

# Dataset Distillation: A Data-Efficient Learning Framework

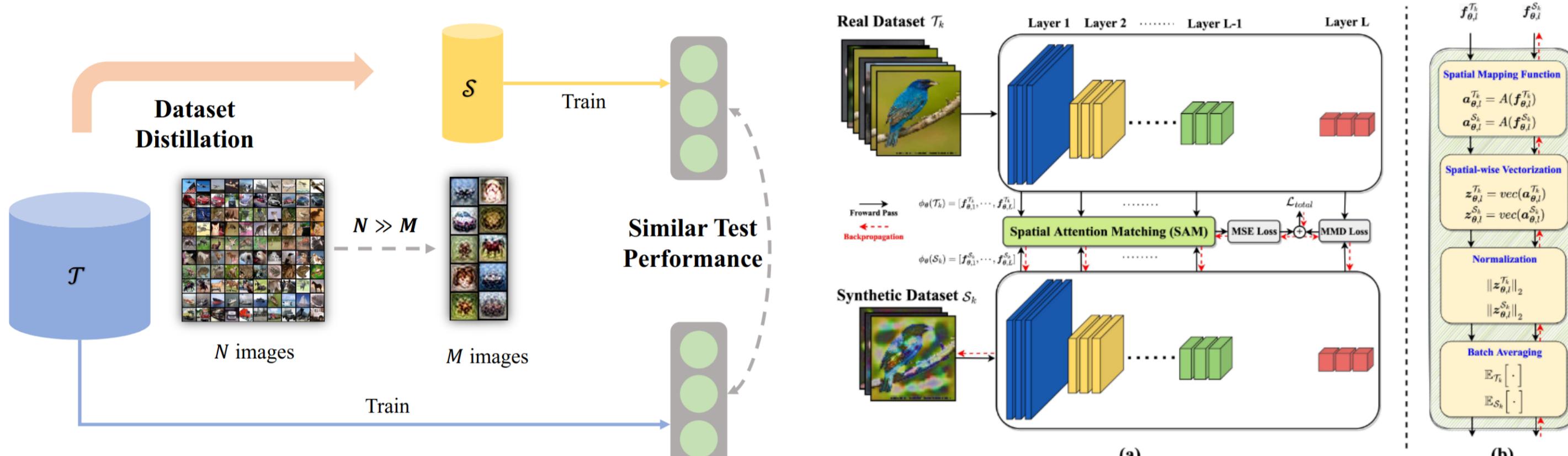
Swapnil Patel

Electrical & Computer Engineering, University of Toronto

## Abstract

This project evaluates dataset distillation using DataDAM, creating compact synthetic datasets. MNIST results show efficiency and representativeness, while MHIST was less promising. DataDAM was also compared with PAD, highlighting trade-offs in accuracy and efficiency.

