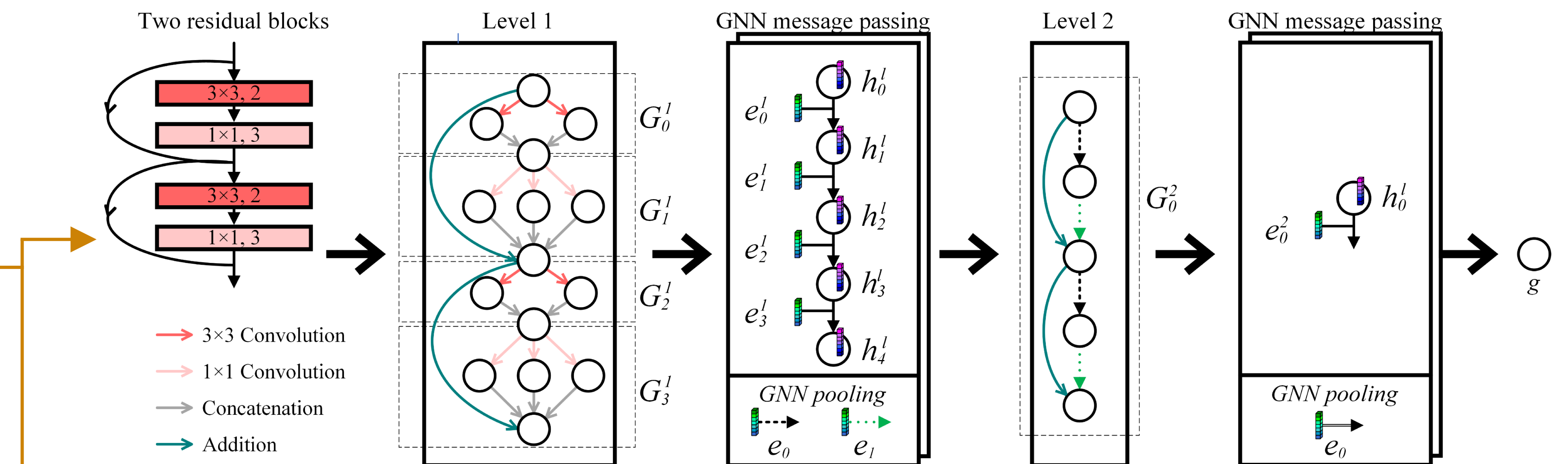


Fig. 1 Example of graph transformation for two residual block.

Overview of underlying process

The GNN is integrated into a reinforcement learning (RL) framework. The baseline or pruned-quantized model is transformed into a graph to serve as the RL **state**, where the agent proposes **actions** to prune and quantize the model, with the resulting accuracy as the **reward**. As episodes progress, the agent, leveraging the GNN to capture layer correlations, learns to propose optimal pruning-quantization actions that maximize accuracy.

Experimental setup

- CIFAR10 (5:1 training-to-validation ratio)
- ResNet20/56
- Batch size 256
- Pytorch's augmentation

Symbol	Representation
b_d	Bitwidth at layer d
C	Complexity reduction
d	Layer of $model_{base}$
e, E	Episode
F	Finetuning epochs
G	Computational graph
h^l	Hidden state at layer l
I	Interval for policy update
K	Epochs for policy update
l	Graph layer
n, N	Timesteps
P_d	Pruning action
Q_d	Quantization action
R	Reward

