GROWING EFFICIENT ACCURATE AND ROBUST NEURAL NETWORKS ON THE EDGE

Anonymous authors

Paper under double-blind review

ABSTRACT

The ubiquitous deployment of deep learning systems on resource-constrained Edge devices is hindered by their high computational complexity coupled with their fragility to out-of-distribution (OOD) data, especially to naturally occurring common corruptions. Current solutions rely on the Cloud to train and compress models before deploying to the Edge. This incurs high energy and latency costs in transmitting locally acquired field data to the Cloud while also raising privacy concerns. We propose Growing Efficient, Accurate, and Robust neural networks (GEARnn) to grow and train robust networks in-situ, i.e., completely on the Edge device. Starting with a low-complexity initial backbone network, GEARnn employs One-Shot Growth (OSG) to grow a network satisfying the memory constraints of the Edge device using clean data, and robustifies the network using Efficient Robust Augmentation (ERA) to obtain the final network. We demonstrate results on a NVIDIA Jetson Xavier NX, and analyze the trade-offs between accuracy, robustness, model size, energy consumption, and training time. Our results demonstrate the construction of efficient, accurate, and robust networks entirely on an Edge device.

1 INTRODUCTION

027 028

004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

024

025 026

029 The ubiquitous practical deployment of deep neural networks is mainly hindered by their lack of robustness and high computational cost. Prior art has shown that these deep networks are ex-031 tremely fragile to adversarial perturbations Szegedy et al. (2013)Goodfellow et al. (2014) and out-of-distribution (OOD) data Hendrycks & Dietterich (2019)Mintun et al. (2021). Natural corruptions Hendrycks & Dietterich (2019) (a specific type of OOD data) are more commonly encountered 033 at the Edge where real-time data is being continually acquired, e.g., video sequences acquired by 034 on-board cameras in autonomous agents (self-driving cars, field robots, drones), which tend to be distorted by weather and blur. The state-of-the-art defense against these corruptions employs robust data augmentation Hendrycks et al. (2019; 2021); Modas et al. (2022) which incurs a huge 037 computational cost when implemented on an Edge device. Fig. 1 indicates that it takes more than 2 days to robustly train a VGG-19 network Simonyan & Zisserman (2014) on a simple CIFAR-10 dataset when implemented on NVIDIA Jetson Xavier NX Edge device NVIDIA (a). Even for a small 040 5% VGG-19 network it takes more than a day, thus highlighting the non-trivial nature of the problem. 041 This is a huge concern because Edge devices are typically battery-powered and such large training 042 costs reduce their operational life-time.

043 Traditional solutions for reducing network complexity such as pruning Han et al. (2015); Li et al. 044 (2016); Diffenderfer et al. (2021), quantization Rastegari et al. (2016); Hubara et al. (2016) and neural architecture search (NAS) Liu et al. (2018); Zoph et al. (2018) mainly target Edge inference, 046 and are not suited for Edge training since they start with hard-to-fit over-parameterized networks 047 that require the large computational resources of the Cloud. However, transmitting local data to 048 the Cloud incurs energy and latency costs while raising privacy concerns, thus requiring training to 049 happen completely on the Edge. Given the above challenges, the primary objective of our work is: To design and train compact robust networks completely in-situ on Edge devices. Our proposed solution 050 Growing Efficient, Accurate, and Robust neural networks (GEARnn) is based on the family of growth 051 algorithms Chen et al. (2015); Wu et al. (2020); Evci et al. (2022); Yuan et al. (2020) that gradually 052 increase the size of an initial backbone network to reach the robust accuracy of a full network but at a fraction of its size, training complexity, and energy consumption.

054 Prior work on network growth Wu et al. (2020; 055 2019); Yuan et al. (2023) do not consider robustness to common corruptions since they use 057 clean data during training, while works that consider robustness train fixed-sized networks using augmented data Hendrycks et al. (2019); Modas et al. (2022) without considering the ef-060 ficiency of robust training. Hence, in order to 061 grow robust networks on the Edge and achieve 062 good robustness vs. training efficiency trade-063 off, we ask the following questions: Q1) should 064 networks be grown using augmented data only 065 (1-Phase), or should they be grown using clean 066 data first and then trained with augmented data 067 (2-Phase)? Q2) for growth, how many steps 068 should be employed? We answer these questions by proposing our method GEARnn to effi-069 ciently grow robust networks. Fig. 1 shows that GEARnn achieves significant improvements in 071 robust accuracy over vanilla trained baselines 072 while requiring much smaller training energy 073 consumption compared to robustly trained base-074 lines. 075

Contributions: We make the following contributions (Fig. 2):



Figure 1: Improvements in robust accuracy, training time, and model size (area of circles) of our proposed GEARnn method measured on NVIDIA Jetson Xavier NX Edge device NVIDIA (a). Robust accuracy is evaluated on CIFAR-10-C for GEARnn, full network baselines (VGG-19), and small network baselines (5% VGG-19 networks with same topology as GEARnn-2). For robust training, we employ AugMix Hendrycks et al. (2019). GEARnn demonstrates significant reduction in training complexity over robust baselines at similar robust and clean accuracies (shown in Section 6.2).

078 079

081

082

084

087

090

091

092

094

095

1. <u>Problem Statement</u>: To the best of our knowledge, our work is the first to *grow* networks robust to common corruptions and the first to train robust networks efficiently on an Edge device.

- 2. Key Questions: We answer Q1 as: 2-Phase (growth with clean data followed by robust training using augmented data) provides improved robustness over a 1-Phase (growth using augmented data) at iso-model size. This result indicates the importance of proper initialization for efficient robust training (Sections 6.1, 6.2 & 6.3). We answer Q2 as: One-Shot Growth (OSG) achieves the best training efficiency, clean and robust accuracies at iso-model size compared to *m*-Shot (m > 1) Growth (Section 6.3).
- 3. Algorithm: We propose two Growing Efficient Accurate and Robust neural networks (GEARnn) algorithms (see Fig. 2 and Section 4.3) by combining 1-Phase/2-Phase with OSG and Efficient Robust Augmentation (ERA). We show that GEARnn generated networks shine on all four metrics simultaneously clean accuracy, robust accuracy, training efficiency and inference efficiency by implementing them on a real-life Edge device, the NVIDIA Jetson Xavier NX (Section 6.2).
 - 4. Interpretability: We explain the network topologies generated during OSG, and also provide rationale for the efficacy of 2-Phase approach (Section 8).

⁰⁹⁶ 2 BACKGROUND AND RELATED WORK

Robust Data Augmentation: This is the most commonly used method for addressing corruptions 098 due to its ease of integration into the training flow and ability to replicate low-level structural 099 distortions. AugMix Hendrycks et al. (2019), PRIME Modas et al. (2022) and FourierMix Sun 100 et al. (2021) combine chains of stochastic image transforms and enforce consistency using a suitable 101 loss function to generate an augmented sample from a clean image. DeepAugment Hendrycks 102 et al. (2021) randomly distorts the parameters of an image-to-image network to generate augmented 103 images. CARDs Diffenderfer et al. (2021) combines data augmentation Hendrycks et al. (2019) and pruning Frankle & Carbin (2018) to find compact robust networks embedded in large over-104 parameterized networks. Adversarial augmentations Zhao et al. (2020); Rusak et al. (2020); Calian 105 et al. (2021) have also been proposed to handle common corruptions. Unlike our proposed GEARnn 106 algorithm, all these techniques significantly increase the complexity over vanilla training and are thus 107 inappropriate for Edge deployment.