## Introduction

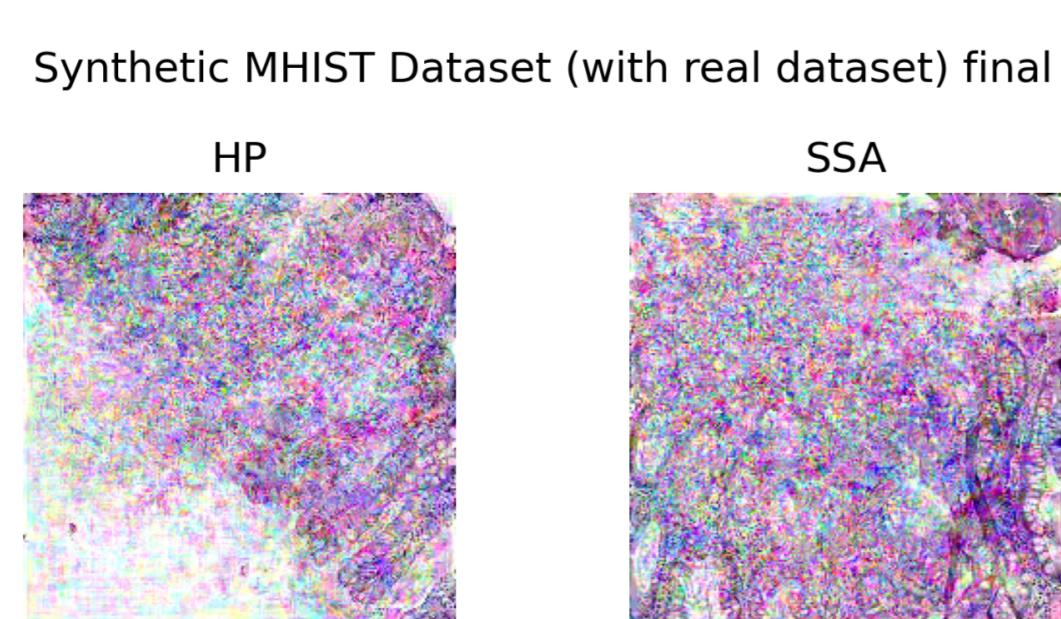
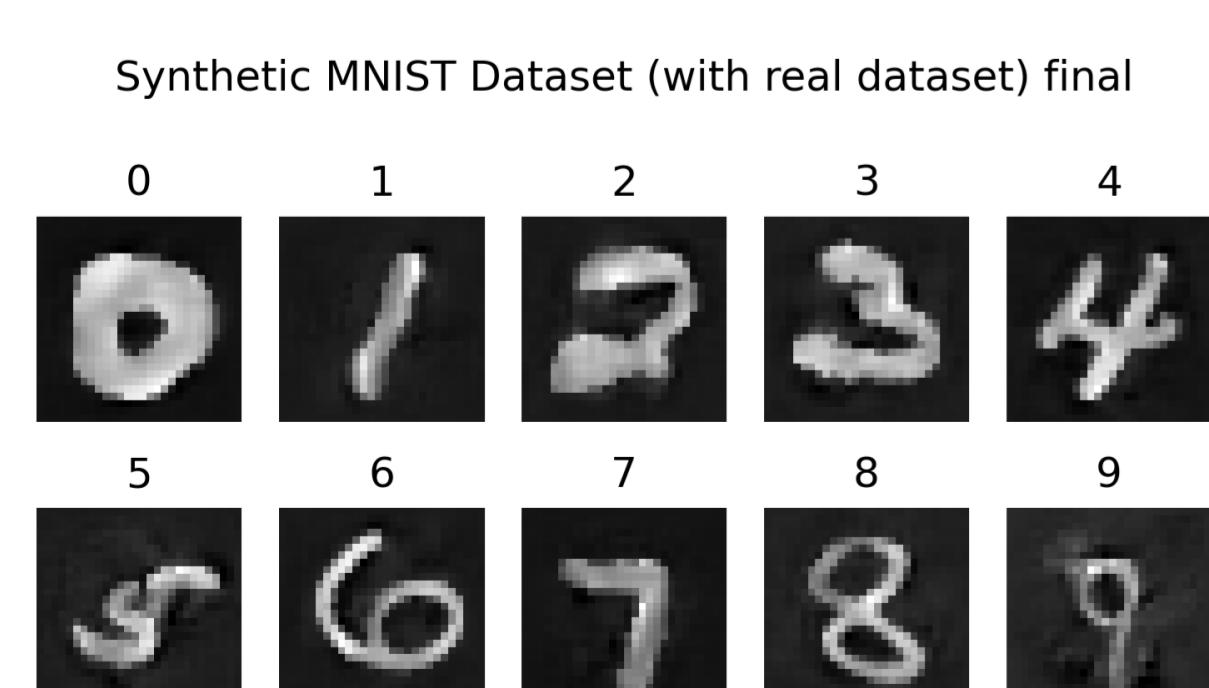
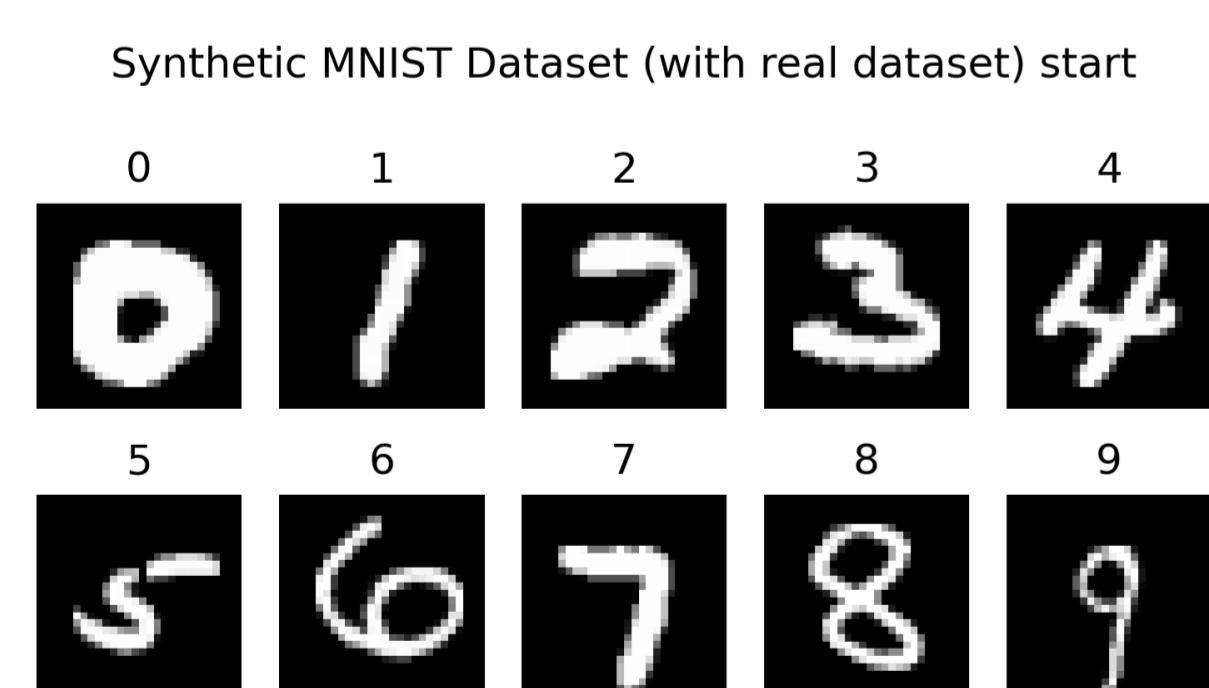


