

# HARNAS: Neural Architecture Search Jointly Optimizing for Hardware Efficiency and Adversarial Robustness of Convolutional and Capsule Networks

Alberto Marchisio<sup>1</sup>, Vojtech Mrazek<sup>2</sup>, Andrea Massa<sup>3</sup>, Beatrice Bussolino<sup>3</sup>, Maurizio Martina<sup>3</sup>, Muhammad Shafique<sup>4</sup>

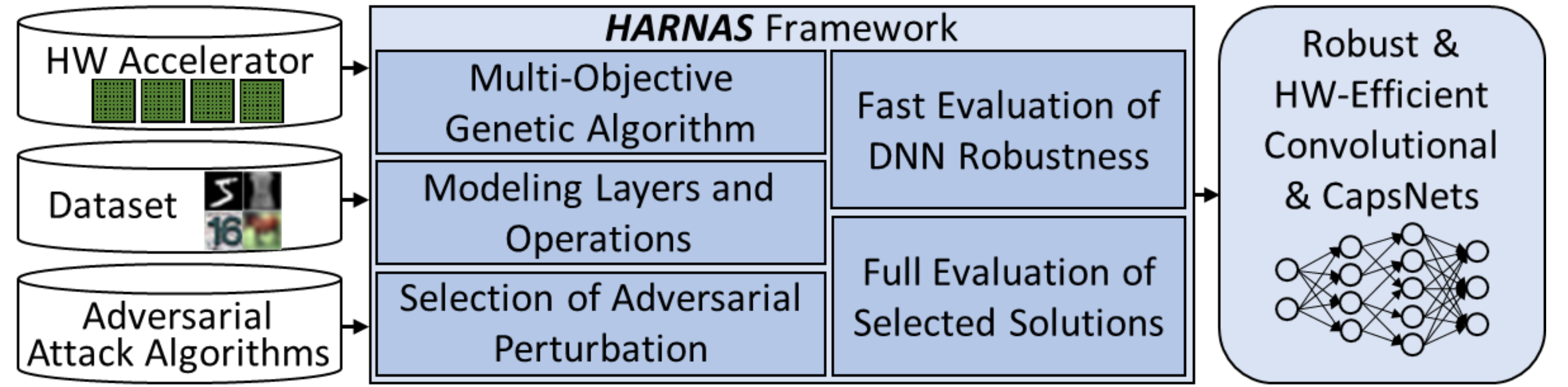
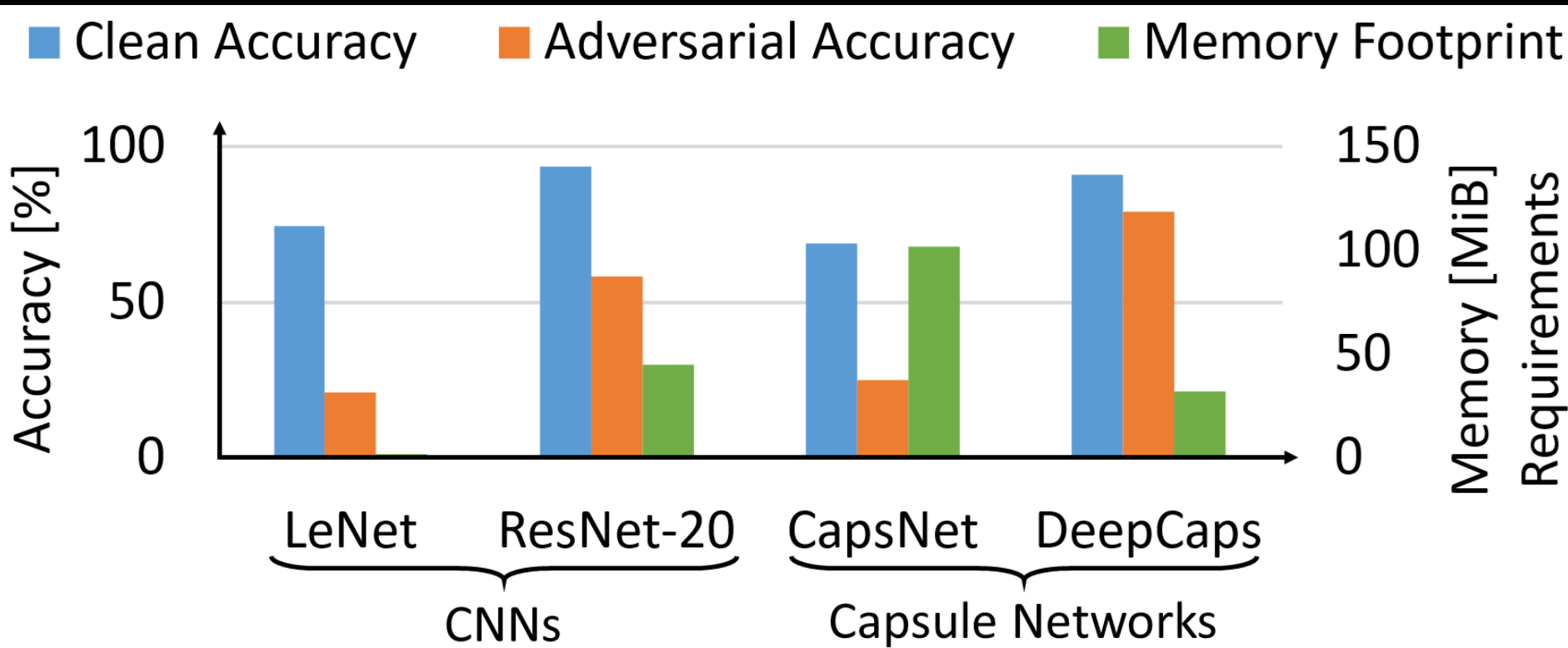
<sup>1</sup>Institute of Computer Engineering, Technische Universität Wien (TU Wien), Vienna, Austria