108 Growth Techniques: A typical growth algorithm 110 starts with a small ini-111 tial backbone model whose size is gradually increased 112 until the desired perfor-113 mance or network topology 114 is reached. Neural network 115 growth has been previously 116 used in optimization Fuku-117 mizu & Amari (2000), con-118 tinual learning Rusu et al. 119 (2016); Hung et al. (2019) 120 and in speeding up the train-121 ing of large networks Chen 122 et al. (2015). Recent 123 works Evci et al. (2022); Yuan et al. (2023) look at 124 improving the training dy-125



Figure 2: Proposed approach: GEARnn-1 performs One-Shot Growth (OSG) on augmented data (\mathcal{D}_{aug}) generated by Efficient Robust Augmentation (ERA) (using clean data (\mathcal{D}_{in})) in a single phase (1-Phase). GEARnn-2 performs OSG using \mathcal{D}_{in} first followed by parametric training on \mathcal{D}_{aug} in two consecutive phases (2-phase). Here \mathcal{L}_{CE} and \mathcal{L}_{aug} denote the cross-entropy loss and augmented loss, respectively.

namics and efficiency for growth by using better neuron initializations. Others find efficient networks by growing the width Wu et al. (2019), depth Wen et al. (2020) or both Wu et al. (2020); Yuan et al. (2020). However, none of these methods address the issue of robustness to common corruptions or demonstrate the utility for training on a resource-constrained Edge setting, which is our focus. Though our work GEARnn builds upon Firefly Wu et al. (2020), it is flexible and can incorporate other growth methods mentioned above.

3 NOTATION AND PROBLEM SETUP

132 133

Notation: Let $f : \mathbb{R}^d \to [C]$ be a hard classifier which classifies input $\mathbf{x} \in \mathbb{R}^d$ into one of C classes. We choose f to be a convolutional neural network (CNN) with L layers (depth), $\{w_l\}_{l=1}^L$ output channels (widths), and (K, K) sized kernels. The network f is trained on n samples $(\mathbf{x}, y) \sim \mathcal{D}_{in}$, where $(\mathbf{x}, y) \in \mathbb{R}^d \times [C]$ and \mathcal{D}_{in} denotes the "in-distribution" or "clean" data. \mathcal{L}_{CE} represents the cross-entropy loss and $\mathcal{L}_{aug} = \mathcal{L}_{CE} + \lambda \mathcal{L}_{JSD}$ represents the augmentation loss where \mathcal{L}_{JSD} is the Jensen-Shannon divergence loss described in Hendrycks et al. (2019).

140 During inference, f can be exposed to samples from both \mathcal{D}_{in} and \mathcal{D}_{out} ("out-of-distribution" or 141 "corrupted" data). In case of common corruptions, $(\mathbf{x}_{out}, y) \sim \mathcal{D}_{out}$ is obtained by $\mathbf{x}_{out} = \kappa(\mathbf{x}_{in}, s)$, 142 where $(\mathbf{x}_{in}, y) \sim \mathcal{D}_{in}$, κ is a corruption filter and s is the severity level of the corruption. We denote 143 $p_e = \Pr(\hat{y} \neq y)$ as the classification error at inference where $\hat{y} = f(\mathbf{x}_{test})$. When $(\mathbf{x}_{test}, y) \sim \mathcal{D}_{in}$, 144 we define $(1 - p_e)$ as clean accuracy \mathcal{A}_{cln} , and when $(\mathbf{x}_{test}, y) \sim \mathcal{D}_{out}$, we define $(1 - p_e)$ as robust 145 accuracy \mathcal{A}_{rob} . The value of p_e is determined empirically in this work.

Problem: Our primary objective is to maximize the empirical clean and robust accuracies (\mathcal{A}_{cln} and \mathcal{A}_{rob}) while ensuring the network complexity ($\sum_{l=1}^{L} w_l$) is small. Along with these two criteria, we also prioritize reduction in training time (t_{tr}) and training energy consumption (E) on hardware.

4 GROWING EFFICIENT ACCURATE AND ROBUST NEURAL NETWORKS (GEARNN)

152 153 As shown in

As shown in Fig. 2, two flavors of GEARnn algorithms are proposed – GEARnn-1 and GEARnn-2.
While GEARnn-1 leverages the 1-Phase (joint growth and robust training) training, GEARnn-2 employs the 2-Phase (sequential growth and robust training) approach. Both flavors incorporate One-Shot Growth (OSG) and Efficient Robust Augmentation (ERA) in different ways. In this section, we first describe OSG and ERA, and then formally present the GEARnn algorithms.

158 4.1 ONE-SHOT GROWTH (OSG)

159

149

150

151

160 One-Shot Growth (OSG) employs labeled data to perform a single growth step sandwiched between 161 two training stages. The initial backbone f_0 is first trained for \mathcal{E}_1 epochs. The resulting network f_1 is 162 grown over \mathcal{E}_q epochs to obtain the grown network f_q , i.e., $f_q = \mathcal{G}(f_1|\gamma, \mathcal{D}, \mathcal{L}, \mathcal{E}_q)$, where \mathcal{G} is the



Figure 3: OSG takes in labeled data (\mathcal{D}) and backbone network f_0 , and performs a training step, 170 a growth step, and a training step in sequence to generate network f_2 . The 2-tuple $(\mathcal{L}, \mathcal{E}) = (loss$ function, number of epochs) employed in each step. 172

growth technique which is nominally Firefly Wu et al. (2020) in our work. The final network f_2 is 173 obtained by training f_q over \mathcal{E}_2 epochs. Either clean (\mathcal{D}_{in}) or augmented (\mathcal{D}_{aug}) data can be used in 174 OSG. For instance, OSG in GEARnn-1 and GEARnn-2 employs augmented data ($\mathcal{D} \sim \mathcal{D}_{aug}$) and 175 clean data ($\mathcal{D} \sim \mathcal{D}_{in}$), respectively. 176

The growth technique \mathcal{G} is described below: 177

$$f_{g} = \underset{f}{\operatorname{arg\,min}} \quad \mathcal{L}(f, \mathcal{D}|f_{1})$$

s.t. $f \in \partial(f_{1}, \epsilon)$
 $\mathcal{C}(f) \leq (1 + \gamma) \mathcal{C}(f_{1})$ (1)

where $\partial(f_1, \epsilon)$ represents the growth neighbourhood for topology search, $\mathcal{C}(f) = \sum_{l=1}^{L} w_l$ represents 183 the complexity estimate of network f and γ denotes the growth ratio. The neighbourhood $\partial(f_1, \epsilon)$ is expanded in two ways - splitting and growing new neurons - as described in Wu et al. (2020; 2019). 185 We perform growth only in the width dimension and keep the number of layers L and the kernel size 186 (K, K) constant for reasons described in Wu et al. (2020) and Simonyan & Zisserman (2014). 187

