# The Dark Side of Dynamic Routing Neural Networks: Towards Efficiency Backdoor Injection

**Simin Chen**[1].      Hanlin Chen[2].     Mirazul Haque[1].
Cong Liu[3].        Wei Yang[1].

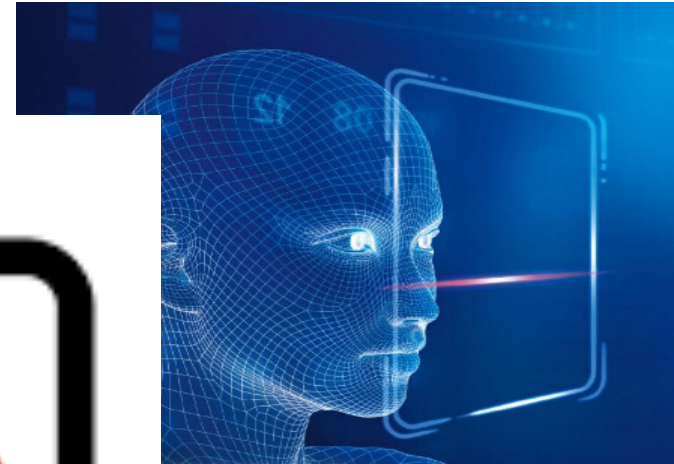[1]University of Texas at Dallas.        [2]Purdue University.      [3]University of California at Riverside
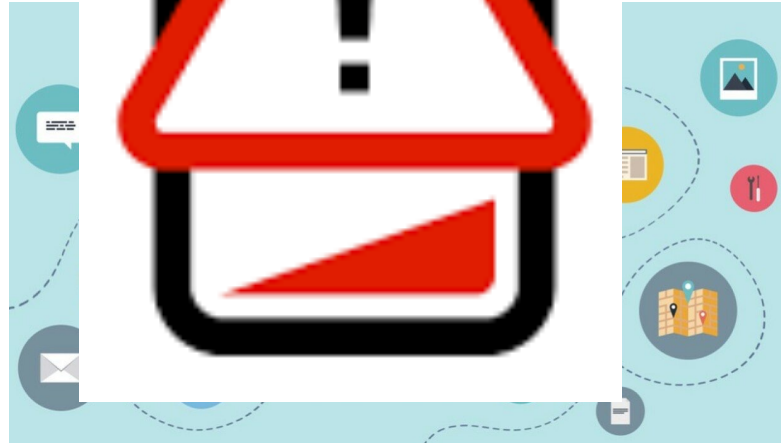
# Deep Learning is Pervasive on Edge Computing
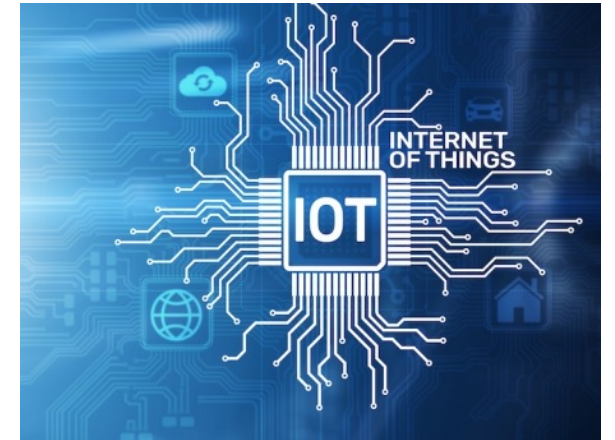