<sup>2</sup>Faculty of Information Technology, Brno University of Technology, Brno, Czechia

<sup>3</sup>Department of Electronics and Telecommunications, Politecnico di Torino, Turin, Italy

<sup>4</sup>eBrain Lab, Division of Engineering, New York University Abu Dhabi, United Arab Emirates

## Motivations and Novel Contributions



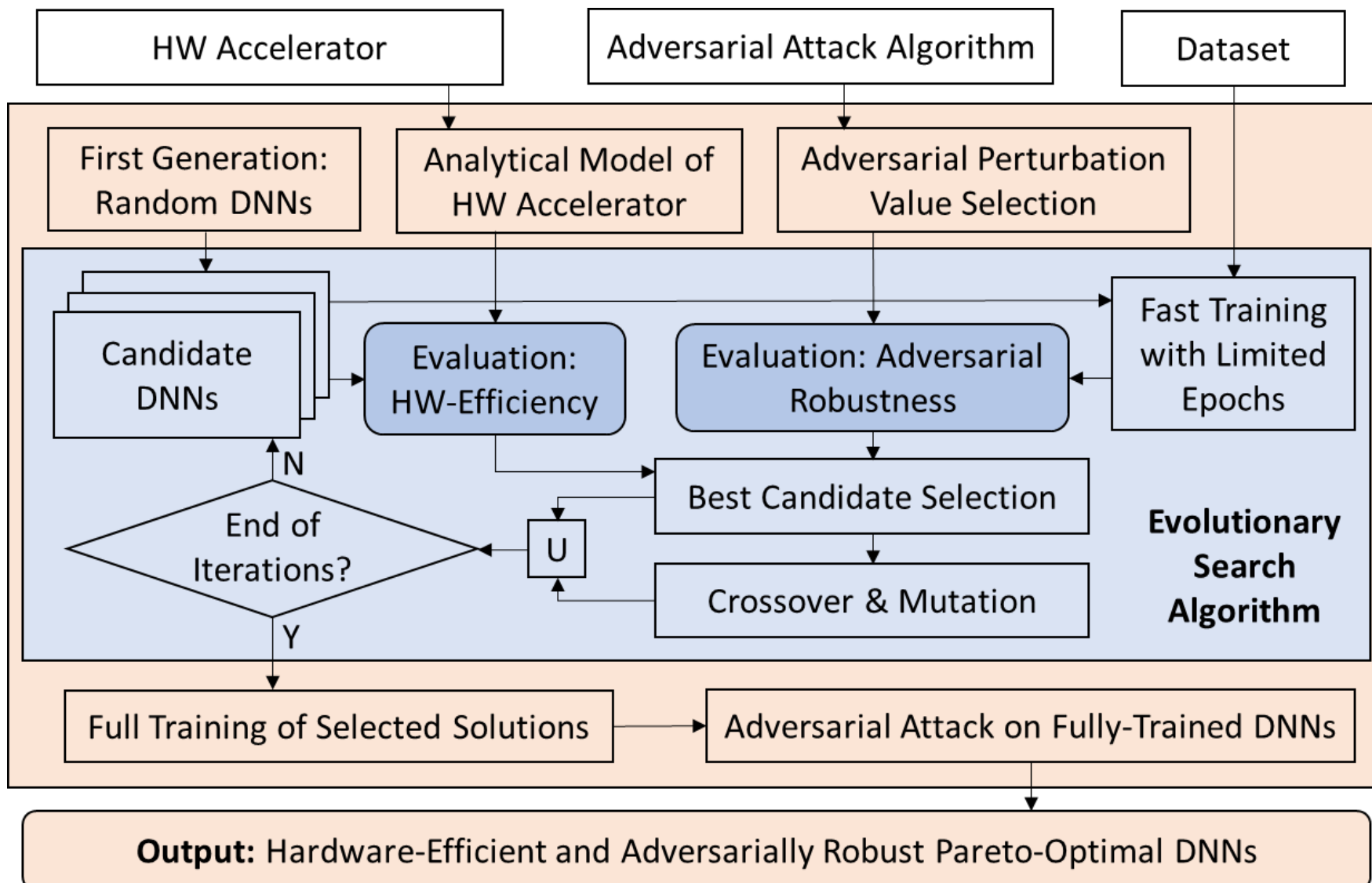
### Motivation Analysis

- The LeNet, which is relatively small and shallow, is hardware efficient due to its low memory footprint, but relatively more vulnerable to attacks.
- A more complex DNN such as the ResNet-20 has a higher memory footprint but it also exhibits higher adversarial accuracy than the LeNet.
- The DeepCaps, despite having a smaller memory footprint than the ResNet-20, is also relatively more robust against adversarial attacks.