- Definition:** Dataset distillation is a technique to create compact synthetic datasets that retain the critical information of the original data, enabling efficient model training.
- Motivation:** As datasets grow in size, training machine learning models becomes computationally expensive. Dataset distillation addresses this by reducing memory and compute requirements without significantly sacrificing performance.

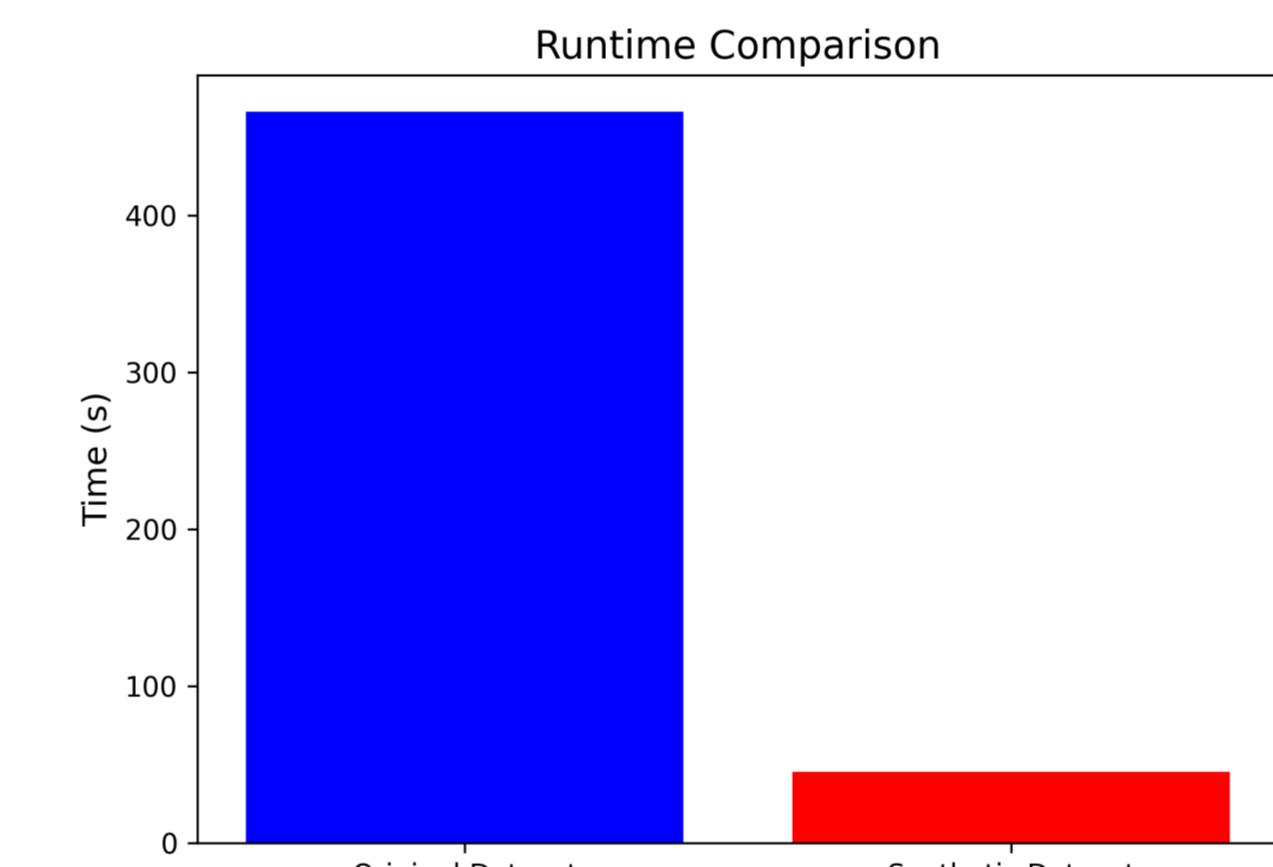
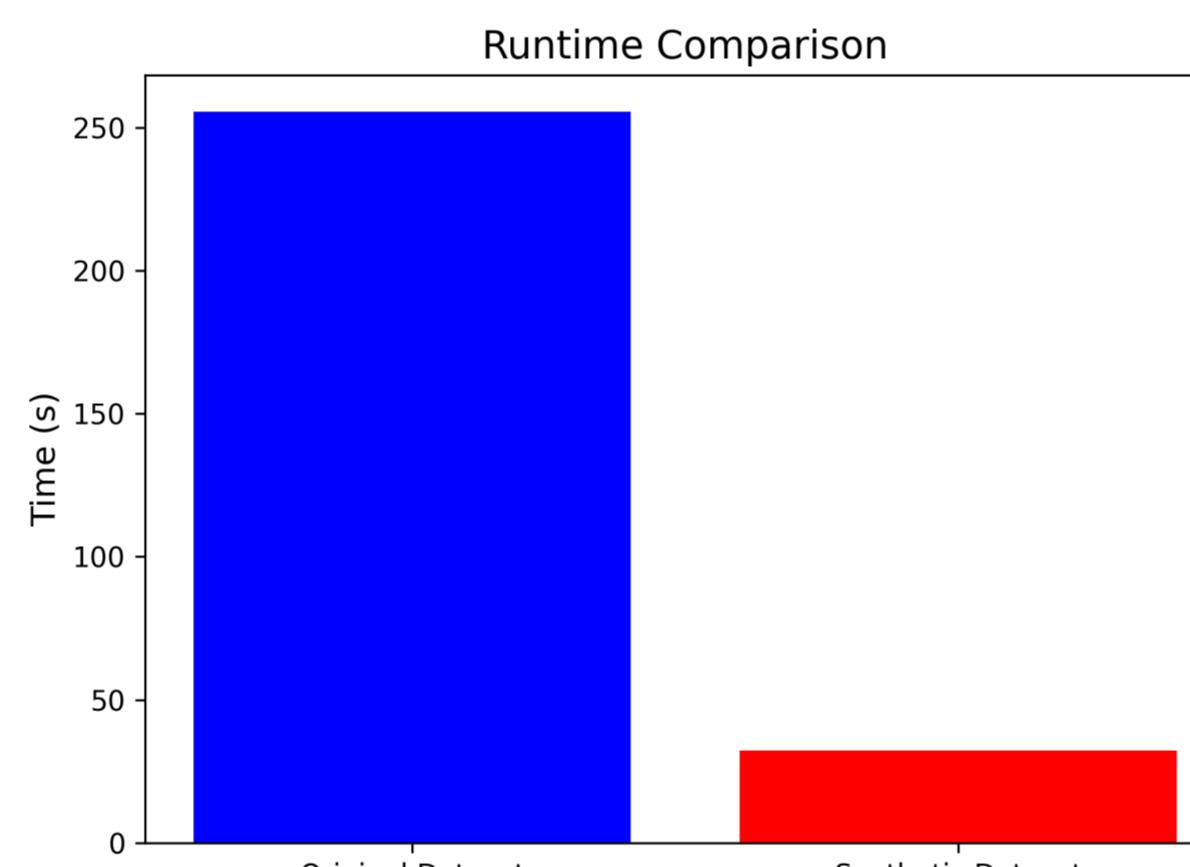
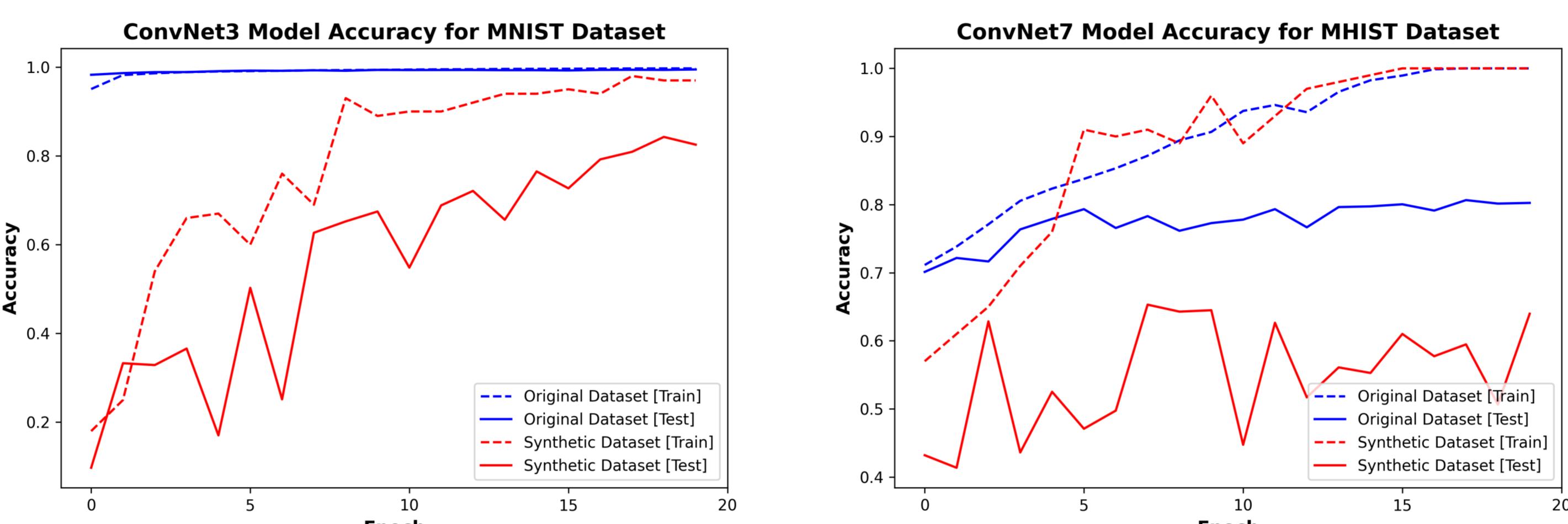
### Dataset Distillation with Attention Matching:

- Attention Matching focuses on aligning attention maps from real and synthetic datasets across layers of a neural network to capture critical features.
- Uses Spatial Attention Mapping (SAM) to extract and align attention distributions.
- Employs randomly initialized networks instead of pre-trained ones, enhancing generalization across architectures.

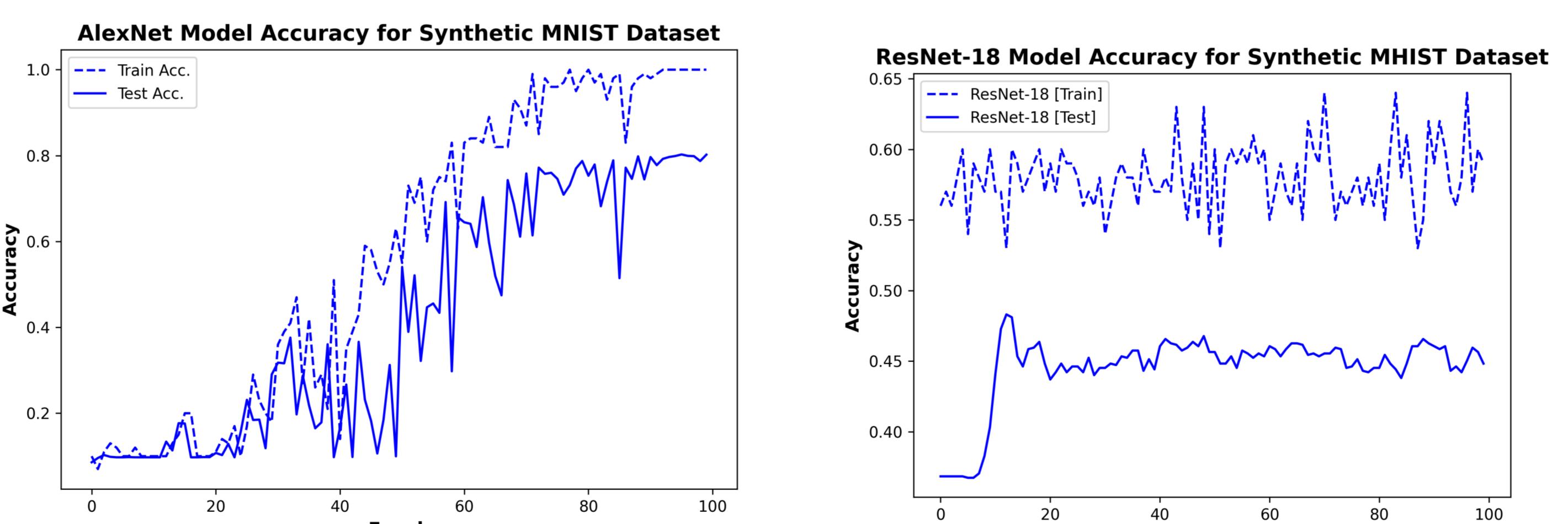