Autonomous-driving



Facial recognition



Robotics
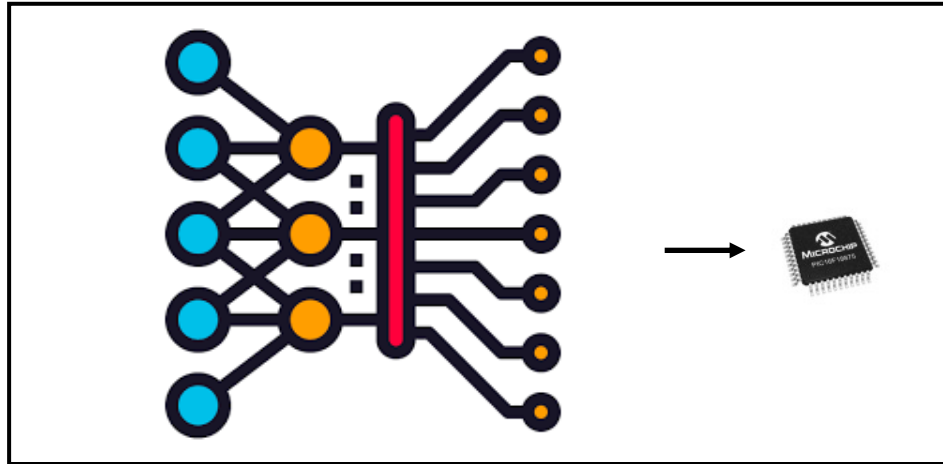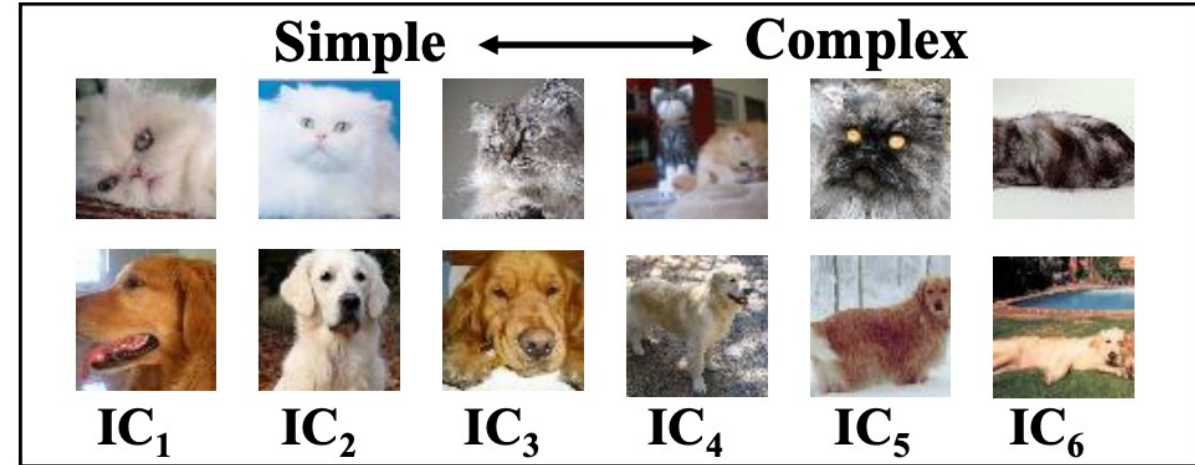


Mobile Application



IoT Application

# Not All Inputs Require the Same Computations

The real-time requirements of on-device DNN applications

Not all inputs require the same computations (The figure is taken from [1])

Dynamic Neural Networks [2]

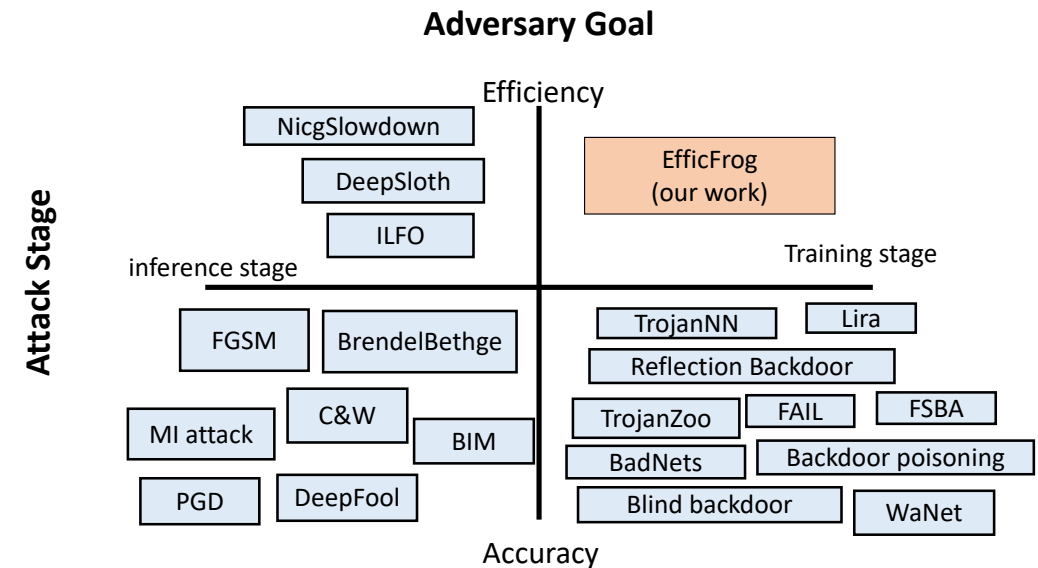[1] Shallow-Deep Networks: Understanding and Mitigating Network Overthinking (ICML 2020)
[2] Dynamic Neural Networks: A Survey   (TPAMI 2021)