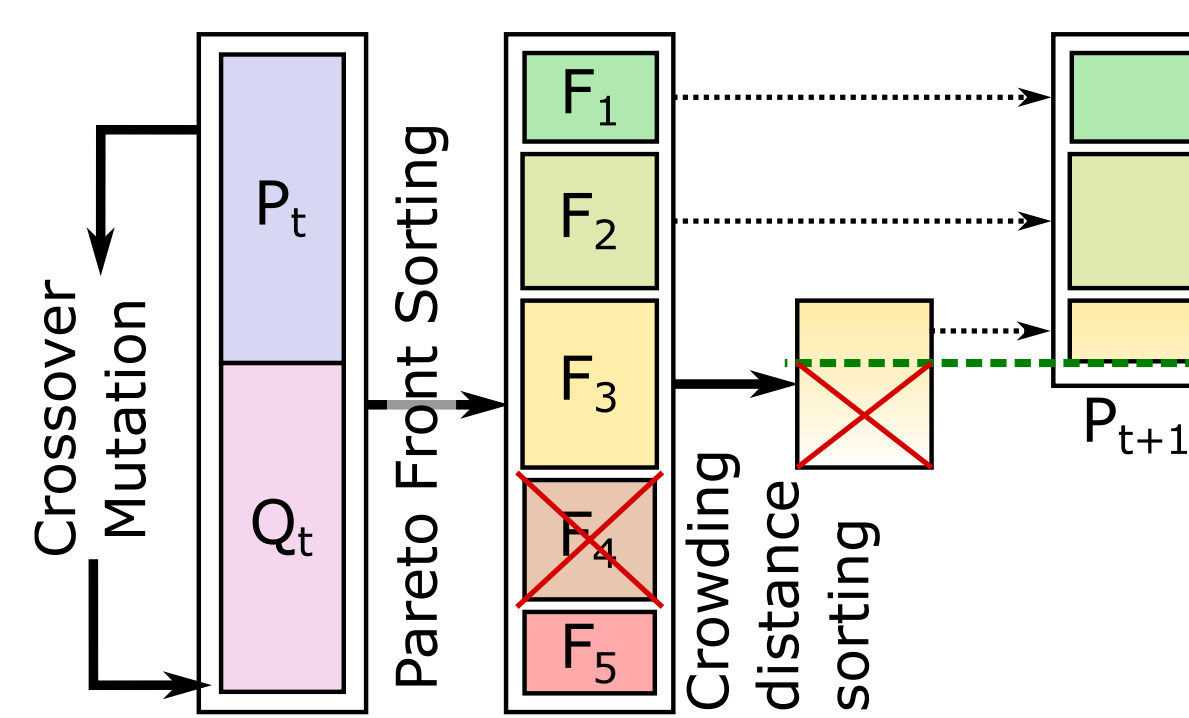
### Our Novel Contributions

- Analytical models of DNN and CapsNet layers and operations for architectural flexibility and fast hardware estimation.
- Analysis and selection of the adversarial perturbations values to employ in the NAS for a fast robustness evaluation.
- Specialized evolutionary algorithm, based on the principles of the NSGA-II method, to perform a multi-objective Pareto frontier selection, with conjoint optimization for adversarial robustness, energy, memory, and latency of DNNs.
- Fast evaluation methodology for DNNs trained for a limited number of epochs to reduce the training time.
- Full-training evaluation of the Pareto-optimal solutions to obtain the exact results.