Existing growth methods Wu et al. (2019; 2020); Evci et al. (2022) use several growth steps (typically 188 10 steps) and large number of training epochs (typically 1600 total epochs) which makes them 189 inefficient for training. This directs us to pick OSG over multi-step growth (validated in Section 6.3) 190 and reduce the training epochs significantly $(20 \times)$ compared to prior growth algorithms. The drop 191 in accuracy observed due to these modifications is compensated for using a 2-Phase approach 192 (Section 6.1). 193

4.2 EFFICIENT ROBUST AUGMENTATION (ERA) 194

Efficient Robust Augmentation (ERA) em-196 ploys clean data (\mathcal{D}_{in}) to generate aug-197 mented data (\mathcal{D}_{aug}) in an efficient manner. The clean sample x (where (x, y) \sim \mathcal{D}_{in}) is passed through a set of transforms 199 a_1, a_2, \dots, a_{d_i} to obtain the transformed 200 sample $A_i(x)$, which is then combined 201 linearly with the clean sample to give 202 the augmented sample $\mathbf{x}_{j}^{\text{aug}}$. We concate-203 nate (J - 1) such augmented samples 204 $\{\mathbf{x}_{j}^{\text{aug}}\}_{j=1}^{J-1}$ along with the clean sample to 205 obtain our Efficient Robust Augmentation 206 $\mathcal{R}((\mathbf{x}, y)|\mathcal{T}).$ 207



Figure 4: ERA takes in clean data (\mathcal{D}_{in}) as input and applies a set of stochastic transforms to generate aug- $A_j(\mathbf{x}) = a_1 \circ a_2 \circ \dots \circ a_{d_j}(\mathbf{x}) \xrightarrow{\text{data}} (\mathcal{D}_{\text{aug}}) \text{ in an efficient manner.}$

$$\mathbf{x}_j^{\text{aug}} = p\mathbf{x} + (1-p)A_j(\mathbf{x})$$

aug

210 211

208

171

178 179

181 182

212 213

214

$$\mathcal{R}((\mathbf{x}, y) | \mathcal{T}) = (\{\mathbf{x}_{1}^{\text{aug}}, ..., \mathbf{x}_{J-1}^{\text{aug}}, \mathbf{x}\}, y) \implies \mathcal{D}_{\text{aug}} := \mathcal{R}(\mathcal{D}_{\text{in}} | \mathcal{T})$$

where $a_{i} \sim \text{Unif}(\mathcal{T}), \ p \sim \beta(1, 1), \ d_{j} \sim \text{Unif}(\{1, ..., D\}), \ j \in \{1, ..., J-1\}$
(2)

where \mathcal{T} denotes the set of transforms, $\beta()$ and Unif() represent the beta and uniform distributions, 215 respectively. SOTA robust data augmentation Hendrycks et al. (2019); Modas et al. (2022); Sun et al.

 $\mathbf{v} \rightarrow \mathcal{D}$

216 (2021) methods for common corruptions also employ stochastic chains of transforms with width W, 217 depth D, and enforce consistency across J-1 augmented and clean samples using the \mathcal{L}_{aug} loss 218 function (described in Section 3). The SOTA augmentation framework increases the training time 219 and energy by $3 \times$ to $4 \times$ compared to vanilla training. We choose (W, D, J) = (1, 3, 4) based on our 220 diagnosis (shown in Appendix B.1) to improve the efficiency without compromising on robustness compared to SOTA approaches Hendrycks et al. (2019); Modas et al. (2022). In GEARnn-2, grown 221 network f_2 (see Fig. 3) obtained using clean data OSG is trained for \mathcal{E}_r epochs using \mathcal{D}_{aug} generated 222 by ERA. 223

224 225

226

227

228

229

230

231

232

233

235

236

237

238

239

240

241

242

260

4.3 GEARNN ALGORITHMS

Algorithm 1 GEARnn-1 1: Input: clean training data \mathcal{D}_{in} , initial backbone network f_0 , growth ratio γ , set of augmentation transforms \mathcal{T} , training epochs $\{\mathcal{E}_1, \mathcal{E}_q, \mathcal{E}_2\}$ 2: **Output:** compact and robust model f_{1p}^* 3: /* Phase-1: OSG */ 4: for $e = 1, ..., \mathcal{E}_1$ do $\mathcal{D}_{aug} := \mathcal{R}(\mathcal{D}_{in}|\mathcal{T}) \quad /\!\!/ \text{ ERA}$ 5: 6: $f_1 \leftarrow \arg\min_{aug} (f, \mathcal{D}_{aug}|f_0)$ // backbone robust training 7: end for 8: $f_g \leftarrow \mathcal{G}(f_1|\gamma, \mathcal{D}_{aug}, \mathcal{L}_{aug}, \mathcal{E}_g)$ // augmented growth 9: for $e = 1, ..., \mathcal{E}_2$ do $\mathcal{D}_{aug} := \mathcal{R}(\mathcal{D}_{in}|\mathcal{T}) / \mathbf{ERA}$ 10: 11: $f_2 \leftarrow \arg \min \mathcal{L}_{aug}(f, \mathcal{D}_{aug}|f_g)$ // grown-network robust training 12: end for 13: $f_{1p}^* \leftarrow f_2$ 14: return f_{1p}^*

0.40	
243	Algorithm 2 GEARnn-2
244	1: Input: clean training data \mathcal{D}_{in} , initial backbone network f_0 , growth ratio
245	γ , set of augmentation transforms \mathcal{T} , training epochs $\{\mathcal{E}_1, \mathcal{E}_g, \mathcal{E}_2, \mathcal{E}_r\}$
246	2: Output: compact and robust model f_{2p}^*
247	3: /* Phase-1: OSG */
2/18	4: for $e = 1,, \mathcal{E}_1$ do
240	5: $f_1 \leftarrow \arg\min_{f} \mathcal{L}_{CE}(f, \mathcal{D}_{in} f_0)$ // backbone clean training
249	6: end for
250	7: $f_a \leftarrow \mathcal{G}(f_1 \gamma, \mathcal{D}_{in}, \mathcal{L}_{CE}, \mathcal{E}_a)$ // clean growth
251	8: for $e = 1,, \mathcal{E}_2$ do
252	9: $f_2 \leftarrow \arg \min \mathcal{L}_{CE}(f, \mathcal{D}_{in} f_q)$ // grown-network clean training
253	10: end for
254	11: /* Phase-2: Train */
255	12: for $e = 1,, \mathcal{E}_r$ do
256	13: $\mathcal{D}_{aug} := \mathcal{R}(\mathcal{D}_{in} \mathcal{T}) // ERA$
230	14: $f_{2p}^* \leftarrow \arg \min \mathcal{L}_{aug}(f, \mathcal{D}_{aug} f_2) // \text{ grown-network robust training}$
257	15: and for
258	16. return f*
259	10. rearring _{22p}

Algorithms 1 and 2 describe GEARnn-1 and GEARnn-2, respectively. Algorithms 1 and 2 output final compact and robust models f_{1p}^* and f_{2p}^* , respectively. For empirical results in Section 6, the growth technique \mathcal{G} and the set of transforms \mathcal{T} are chosen from Firefly Wu et al. (2020) and AugMix Hendrycks et al. (2019), respectively, though other growth Yuan et al. (2023); Wu et al. (2019) and augmentation Modas et al. (2022); Sun et al. (2021) methods can be substituted to obtain different GEARnn variants.