# Background

Dynamic Neural Networks (DyNNs) allocate different computational resources for different inputs.

Confident Enough (Early Exit)

| DNN Block | → | Decision Unit | → | DNN Block | → | Decision Unit | → | DNN Block |

Y

N

Not Used

Can we inject efficiency backdoors into DyNNs that only affect DyNNs' computational efficiency on triggered inputs, while keeping DyNNs' behavior in terms of accuracy and efficiency unchanged on benign inputs?

**Adversary Goal**

Efficiency

**Attack Stage**

| NicgSlowdown | | EfficFrog (our work) |

| DeepSloth |

| ILFO |

inference stage

Training stage

| FGSM | BrendelBethge |

| TrojanNN | Lira |

| Reflection Backdoor |

| C&W |

| TrojanZoo | FAIL | FSBA |

| MI attack | BIM |

| BadNets | Backdoor poisoning |

| PGD | DeepFool |

| Blind backdoor | WaNet |

Accuracy

4

# Problem Formulation

We summarize three properties of the efficiency-based backdoor attacks

*a.Unnoticeable to users*
*b.Effective to degrade model efficiency on triggered inputs*
*c.Behave normally on benign inputs*

$$\theta^* = \boxed{\text{argmax}_\theta \quad \mathbb{E}_{x \in \mathcal{D}}[\text{FLOPs}(\mathcal{F}, \theta, x \oplus r)]}$$

$$s.t. \quad \boxed{||r|| \quad \leq \quad \epsilon}$$

$$\boxed{\begin{array}{c} \text{Acc}(\mathcal{F}, \theta, x) \approx \text{Acc}(\mathcal{F}, \hat{\theta}, x) \\ \text{FLOPs}(\mathcal{F}, \theta, x) \approx \text{FLOPs}(\mathcal{F}, \hat{\theta}, x) \end{array}}$$

# Design Overview

The uniform distribution is the optimal target distribution that will result in the DyNN model consuming the most computational resources among all distributions



**Algorithm 1:** Algorithm to inject backdoor.

**Require:**

A set of labeled training data $\mathcal{D}$;
A pre-defined adversarial budget $\epsilon$;
A pre-defined poisoning ratio $p$;
balance hyper-parameters $\lambda_1, \lambda_2$;

1:   $r = \text{GenerateRandom}()$
2:   Load parameters $\theta$ from a clean model $\mathcal{F}$
3:   **for** each epoch **do**
4:      Compute $Loss_{per}$ on $(r, \epsilon)$ based on Eq. 3
5:      Compute $Loss_{uncertain}$ on $(x, \mathcal{U})$ based on Eq. 4
6:      $L = \lambda_1 \times Loss_{per} + \lambda_2 \times Loss_{uncertain}$
7:      $r-= \frac{\partial L}{\partial r}$
8:   **end for**
9:   **for** each epoch **do**
10:      Get batch $(x, y)$ from $\mathcal{D}$
11:      **if** RANDOM$() \leq p$ **then**
12:         $x^* = x \oplus r$
13:      **end if**
14:      Compute $Loss_1$ on $(x, y)$ based on Eq. 5
15:      Compute $Loss_2$ on $(x^*, \mathcal{U})$ based on Eq. 6
16:      $L = \lambda_1 \times Loss_1 + \lambda_2 \times Loss_2$
17:      $\theta-= \frac{\partial L}{\partial \theta}$
18:   **end for**
19:   Return $\theta$