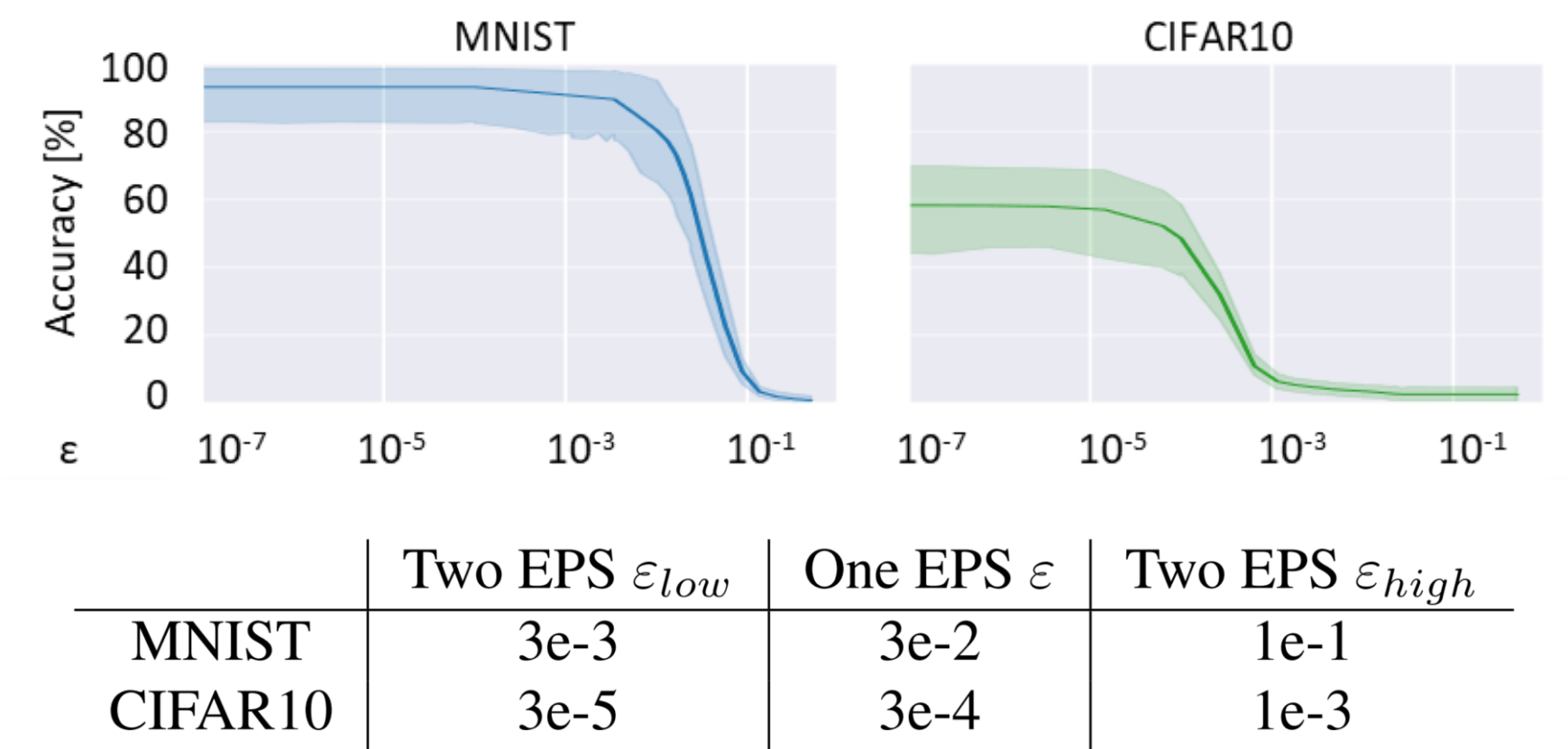
## Proposed HARNAS Framework



### One Iteration of the