5 EXPERIMENTAL SETUP

Datasets and Architectures: All results are shown on CIFAR-10, CIFAR-100 Krizhevsky et al. (2009) and Tiny ImageNet Le & Yang (2015) (\mathcal{D}_{in}) datasets. CIFAR-10-C, CIFAR-100-C and Tiny ImageNet-C Hendrycks & Dietterich (2019) (\mathcal{D}_{out}) are used to benchmark corruption robustness. Convolutional neural network architectures MobileNet-V1Howard et al. (2017), VGG-19Simonyan & Zisserman (2014), ResNet-18He et al. (2016) are employed to demonstrate the results.

Hardware: For the server-based experiments, we use a single NVIDIA Quadro RTX 6000 GPU with 24GB RAM, 16.3 TFLOPS peak performance and an Intel Xeon Silver 4214R CPU. This machine is referred to as "Quadro". For the Edgebased experiments, we use the NVIDIA Jetson Xavier NX NVIDIA (a) which has

a Volta GPU with 8GB RAM, 21 TOPS peak performance and a Carmel CPU. We refer to this deviceas "Jetson".

Metrics: Clean accuracy $\mathcal{A}_{cln}(\%)$ measured on clean test data \mathcal{D}_{in} , and robust accuracy $\mathcal{A}_{rob}(\%)$ measured on corrupted test data \mathcal{D}_{out} , are used as accuracy metrics (both computed using Robust-Bench Croce et al. (2021)). The number of floating-point parameters (model size), wall-clock training time t_{tr} (in minutes), per-sample wall-clock inference time t_{inf} (in seconds) and energy consumption E (in Joule) are used as the efficiency metrics. Size (%) represents the fraction of the full model size. In case of growth algorithms, training times include both the time taken for training and growth. The power is measured from the Quadro and Jetson using Nvidia-SMI NVIDIA (b) and Jetson Stats Bonghi, respectively, and the energy E is computed by summing the mean power values polled. 270 **Baselines:** In the absence of prior work on robust growth, we propose our own baselines Small (\mathcal{D}_{in}) 271 and Small (\mathcal{D}_{aug}), both of which use 160 training epochs to be consistent with Diffenderfer et al. 272 (2021). They are networks with the same size and topology as the final GEARnn-2 network (f_{2n}^*) in 273 Fig. 2) trained with random initialization on clean data and augmented data (AugMix Hendrycks et al. 274 (2019), unless specified otherwise), respectively.

275 We pick Small (\mathcal{D}_{aug}) as the main baseline for a fair comparison with GEARnn as it depicts a typical 276 private-Edge training scenario. We do not compare with compression techniques since they have 277 been shown to have worse training efficiency compared to growth Yuan et al. (2020), and require a 278 robust-trained full baseline, and this is clearly more expensive than training Small (\mathcal{D}_{aug}) (see Fig. 1). 279

MAIN RESULTS 6

284

285

286

287

289

In this section, we first compare the performance of GEARnn across different network architectures and datasets on Quadro. We then show results for CIFAR-10 and CIFAR-100 using VGG-19 and MobileNet on Jetson. Finally, we compare OSG with 288 *m*-shot growth methods on Jetson.

6.1 **RESULTS ACROSS NETWORK** 290 ARCHITECTURES AND DATASETS 291

292 Table 12 shows GEARnn 293 is consistently better in 294 terms of training time and 295 training energy consump-296 tion over the best baseline 297 Small (\mathcal{D}_{aug}) over multiple 298 network architectures and 299 datasets. Specifically, an av-300 erage reduction in training time (energy consumption) 301 of 3.5×, 2.9× and 1.8× 302 $(3.7\times, 2.0\times \text{ and } 2.0\times)$ 303 is observed for CIFAR-10, 304 CIFAR-100 and Tiny Im-305 ageNet, respectively. Fur-306 thermore, we find GEARnn-307 1 is inferior to GEARnn-2

Table 1: GEARnn hyperparameters for different networks and datasets.

Detecat	Growth Ratio (γ)			Small(D)	GE	AR	nn-1	G	EA	Rnn	-2
Dataset	Mob.	VGG	Res.	ε	\mathcal{E}_1	\mathcal{E}_{g}	\mathcal{E}_2	\mathcal{E}_1	\mathcal{E}_{g}	\mathcal{E}_2	\mathcal{E}_r
CIFAR-10	1.8	0.9	0.6	160	40	1	40	40	1	40	40
CIFAR-100	2.0	1.5	0.8	160	50	1	50	40	1	40	50
Tiny ImageNet	2.0	1.5	0.8	160	50	1	50	40	1	40	50



Figure 5: GEARnn-2 achieves higher robustness at the same: (a) number of robust training epochs at final model size, and (b) training time, for VGG-19/CIFAR-100 on Quadro.

308 on all the four metrics thereby answering Q1 in Section 1 - 2-Phase approach is better than 1-Phase 309 approach for efficiently growing robust networks.

310 A key reason underlying GEARnn-2's training efficiency is the reduction in the number of robust 311 training epochs \mathcal{E}_r made possible by the OSG initialization in Phase-1. Fig. 5 shows that for the 312 same training time, GEARnn-2 provides better robustness than Small (\mathcal{D}_{aug}) and GEARnn-1. Similar 313 results were obtained for CIFAR-10 and other network architectures as shown in Appendix C.1. 314

6.2 **RESULTS ON THE EDGE** 315

316 We now study GEARnn when mapped onto the Edge device NVIDIA Jetson Xavier NX. The training hyperparameters for Jetson are described in Appendix A. Results on Jetson (Table 13) show similar 317 318 trends to those on Quadro (Table 12).

319 Specifically, Table 13 shows that GEARnn-2 achieves comparable clean and robust accuracies to 320 the baseline Small (\mathcal{D}_{aug}) but at a fraction of its training cost – a 2.3× (2.8×) reduction in training 321 time (training energy) when averaged across both networks and datasets. Additionally, GEARnn-2 322 beats GEARnn-1 on almost all metrics, again confirming our answer to Q1 in favour of 2-Phase. 323 Interestingly, GEARnn-2 achieves a clean accuracy within 1% of Small (\mathcal{D}_{in}) at a similar training cost. These results confirm that it is possible to grow efficient and robust networks on the Edge.

GEARnn-2

8 91.35

Small (\mathcal{D}_{in}) 5 92.69

81.96

70.57

56

31

270 8 67.95

Architecture			CIFAR-10)			CIFAR-10	00		1	Finy ImageN	et
(full model size)	Method	Size (%)	$\left. \begin{array}{c} \text{Accuracy} \\ \mathcal{A}_{cln}(\%) \ \mathcal{A}_{rob}(\%) \end{array} \right $	Train $t_{tr}(min)$	E(kJ) Size	Acc $A_{cln}(\%)$	uracy $\mathcal{A}_{rob}(\%)$	$ $ Train $t_{tr}(min)$	$\begin{bmatrix} ing \\ E(kJ) \end{bmatrix} \begin{bmatrix} Si \\ (4) \end{bmatrix}$	$ze Acc (6) \mathcal{A}_{cln}(\%)$	$\mathcal{A}_{rob}(\%) \left t \right $	Trair _{tr} (min)
	Small (\mathcal{D}_{in})	8	92.28 66.31	42	192 8	67.66	39.04	45	274	3 55.13	18.48 \downarrow	262

Table 2: Comparison of accuracy, robustness, and efficiency between the baselines and GEARnn across various network architectures for CIFAR-10, CIFAR-100 and Tiny ImageNet on Quadro.