Results 1

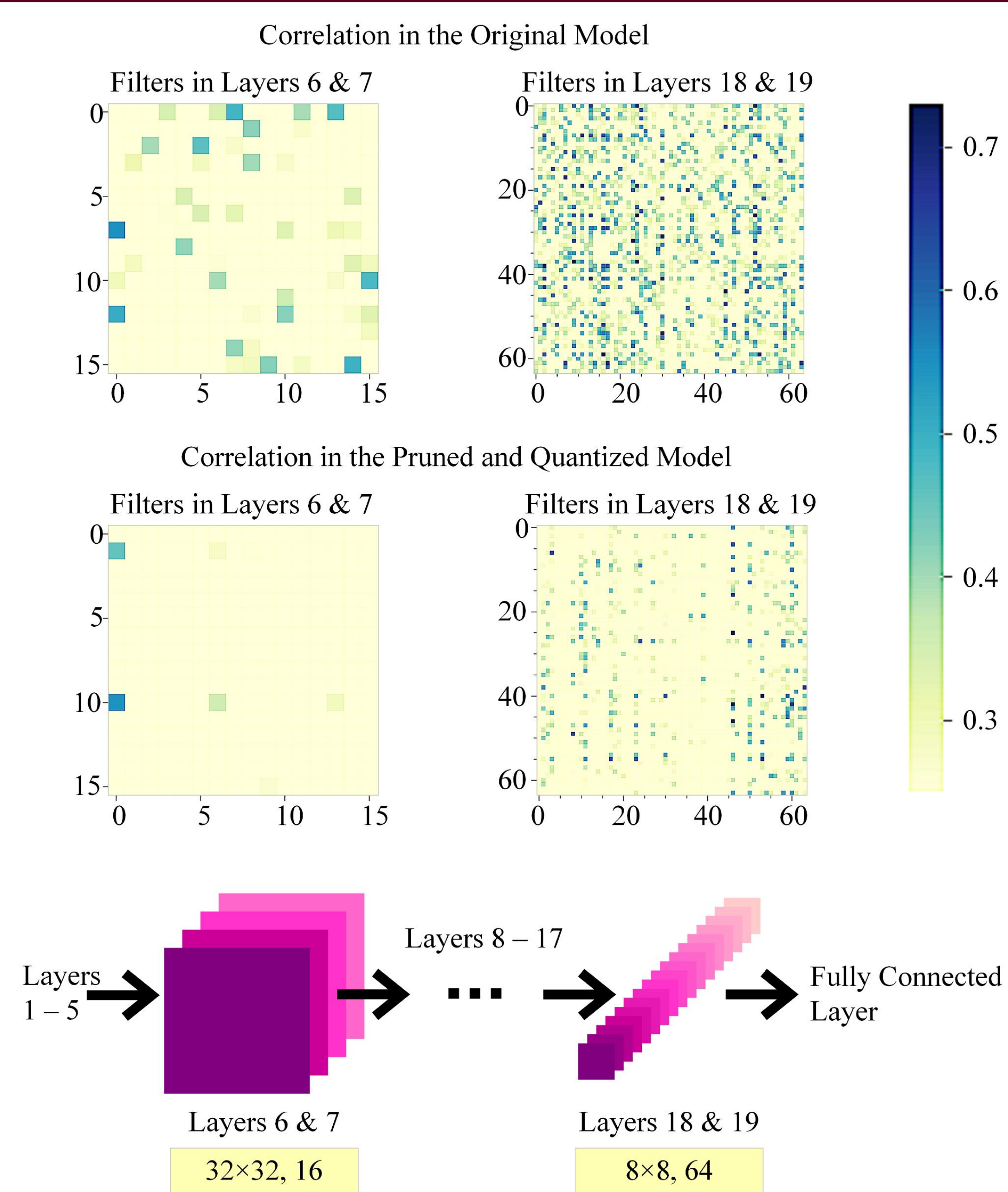


Fig. 2 Pearson correlation matrices

1 Lower correlation in the joint pruning-quantization model suggests that each layer contributes unique computations.

Results 2

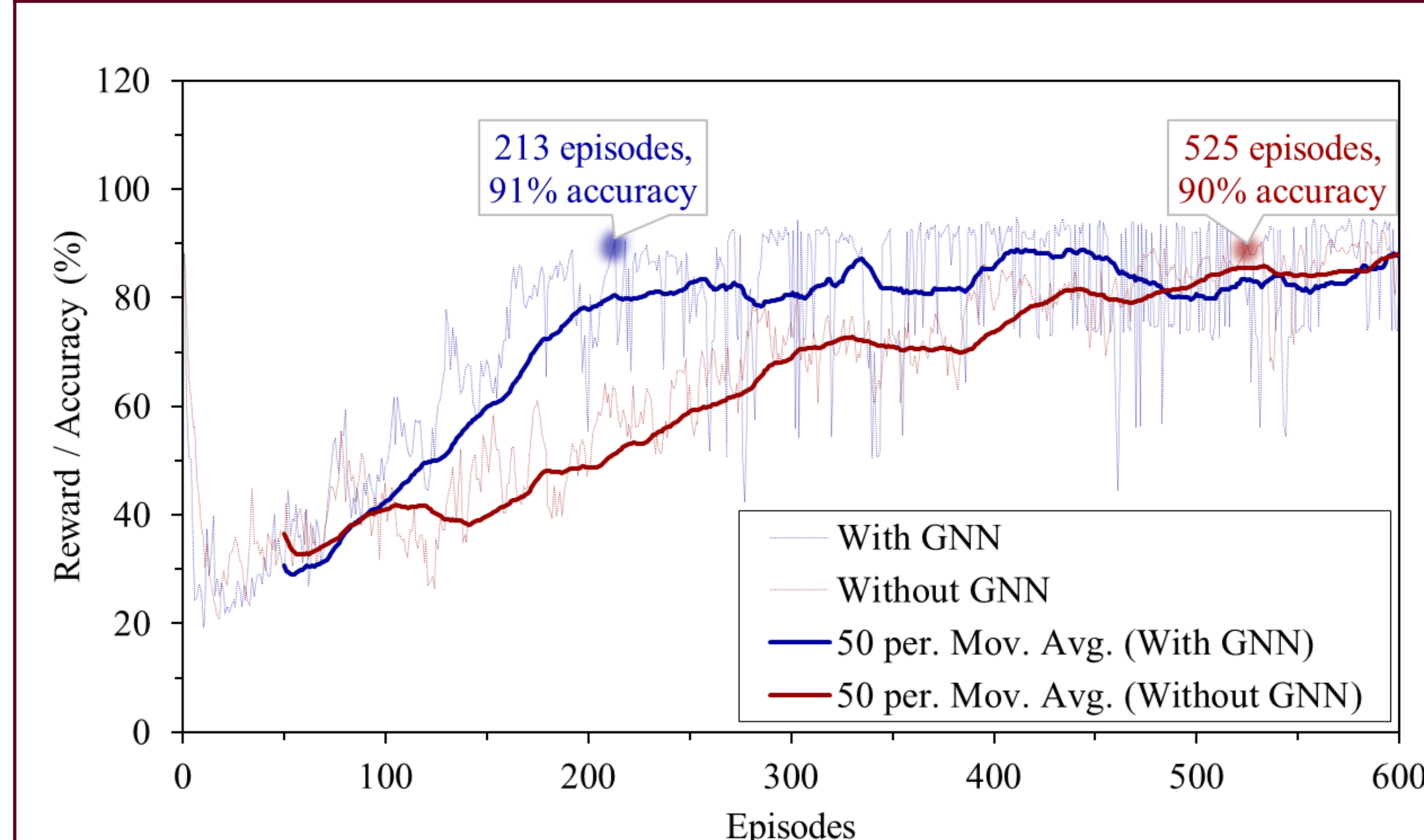


Fig. 3 Convergence rate

2 The ablation study demonstrates that the GNN reduces the number of iterations by an average of 2.46x.

Results 3

DL Model	Approach	Complexity Reduction	Top-1 Accuracy
ResNet20	Baseline	0%	91.73%
	AGMC [1]	50%	91.42%
	GNNRL [2]	49%	91.31%
	DBNN [3]	99.1%	91.60%
	Proposed	96.36%	91.62%
ResNet56	Baseline	0%	93.39%
	AGMC [1]	50%	92.76%
	GNNRL [2]	50%	93.49%
	Proposed	98.51%	92.80%

3 96–98% complexity reduction with 0.1–1.1% accuracy trade-offs.

Future works

- 1** Investigate how other joint pruning-quantization methods impact inter-layer correlations to gain further insights.
- 2** Optimize the depth for the best trade-off, as fewer iterations increase per-iteration time.
- 3** Validate with additional models and datasets.

References

- [1] Sixing Yu, Arya Mazaheri, and Ali Jannesari, "Auto graph encoder-decoder for neural network pruning," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 6362–6372.
- [2] Sixing Yu, Arya Mazaheri, and Ali Jannesari, "Topology-aware network pruning using multi-stage graph embedding and reinforcement learning," in *International Conference on Machine Learning*. PMLR, 2022, pp. 25656–25667.
- [3] Yuhua Lin, Lingfeng Niu, Yang Xiao, and Ruizhi Zhou, "Diluted binary neural network," *Pattern Recognition*, vol. 140, pp. 109556, 2023.