NSGA-II Algorithm

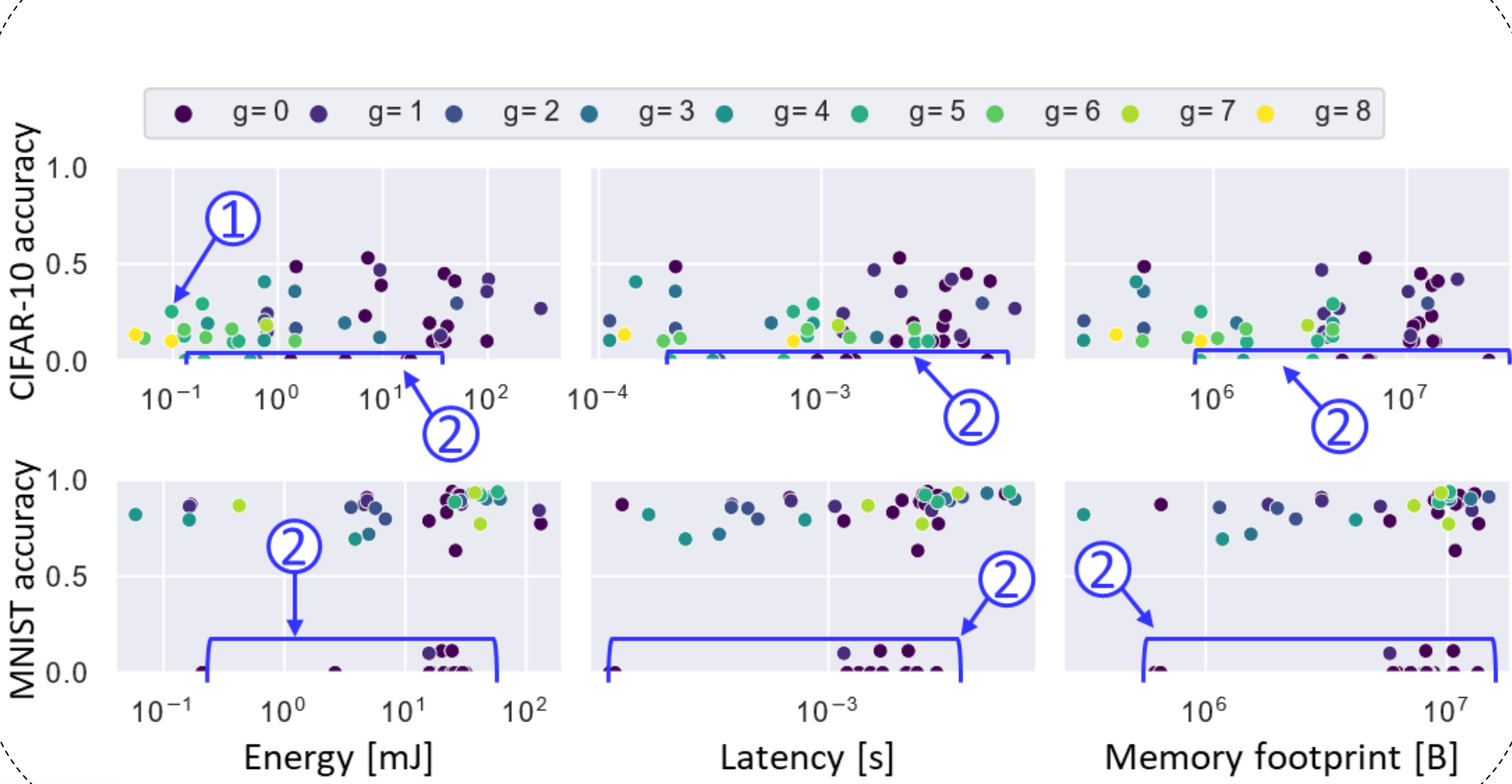


### Selection of Adversarial