6

**Datasets**: CIFAR-10, and Tiny-ImageNet.

**Backbone Networks**: VGG19, MobileNet, and ResNet5.

**DyNN Training Algorithms**: IC-Training and ShallowDeep.

**Effectiveness Evaluation:**

   (i) Computational complexity on triggered inputs and (ii) EEC Scores

**Stealthiness Evaluation:**

   (i)Symmetric Segment-Path Distance (SSPD) distance  (ii) the Hausdorff distance

# Evaluation (Effectiveness)

Table 1. Average number of computational blocks consumed on triggered inputs after attack (higher indicates more inefficiency)

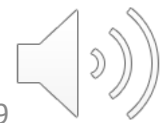| Backbone | Percentage | C10 | | | TI | | |
|---|---|---|---|---|---|---|---|
| | | BadNets | TrojanNN | EfficFrog | BadNets | TrojanNN | EfficFrog |
| VGG19 | 5% | 1.02 | 1.08 | **3.42** | 1.09 | 1.13 | **3.94** |
| | 10% | 1.02 | 1.06 | **3.92** | 1.09 | 1.13 | **4.12** |
| | 15% | 1.02 | 1.03 | **4.10** | 1.07 | 1.10 | **4.32** |
| MobileNet | 5% | 1.01 | 1.07 | **2.92** | 1.04 | 1.05 | **3.25** |
| | 10% | 1.01 | 1.04 | **3.51** | 1.04 | 1.08 | **3.56** |
| | 15% | 1.01 | 1.03 | **3.74** | 1.03 | 1.06 | **3.88** |
| ResNet56 | 5% | 1.06 | 1.08 | **4.04** | 1.07 | 1.09 | **4.01** |
| | 10% | 1.04 | 1.09 | **4.39** | 1.06 | 1.08 | **4.21** |
| | 15% | 1.04 | 1.04 | **4.48** | 1.04 | 1.09 | **4.56** |

Table 2. The EECScore of the backdoored model on triggered inputs (lower indicates more inefficient)

| Backbone | Percentage | C10 | | | TI | | |
|---|---|---|---|---|---|---|---|
| | | BadNets | TrojanNN | EfficFrog | BadNets | TrojanNN | EfficFrog |
| VGG19 | 5% | 0.93 | 0.93 | **0.55** | 0.92 | 0.92 | **0.50** |
| | 10% | 0.93 | 0.93 | **0.55** | 0.92 | 0.92 | **0.50** |
| | 15% | 0.93 | 0.93 | **0.56** | 0.92 | 0.92 | **0.51** |
| MobileNet | 5% | 0.91 | 0.91 | **0.68** | 0.91 | 0.91 | **0.53** |
| | 10% | 0.92 | 0.92 | **0.68** | 0.91 | 0.91 | **0.54** |
| | 15% | 0.92 | 0.92 | **0.68** | 0.91 | 0.91 | **0.54** |
| ResNet56 | 5% | 0.92 | 0.92 | **0.55** | 0.92 | 0.92 | **0.49** |
| | 10% | 0.92 | 0.92 | **0.55** | 0.92 | 0.92 | **0.49** |
| | 15% | 0.92 | 0.92 | **0.55** | 0.92 | 0.92 | **0.49** |

# Evaluation (Stealthiness)

Table 3. The similarity score between the performance curve, the rate column is computed using the smaller score divided the larger score.