## Distilled Dataset



## Synthetic Dataset Evaluation

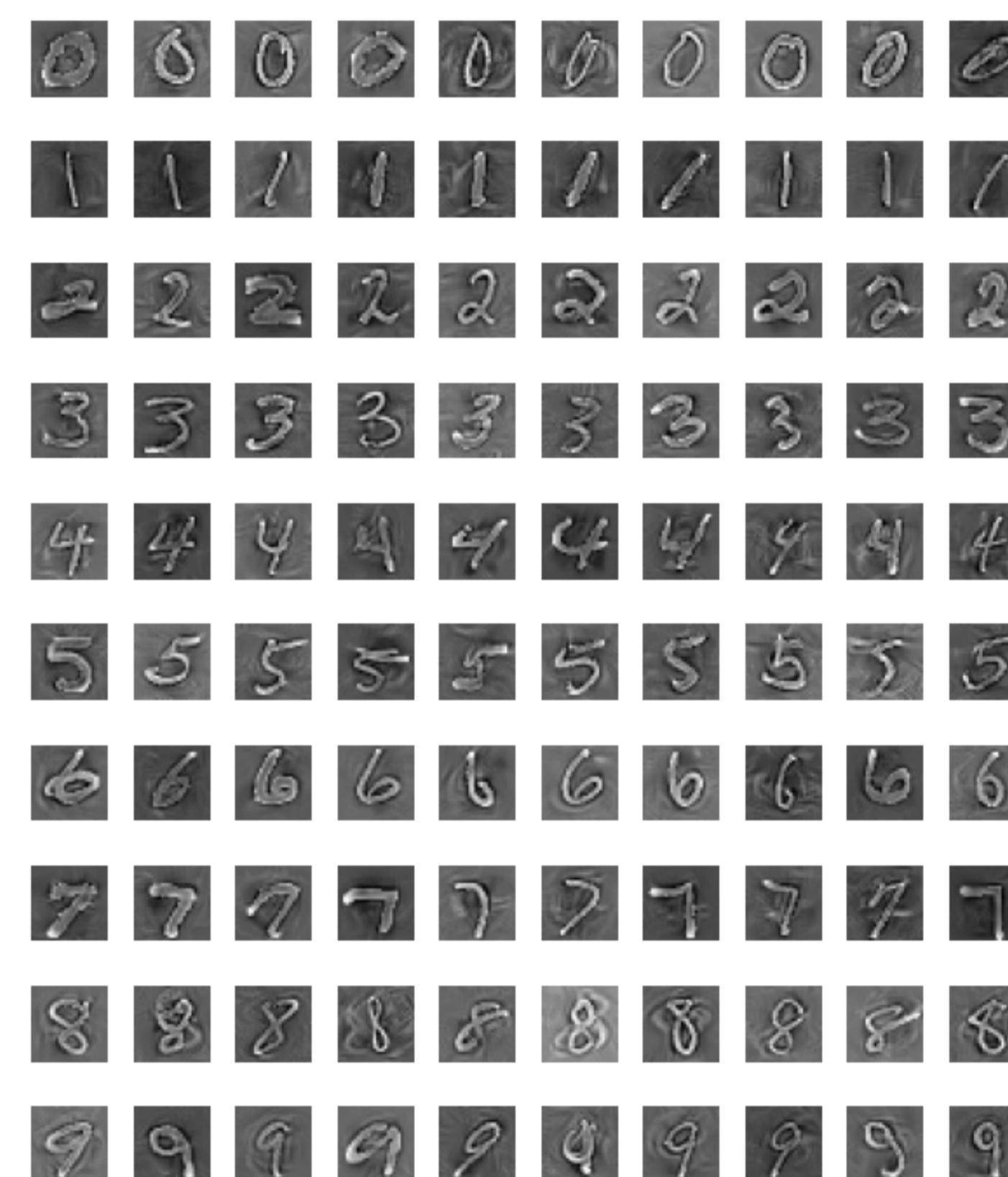


## Cross-Architecture Generalization

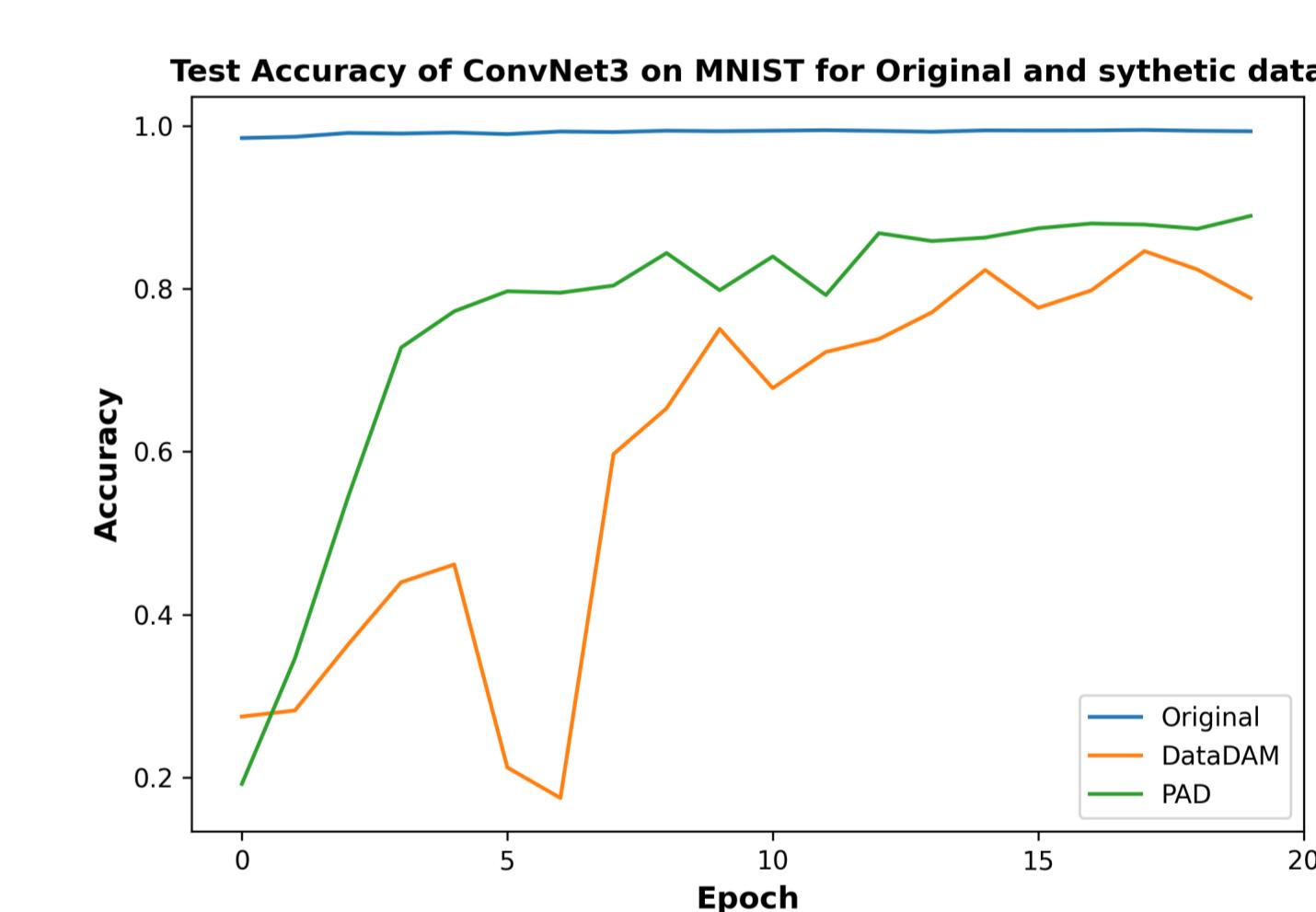
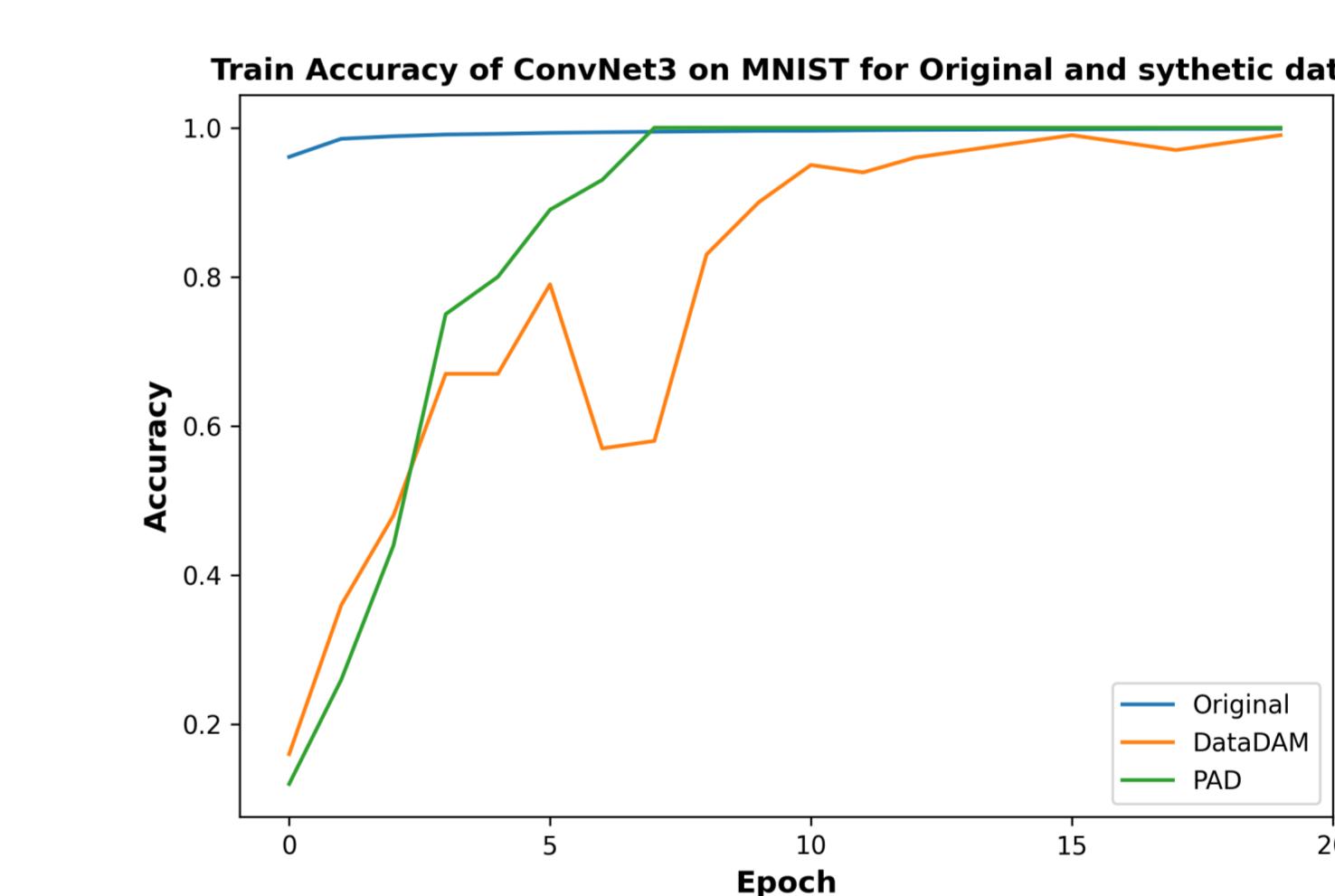