Perturbation for the NAS

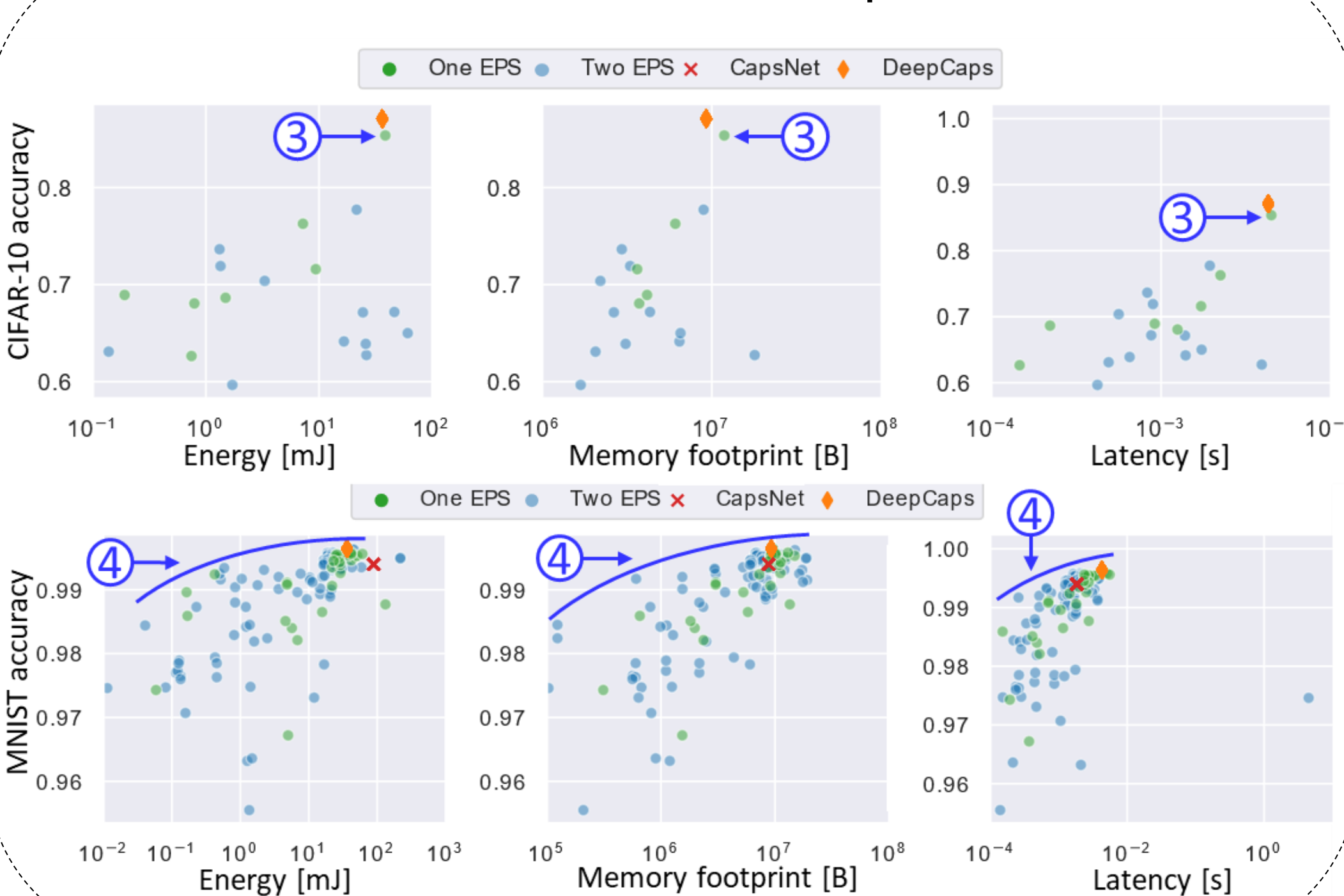


## Evaluation and Related Work