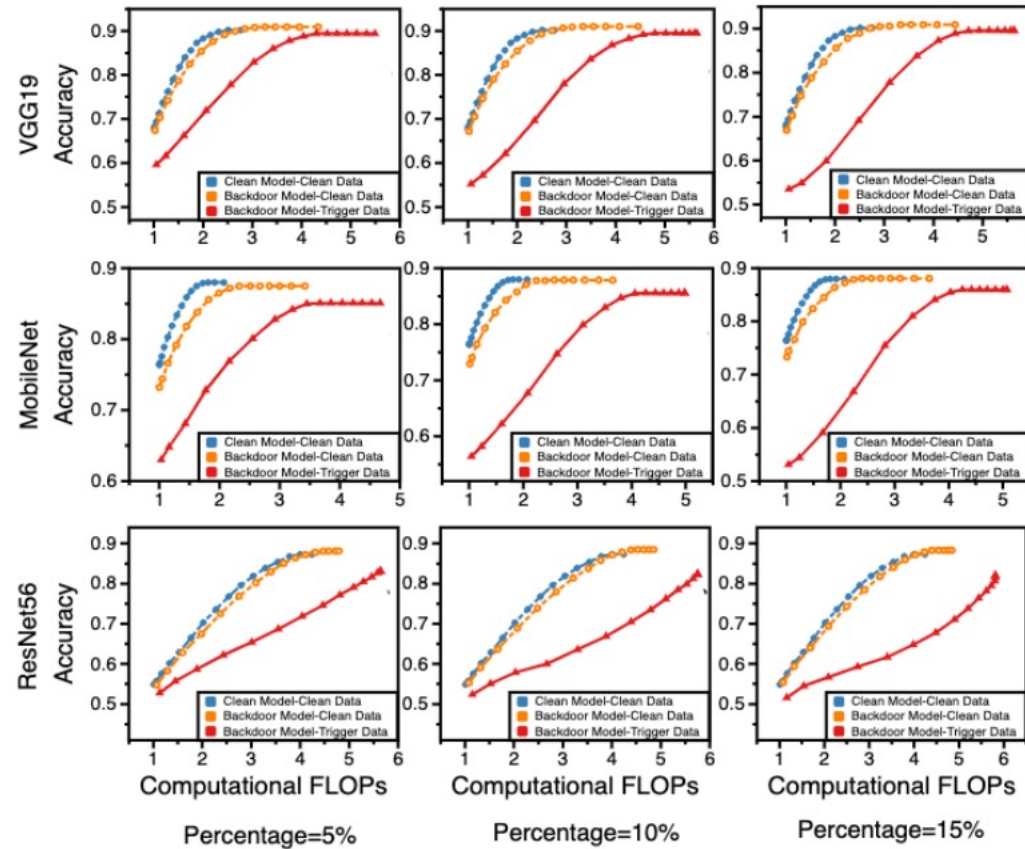| Backbone | Percentage | SSPD | | | Hausdorff | | | SSPD | | | Hausdorff | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | CC-BC | CC-BB | Rate | CC-BC | CC-BB | Rate | CC-BC | CC-BB | Rate | CC-BC | CC-BB | Rate |
| VGG19 | 5% | 0.18 | 1.58 | 0.11 | 20.83 | 2.73 | 0.13 | 0.19 | 1.52 | 0.13 | 29.28 | 3.24 | 0.11 |
| | 10% | 0.20 | 1.70 | 0.12 | 26.64 | 2.89 | 0.11 | 0.17 | 1.32 | 0.13 | 33.12 | 3.26 | 0.10 |
| | 15% | 0.20 | 1.68 | 0.12 | 28.33 | 2.88 | 0.10 | 0.16 | 1.27 | 0.13 | 34.42 | 3.21 | 0.09 |
| MobileNet | 5% | 0.20 | 1.35 | 0.15 | 21.30 | 2.59 | 0.12 | 0.22 | 1.48 | 0.15 | 28.15 | 2.93 | 0.10 |
| | 10% | 0.23 | 1.57 | 0.14 | 28.71 | 2.90 | 0.10 | 0.22 | 1.52 | 0.15 | 33.85 | 3.14 | 0.09 |
| | 15% | 0.22 | 1.57 | 0.14 | 31.26 | 2.99 | 0.10 | 0.23 | 1.59 | 0.15 | 35.61 | 3.23 | 0.09 |
| ResNet56 | 5% | 0.07 | 0.54 | 0.12 | 12.01 | 1.40 | 0.12 | 0.15 | 1.13 | 0.13 | 19.10 | 2.25 | 0.12 |
| | 10% | 0.08 | 0.61 | 0.12 | 14.44 | 1.51 | 0.10 | 0.17 | 1.20 | 0.14 | 24.73 | 2.58 | 0.10 |
| | 15% | 0.08 | 0.59 | 0.13 | 15.40 | 1.57 | 0.10 | 0.18 | 1.18 | 0.15 | 27.05 | 2.78 | 0.10 |

Figure 5. Efficiency and Accuracy degradation plot before and after `EfficFrog` launched

# Conclusion

- We characterize the efficiency backdoor vulnerability in DyNN models

- We propose an attack algorithm to inject efficiency backdoors into DyNN models

- Evaluation results suggest the effectiveness of our proposed methods.

# Thank You