			• 1		•		• 1	
VGG-19 (20M)	$\begin{array}{c c c c c c c c c c c c c c c c c c c $	93.08 91.25 92.18	85.73 215 82.86 86 83.77 53	1140 9 70.01 552 9 65.73 208 0 68.44	56.94 219 52.68 111 54.31 65	927 9 55.51 779 9 54.38 566 0 56 19	<u>30.01</u> 668 28.56 428 29.79 357	7120 4220 3210
	ULAKIII-2 J	92.10	<u>03.77 <u>35</u></u>	<u>238</u> 9 08.44	<u>34.31 05</u>	<u>300</u> 3 <u>30.13</u>	23.13 <u>351</u>	3413
	Small $(\mathcal{D}_{in}) \parallel 6$	93.34	68.85↓ 61	546 7 68.74	40.83↓ 67	490 7 54.72	18.11 ↓ 381	3390
ResNet-18	Small $(\mathcal{D}_{aug}) = 6$	<u>94.18</u>	86.50 217	1730 7 <u>71.97</u>	57.30 219	1250 7 54.50	25.74 1103	12400
(12M)	GEARnn-1 6	92.36	83.86 108	747 8 69.15	55.62 142	1020 7 53.17	24.93 898	9100
	GEARnn-2 6	93.14	84.45 <u>77</u>	<u>567</u> 7 70.94	56.54 <u>97</u>	<u>905</u> 7 <u>54.79</u>	<u>26.64</u> <u>649</u>	<u>7270</u>
Table 3: C	Comparison o	f accu	racy, robust	ness, inference	and training	efficiency bety	ween the bas	elines

241 9 68.07

72

38

432

53.28

41.24

56.16

335 9 53.9

28.56

17.78

429

218

<u>3565</u>

2040

and GEARnn for CIFAR-10 and CIFAR-100 using MobileNet-V1 and VGG-19 on Jetson. Due to computational limitations, the results for Tiny ImageNet and ResNet-18 are excluded for Jetson.

			C	CIFAR	-10					CIFAR	R-100		
Network	Method	Accu Accu Accu Accu Accu Accu Accu Accu	$\left \begin{array}{c} \operatorname{tracy} \\ \mathcal{A}_{\operatorname{rob}}(\%) \end{array} \right \operatorname{Si}$	Infer ize%	ence t _{inf} (ms)	Train $t_{tr}(min)$	$\begin{bmatrix} ing \\ E(kJ) \end{bmatrix}$	$Accu \mathcal{A}_{cln}(\%)$	$\mathcal{A}_{rob}(\%)$	Infe Size%	rence $t_{inf}(ms)$	$\begin{vmatrix} \text{Train} \\ t_{\rm tr}(\min) \end{vmatrix}$	ing E(kJ)
	$\ $ Small $(\mathcal{D}_{in}) \ $	91.88	68.35↓	7	0.9	675	175	68.59	39.47	8	0.9	744	166
MobileNet-V1	$ \begin{array}{ l l l l l l l l l l l l l l l l l l l$	<u>92.58</u> 90.20 91.43	83.84 79.65 81.64	7 7 7	<u>0.9</u> <u>0.9</u> <u>0.9</u>	1216 560 <u>553</u>	511 238 <u>162</u>	<u>69.24</u> 65.48 67.42	50.46 52.39	8 8 8	<u>0.9</u> 1.0 <u>0.9</u>	1333 704 <u>690</u>	586 226 <u>216</u>
	$\ $ Small $(\mathcal{D}_{in}) \ $	92.97	71.08↓	5	1.0	533	128	67.92	40.49	9	1.4	714	187
VGG-19	$\begin{array}{c} \text{Small} \left(\mathcal{D}_{\text{aug}} \right) \\ \text{GEARnn-1} \\ \text{GEARnn-2} \end{array}$	<u>93.36</u> 90.94 92.07	85.73 82.25 83.45	5 5 5	<u>1.0</u> 1.2 <u>1.0</u>	1543 652 <u>596</u>	522 207 155	70.07 62.89 67.59	<u>56.68</u> 49.63 53.64	9 9 9	<u>1.4</u> 1.5 <u>1.4</u>	2016 936 <u>884</u>	678 <u>281</u> 328

6.3 **ONE-SHOT VS. MULTI-SHOT GROWTH**

354 Since GEARnn employs OSG (One-Shot Growth) for growing networks, it begs the question if we 355 are missing anything if multiple growth steps (*m*-Shot Growth) were to be permitted, i.e., question 356 **Q2** from Section 1. To answer this question, we compare the clean and robust accuracies along with 357 training time and energy for different growth steps between GEARnn-1 and GEARnn-2 in Table 4. 358 All *m*-Shot Growth methods start with the same initial backbone f_0 (1.4% of full model size) and perform growth to reach f_2 (5% of full model size) using different growth ratios. All methods use 359 VGG-19 model and perform 80 epochs parametric training during the growth phase. The experiments 360 are done on CIFAR-10 data and the hardware measurements are taken from Jetson. 361

362 Table 4 indicates that OSG is comparable or better than the other m-Shot Growth 364 methods in all the metrics, thereby answering **Q2**. This 366 result can be attributed to 367 the lower training overhead 368 of growth stage in OSG 369 compared to the m-Shot 370 Growth methods. It should 371 be noted that as the growth 372 steps increase, the accura-373 cies go down and training

Table 4: Comparison of training complexities, clean and robust accuracies for different growth methods implemented using VGG-19 and CIFAR-10 on Jetson. 2-Phase approach and OSG provide the best solution for growing robust networks on the Edge.

Growth		GEAR	.nn-1		GEARnn-2					
Steps	$\mathcal{A}_{cln}(\%)$	$\mathcal{A}_{rob}(\%)$	$t_{\rm tr}~({\rm min})$	E (kJ)	$\mathcal{A}_{cln}(\%)$	$\mathcal{A}_{rob}(\%)$	$t_{\rm tr}~({\rm min})$	E (kJ)		
1	90.94	82.25	652	207	92.07	83.45	596	155		
2	90.01	81.92	640	191	91.94	83.34	593	157		
3	89.73	80.86	653	194	91.79	83.05	624	177		
4	89.90	81.08	845	223	91.65	82.75	645	173		

374 cost goes up, thus indicating that the optimal solution cannot be found by further increasing the 375 growth steps. Another comparison that is highlighted by Table 4 is the one between GEARnn-1 and GEARnn-2. For each growth step, GEARnn-2 is better than the corresponding GEARnn-1 376 solution on all the metrics. The numbers highlighted in red indicate the best solution across the 377 table. Thus, Table 4 clearly highlights that 2-Phase approach using One-Shot Growth is the best

324

325

333

347 348

341

342 343

349 350

378 combination to grow robust networks efficiently on the Edge. More comparisons between OSG and
 379 Multi-Shot growth are shown in Appendix B.2.
 380

7 ABLATION STUDY

In this section, we look at the generalization of GEARnn to other robust augmentations and then
 understand the robustness and efficiency breakdowns for GEARnn.