Comparison

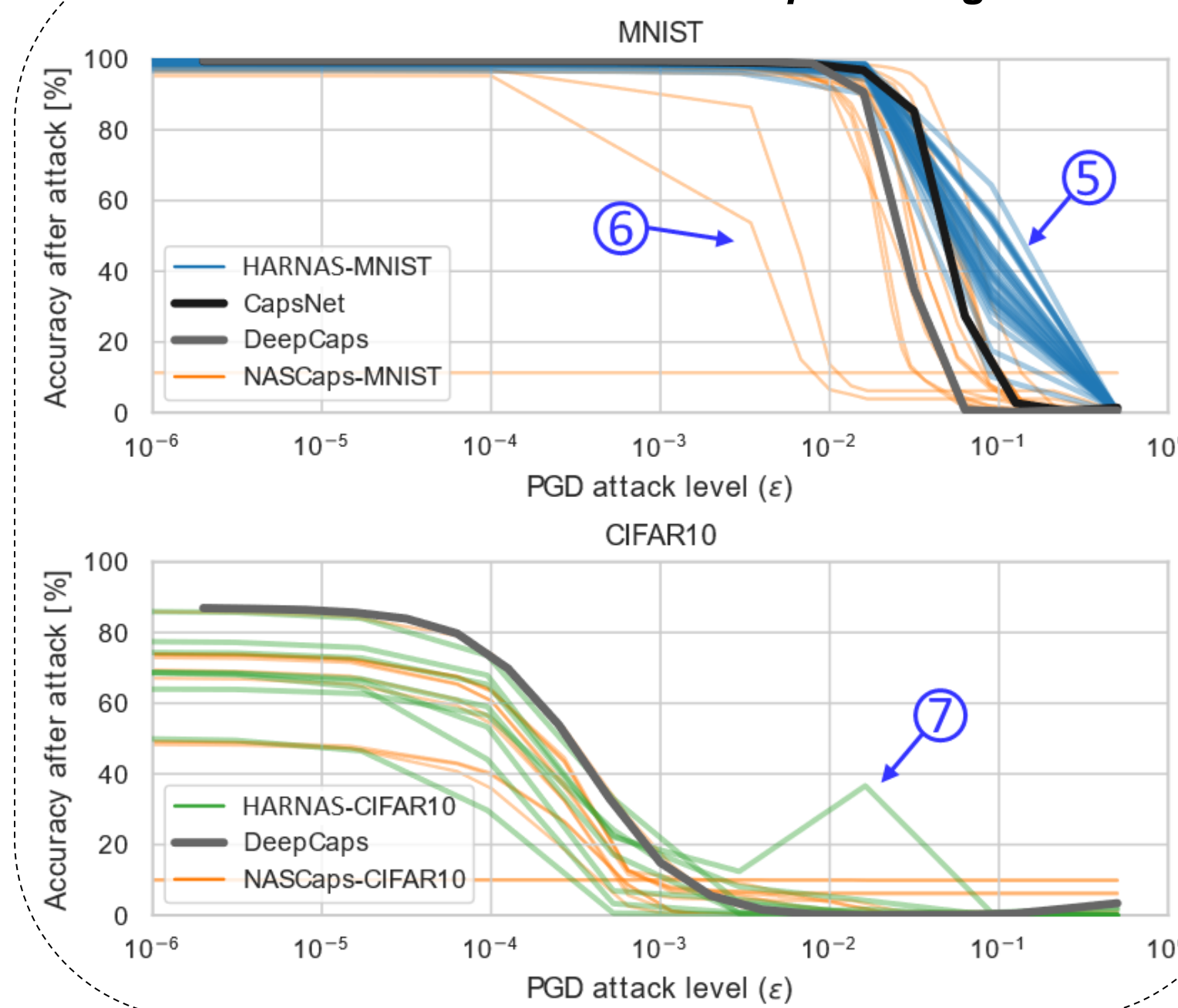
### HARNAS Results with Fast DNN Robustness