## Prioritize Alignment in Dataset Distillation

Synthetic images Generated using PAD



**Synthetic Dataset Generation:** The figure showcases synthetic images generated with DataDAM, demonstrating the retention of critical dataset features in compact forms. While visually different from the original, these images effectively train models with reduced computational costs.



## Applications

- Neural Architecture Search (NAS)** Provides proxy datasets for rapid exploration of architectures.
- Privacy Preservation** Limits exposure to sensitive information by condensing data into unidentifiable formats.

## References

- [1] Li Deng. The mnist database of handwritten digit images for machine learning research. *IEEE Signal Processing Magazine*, 29(6):141–142, 2012.
- [2] Ahmad Sajedi, Samir Khaki, Ehsan Amjadian, Lucy Z. Liu, Yuri A. Lawryshyn, and Konstantinos N. Plataniotis. Dataadam: Efficient dataset distillation with attention matching, 2023.
- [3] Jerry Wei, Arief Suriawinata, Bing Ren, Xiaoying Liu, Mikhail Lisovsky, Louis Vaickus, Charles Brown, Michael Baker, Naofumi Tomita, Lorenzo Torresani, et al. A petri dish for histopathology image analysis. In *International Conference on Artificial Intelligence in Medicine*, pages 11–24. Springer, 2021.