384 385

381

7.1 GENERALIZATION ACROSS ROBUST AUGMENTATION METHODS

The results thus far employed AugMix Hendrycks et al. (2019) tranforms (\mathcal{T}) to generate \mathcal{D}_{aug} for robust training. In this section, we see if the benefits of GEARnn are maintained across other augmentation transforms. Table 5 compares the implementation of PRIME Modas et al. (2022) augmentation across different methods. The accuracy and efficiency trend observed are similar to the results in Table 12. The important aspect to notice is the increase in training complexity gap ($\sim 2\times$) between GEARnn-1 and GEARnn-2. This is because OSG with PRIME is more expensive than OSG with AugMix.

393Table 5: Accuracy and Efficiency comparisons394for PRIME (\mathcal{D}_{aug}) augmentation implemented395for VGG-19 and CIFAR-10 on Quadro.

Table 6: Training time and energy breakdown for GEARnn on CIFAR-10 using VGG-19 on Quadro.

Method	$\mathcal{A}_{cln}(\%)$	$\mathcal{A}_{rob}~(\%)$	$t_{\rm tr}~({\rm min})$	E (kJ)
$\text{Small}\left(\mathcal{D}_{in}\right)$	92.69	70.57	31	241
Small (\mathcal{D}_{aug})	<u>91.30</u>	<u>87.01</u>	829	2550
GEARnn-Ĭ	88.37	83.18	458	1410
GEARnn-2	90.26	84.45	<u>234</u>	<u>856</u>

Ouantity	G	EARnn-1	l	GEARnn-2					
C ,	OSG-1	OSG-2	Total	OSG-1	OSG-2	ERA	Total		
training	38	48	86	5	10	38	53		
time (min)	44%	56%	100%	9%	19%	72%	100%		
energy	180	372	552	26	71	201	298		
(kJ)	33%	67%	100%	9%	24%	67%	100%		

397

399

7.2 EFFICIENCY AND ROBUSTNESS BREAKDOWN

404Table 6 shows the breakdown of energy and training405time for different stages of GEARnn-1 and GEARnn-4062. OSG-1 involves the training of backbone f_0 and407OSG-2 includes both the growth stage and training408of f_g . The key aspect to notice in Table 6 is the409small fraction of training cost required by OSG-1 and410OSG-2 in GEARnn-2 to provide a good initialization.

Table 7 shows the ablation studies of different components used in GEARnn-2 and compares it with a fixed network robust training. Firstly, we notice that OSG is more efficient than vanilla (fixed network) training, both in terms of training time and energy

Table 7: Impact of using OSG and ERA for CIFAR-100 and VGG-19 on Quadro.

Phase-	$1 (\mathcal{D}_{in})$	Phase-2	(\mathcal{D}_{aug})	Arch(%)	t _{tr} (min)	E (kJ)
vanilla	OSG	AugMix	ERA	• 400(7.0)	••••	- ()
\checkmark		1		38.72	18	161
	\checkmark			38.01	16	118
		\checkmark		46.50	62	385
			\checkmark	46.13	46	406
\checkmark		\checkmark		53.74	79	534
	\checkmark		\checkmark	54.31	64	515

while achieving comparable accuracy. Similar observations can be made for ERA over AugMix.
Performing 2-Phase approach by using either vanilla or OSG as initialization provides a significant
boost in robustness while incurring a minimal overhead in training cost. Thus the 2-Phase approach
is a clear winner over the 1-Phase approach, and in particular the combination of OSG and ERA
used for GEARnn-2 is optimal. More comparisons between AugMix and ERA on Jetson are shown
in Appendix B.3.

422 8 DISCUSSION

Until now we have looked at extensive empirical simulations that highlight the efficacy of GEARnn-2.
 In this section, we will look at the inner workings of this algorithm. Specifically, we will see what network topologies are generated when OSG designs compact networks, and also understand why clean data initialization benefits robust training.

427 428 8.1 IMPACT OF OSG ON NETWORK TOPOLOGY

In this section, we look at the growth topology patterns $(\{w_l\}_{l=1}^L)$ as a function of layer index *l*. Specifically, we investigate these patterns in the simple setting of OSG (\mathcal{D}_{in}) implemented on CIFAR-10 for ($\mathcal{E}_1, \mathcal{E}_2$) = (40, 40) and an initial backbone f_0 with $\{w_l\}_{l=1}^L = 45$. The bar plots represent the mean width ($\mathbb{E}[w_l]$) across four random seeds.



Figure 6: Average output channels vs. layer index for CIFAR-10 on Quadro is shown. Plot (a) looks at the impact of network architecture and highlights the non-uniform growth pattern in plain CNNs versus steady zigzag pattern in residual CNNs. Plots (b) and (c) indicate that modifying the number of growth epochs (\mathcal{E}_g) or performing 1-Phase robust growth does not affect the topology pattern much.

Backbone architecture: For plain CNNs like VGG-19 Simonyan & Zisserman (2014) - the initial layers have higher number of convolutional filters compared to final layers. This correlates with the observations seen in quantization Sakr & Shanbhag (2018) where the initial layers require higher precision compared to the final layers. However, in case of residual networks like ResNet-18, the pattern is largely invariant to network depth and is oscillating as shown in Fig. 6a. The invariance in depth can be attributed to the direct gradient flow facilitated by the shortcut connections making each residual block act independently of the depth. In each residual block, the macro-level pattern in plain CNNs is observed at a micro-level, i.e. initial layer has more output channels than the final layer.

Growth Epochs and Data: All the above experiments were performed for a single growth epoch ($\mathcal{E}_g = 1$) and on clean data. The effect of increasing \mathcal{E}_g to 50 and using ERA data for growth (GEARnn-1) is shown in Fig. 6b and Fig. 6c. The topology pattern in both cases remains roughly the same as OSG (\mathcal{D}_{in}) $\mathcal{E}_g = 1$.