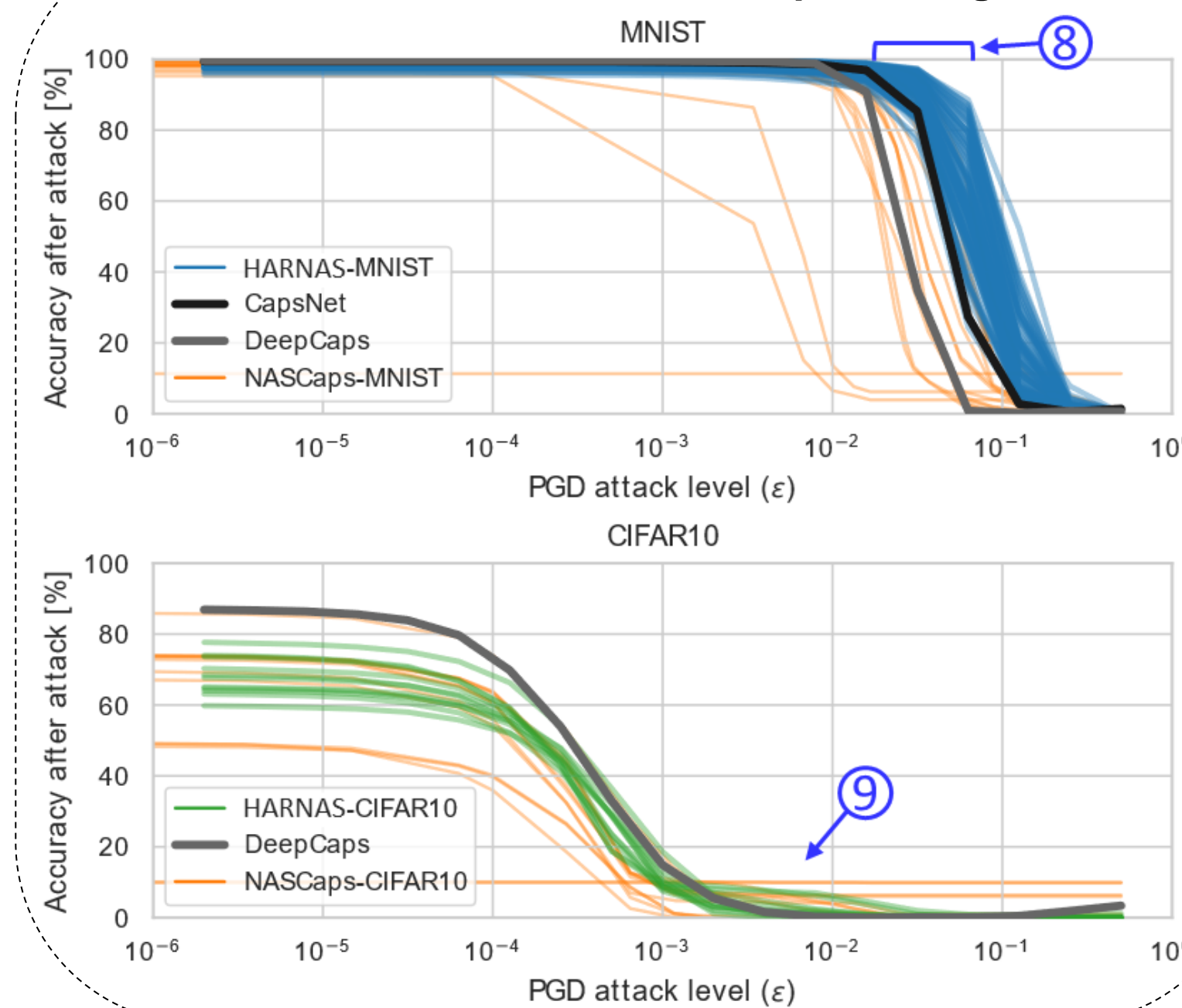
### HARNAS Exact Results for Pareto-Optimal DNNs



### HARNAS Results in One Eps Setting



### HARNAS Results in Two Eps Setting



### Key Observations

- For the HARNAS evaluated on the CIFAR10 dataset, the latest generations find DNNs that are less robust to the PGD attack, but still belong to the Pareto-frontier due to the low energy consumption.
- Several candidate DNNs found in the earliest generations are automatically discarded by the Pareto-frontier selection, since they are highly vulnerable to the PGD attack.
- A Pareto-optimal solution found by the HARNAS framework for the CIFAR10 dataset achieves 86.07% accuracy while having an energy consumption of 38.63 mJ, a memory footprint of 11.85 MiB, and a latency of 4.47 ms.
- The Pareto-optimal DNN search for MNIST covers a wider range of values, leveraging tradeoffs between different objectives.
- For the MNIST dataset, the Pareto-optimal solutions obtained with the HARNAS framework are particularly robust for a high range of perturbation  $\epsilon$ .
- The accuracy starts dropping at around one order of magnitude higher  $\epsilon$  than NASCaps.
- For the CIFAR10 dataset, the HARNAS DNNs' behavior is similar to the DeepCaps for low values of  $\epsilon$ , while a Pareto-optimal HARNAS solution offer a respectable robustness also with higher adversarial perturbation.
- The HARNAS framework with the Two EPS setting, compared to the One EPS setting, produces different levels of robustness w.r.t.  $\epsilon$  for the MNIST dataset.
- For the CIFAR10 dataset, the Two EPS search leads to worse results than the One EPS counterpart.