460 8.2 RATIONALE FOR 2-PHASE APPROACH

461 In this section, we provide insights for the efficacy of GEARnn-2 and the 2-Phase approach. In 462 particular, we highlight why training or growth done on clean data provides a good initialization 463 for robust training. We look at the loss curves for the 1-Phase approaches (Small (AugMix), Small 464 (ERA), GEARnn-1) and the 2-Phase approach (GEARnn-2) in Fig. 7. The initial dip in GEARnn-2 465 loss function in Fig. 7a is due to the loss landscape being different for Phase-1 done on clean data 466 compared to Phase-2 done on augmented data. One can clearly see that GEARnn-2 achieves a lower 467 loss at a faster rate compared to the other 1-Phase approaches, thus justifying the importance of clean growth initialization. We also plot the filter normalized loss curves Li et al. (2018) in Fig. 7b to 468 observe the loss landscapes around the converged weights. GEARnn-2 finds the smallest minima 469 while also having a wide curve which enables better generalization Li et al. (2018). 470

471 The above explanation illustrates why GEARnn-2 has a good training and generalization performance. 472 However, in order to understand why initialization with *clean data* aids faster convergence of robust training, we look at the Fourier spectrums of the clean, augmented and corrupted images 473 474 in Fig. 8. Fig. 8a indicates that the clean images lie in the low frequency domain, while the corrupted samples occupy a wide range of frequencies (Figs. 8b & 8c). Crucially, the spectrum containing all 475 the augmentations (in AugMix) Fig. 8d and all the corruptions (in CIFAR-10-C) Fig. 8e is also in 476 the low-frequency domain, similar to the clean image spectrum Fig. 8a. This is unlike the scenario 477 of adversarial or Gaussian noise perturbations, which lie in the high-frequency domain Yin et al. 478 (2019) and hence may not benefit from clean data initialization. Thus, robust training for common 479 corruptions benefits from initialization with clean data. 480

481 482

443

444

445

446 447

9 LIMITATIONS AND BROADER IMPACTS

While our work has conclusively shown that a 2-Phase approach for growing robust networks is
 computationally efficient, a theoretical convergence analysis for this result is currently lacking. Such
 a result would help identify favorable initial conditions for robust training to achieve high accuracy in
 fewer epochs.



Figure 7: Loss comparisons for 2-Phase (GEARnn-2) and 1-Phase (rest) approaches for CIFAR-100 and VGG-19 on Quadro with 50 epochs of robust training at final model size. Fig. 7a highlights that GEARnn-2 loss converges to the minimum faster than other approaches. Fig. 7b shows the loss landscapes where GEARnn-2 achieves the smallest minima with a wide curve, thus aiding better generalization Li et al. (2018).



Figure 8: The Fourier spectrum of clean images (from CIFAR-10), their corresponding augmented (AugMix) and corrupted versions (CIFAR-10-C at severity 3) are shown. The augmentation and corruption spectrums (Figs. 8b, 8c, 8d & 8e) are obtained by taking Fourier Transform of the difference with the clean image (Eg: $\kappa(\mathbf{x}_{in}, 3) - \mathbf{x}_{in}$). Snow(Fig. 8b) and JPEG compression(Fig. 8c) corruptions are shown to highlight the range of possible frequencies in the corrupted spectrums. The similarity in the spectrums of clean (Fig. 8a), augmented (Fig. 8d) and all-corrupted (Fig. 8e) images highlights the importance of OSG initialization using clean data.

The impact of our work is broadly positive since it enables efficient robust training on Edge devices. We do not see any direct negative impact of our work.

523 10 CONCLUSION

We addressed the problem of growing robust networks efficiently on Edge devices. Specifically, we concluded that a 2-Phase approach with distinct clean growth and robust training phases is significantly more efficient than a 1-Phase approach which employs augmented data for growth. We encapsulated this result into the GEARnn algorithm and experimentally demonstrated its benefits on a real-life Edge device. An interesting and non-trivial extension of our work would be to use unlabeled data for growing efficient and robust networks. Another extension would be to design robust networks for complex tasks such as object detection on highly resource-constrained Edge platforms.

532 11 REPRODUCIBILITY STATEMENT

We list all the experimental setup details in Section 5 and Appendix A. We use fixed seeds during the
 simulation runs so that our results can be reproduced. We will make the code public along with the
 software versions if the paper is accepted so that the community can use and reproduce our results.

540 REFERENCES

- Javad Zolfaghari Bengar, Joost van de Weijer, Bartlomiej Twardowski, and Bogdan Raducanu.
 Reducing label effort: Self-supervised meets active learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 1631–1639, 2021.
- Raffaello Bonghi. Jetson-stats. https://pypi.org/project/jetson-stats/.
- Dan A Calian, Florian Stimberg, Olivia Wiles, Sylvestre-Alvise Rebuffi, Andras Gyorgy, Timothy
 Mann, and Sven Gowal. Defending against image corruptions through adversarial augmentations.
 arXiv preprint arXiv:2104.01086, 2021.
- Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine
 Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar Erlingsson, et al. Extracting training data
 from large language models. In *30th USENIX Security Symposium (USENIX Security 21)*, pp. 2633–2650, 2021.
- Tianqi Chen, Ian Goodfellow, and Jonathon Shlens. Net2net: Accelerating learning via knowledge transfer. *arXiv preprint arXiv:1511.05641*, 2015.
- Francesco Croce, Maksym Andriushchenko, Vikash Sehwag, Edoardo Debenedetti, Nicolas Flammarion, Mung Chiang, Prateek Mittal, and Matthias Hein. Robustbench: a standardized adversarial robustness benchmark. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2021. URL https://openreview.net/forum?id= SSKZPJCt7B.
- James Diffenderfer, Brian Bartoldson, Shreya Chaganti, Jize Zhang, and Bhavya Kailkhura. A
 winning hand: Compressing deep networks can improve out-of-distribution robustness. *Advances in neural information processing systems*, 34:664–676, 2021.
- 565 Utku Evci, Bart van Merrienboer, Thomas Unterthiner, Max Vladymyrov, and Fabian Pedregosa.
 566 Gradmax: Growing neural networks using gradient information. *arXiv preprint arXiv:2201.05125*, 2022.
- Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. *arXiv preprint arXiv:1803.03635*, 2018.
- Kenji Fukumizu and Shun-ichi Amari. Local minima and plateaus in hierarchical structures of
 multilayer perceptrons. *Neural networks*, 13(3):317–327, 2000.
- Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial
 examples. *arXiv preprint arXiv:1412.6572*, 2014.
- Song Han, Jeff Pool, John Tran, and William Dally. Learning both weights and connections for
 efficient neural network. *Advances in neural information processing systems*, 28, 2015.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- 582 Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. *arXiv preprint arXiv:1903.12261*, 2019.
- Dan Hendrycks, Norman Mu, Ekin D Cubuk, Barret Zoph, Justin Gilmer, and Balaji Lakshmi narayanan. Augmix: A simple data processing method to improve robustness and uncertainty.
 arXiv preprint arXiv:1912.02781, 2019.
- Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul
 Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, et al. The many faces of robustness: A critical analysis of out-of-distribution generalization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 8340–8349, 2021.
- Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand,
 Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for
 mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017.

594 595 596	Itay Hubara, Matthieu Courbariaux, Daniel Soudry, Ran El-Yaniv, and Yoshua Bengio. Binarized neural networks. <i>Advances in neural information processing systems</i> , 29, 2016.
597 598 599	Ching-Yi Hung, Cheng-Hao Tu, Cheng-En Wu, Chien-Hung Chen, Yi-Ming Chan, and Chu-Song Chen. Compacting, picking and growing for unforgetting continual learning. <i>Advances in Neural Information Processing Systems</i> , 32, 2019.
600 601	Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
602 603	Ya Le and Xuan Yang. Tiny imagenet visual recognition challenge. CS 231N, 7(7):3, 2015.
604 605	Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. Pruning filters for efficient convnets. <i>arXiv preprint arXiv:1608.08710</i> , 2016.
606 607 608	Hao Li, Zheng Xu, Gavin Taylor, Christoph Studer, and Tom Goldstein. Visualizing the loss landscape of neural nets. <i>Advances in neural information processing systems</i> , 31, 2018.
609 610 611	Chenxi Liu, Barret Zoph, Maxim Neumann, Jonathon Shlens, Wei Hua, Li-Jia Li, Li Fei-Fei, Alan Yuille, Jonathan Huang, and Kevin Murphy. Progressive neural architecture search. In <i>Proceedings</i> of the European conference on computer vision (ECCV), pp. 19–34, 2018.
612 613 614 615	Dongzhu Liu, Guangxu Zhu, Jun Zhang, and Kaibin Huang. Wireless data acquisition for edge learning: Importance-aware retransmission. In 2019 IEEE 20th International Workshop on Signal Processing Advances in Wireless Communications (SPAWC), pp. 1–5. IEEE, 2019.
616 617 618	Eric Mintun, Alexander Kirillov, and Saining Xie. On interaction between augmentations and corruptions in natural corruption robustness. <i>Advances in Neural Information Processing Systems</i> , 34:3571–3583, 2021.
619 620 621 622	Apostolos Modas, Rahul Rade, Guillermo Ortiz-Jiménez, Seyed-Mohsen Moosavi-Dezfooli, and Pascal Frossard. Prime: A few primitives can boost robustness to common corruptions. In <i>European Conference on Computer Vision</i> , pp. 623–640. Springer, 2022.
623 624 625 626	Santosh Kumar Nukavarapu and Tamer Nadeem. Securing edge-based iot networks with semi- supervised gans. In 2021 IEEE International Conference on Pervasive Computing and Com- munications Workshops and other Affiliated Events (PerCom Workshops), pp. 579–584. IEEE, 2021.
627 628 629	Corporation NVIDIA. Nvidia jetson xavier nx for embedded and edge systems. https: //www.nvidia.com/en-sg/autonomous-machines/embedded-systems/ jetson-xavier-nx/, a.
630 631 632	Corporation NVIDIA. System management interface smi. https://developer.nvidia.com/nvidia-system-management-interface, b.
633 634 635 636 637	Junting Pan, Adrian Bulat, Fuwen Tan, Xiatian Zhu, Lukasz Dudziak, Hongsheng Li, Georgios Tzimiropoulos, and Brais Martinez. Edgevits: Competing light-weight cnns on mobile devices with vision transformers. In <i>European Conference on Computer Vision</i> , pp. 294–311. Springer, 2022.
638 639 640 641	Ruiyang Qin, Jun Xia, Zhenge Jia, Meng Jiang, Ahmed Abbasi, Peipei Zhou, Jingtong Hu, and Yiyu Shi. Enabling on-device large language model personalization with self-supervised data selection and synthesis. In <i>Proceedings of the 61st ACM/IEEE Design Automation Conference</i> , pp. 1–6, 2024.
642 643 644	Mohammad Rastegari, Vicente Ordonez, Joseph Redmon, and Ali Farhadi. Xnor-net: Imagenet classification using binary convolutional neural networks. In <i>European conference on computer vision</i> , pp. 525–542. Springer, 2016.
646 647	Evgenia Rusak, Lukas Schott, Roland Zimmermann, Julian Bitterwolfb, Oliver Bringmann, Matthias Bethge, and Wieland Brendel. Increasing the robustness of dnns against im-age corruptions by playing the game of noise. 2020.

648 649 650	Andrei A Rusu, Neil C Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick, Koray Kavukcuoglu, Razvan Pascanu, and Raia Hadsell. Progressive neural networks. <i>arXiv preprint arXiv:1606.04671</i> , 2016.
651 652 653 654	Charbel Sakr and Naresh Shanbhag. An analytical method to determine minimum per-layer precision of deep neural networks. In 2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 1090–1094. IEEE, 2018.
655 656	Lianlei Shan, Weiqiang Wang, Ke Lv, and Bin Luo. Edge-guided and class-balanced active learning for semantic segmentation of aerial images. <i>arXiv preprint arXiv:2405.18078</i> , 2024.
657 658 659	Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. <i>arXiv preprint arXiv:1409.1556</i> , 2014.
660 661 662	Jiachen Sun, Akshay Mehra, Bhavya Kailkhura, Pin-Yu Chen, Dan Hendrycks, Jihun Hamm, and Z Morley Mao. Certified adversarial defenses meet out-of-distribution corruptions: Benchmarking robustness and simple baselines. <i>arXiv preprint arXiv:2112.00659</i> , 2021.
663 664 665	Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. <i>arXiv preprint arXiv:1312.6199</i> , 2013.
666 667 668	Hanlin Tang, Xiangru Lian, Ming Yan, Ce Zhang, and Ji Liu. (2): Decentralized training over decentralized data. In <i>International Conference on Machine Learning</i> , pp. 4848–4856. PMLR, 2018.
669 670 671 672	Wei Wen, Feng Yan, Yiran Chen, and Hai Li. Autogrow: Automatic layer growing in deep con- volutional networks. In <i>Proceedings of the 26th ACM SIGKDD International Conference on</i> <i>Knowledge Discovery & Data Mining</i> , pp. 833–841, 2020.
673 674	Lemeng Wu, Dilin Wang, and Qiang Liu. Splitting steepest descent for growing neural architectures. <i>Advances in neural information processing systems</i> , 32, 2019.
675 676 677 678	Lemeng Wu, Bo Liu, Peter Stone, and Qiang Liu. Firefly neural architecture descent: a general approach for growing neural networks. <i>Advances in neural information processing systems</i> , 33: 22373–22383, 2020.
679 680 681	Dong Yin, Raphael Gontijo Lopes, Jon Shlens, Ekin Dogus Cubuk, and Justin Gilmer. A fourier perspective on model robustness in computer vision. <i>Advances in Neural Information Processing Systems</i> , 32, 2019.
682 683	Xin Yuan, Pedro Savarese, and Michael Maire. Growing efficient deep networks by structured continuous sparsification. <i>arXiv preprint arXiv:2007.15353</i> , 2020.
685 686 687 688	Xin Yuan, Pedro Henrique Pamplona Savarese, and Michael Maire. Accelerated training via in- crementally growing neural networks using variance transfer and learning rate adaptation. In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> , 2023. URL https: //openreview.net/forum?id=Hla7bVVnPK.
689 690 691	Long Zhao, Ting Liu, Xi Peng, and Dimitris Metaxas. Maximum-entropy adversarial data augmen- tation for improved generalization and robustness. <i>Advances in Neural Information Processing</i> <i>Systems</i> , 33:14435–14447, 2020.
692 693 694 695	Barret Zoph, Vijay Vasudevan, Jonathon Shlens, and Quoc V Le. Learning transferable architectures for scalable image recognition. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pp. 8697–8710, 2018.
696 697	
698 699 700	

702 APPENDIX / SUPPLEMENTAL MATERIAL

704 A TRAINING SETUP

705 **Hyperparameters:** The setup for growth and robust augmentation follows closely with what is 706 described in Firefly Wu et al. (2020) and AugMix Hendrycks et al. (2019), respectively. The parametric training is done for 160 epochs using a batch-size of 128 and an initial learning rate of 708 0.1. The learning rate scheduler decays by 0.1 at half and three-fourths of the total number of epochs. 709 We use the Swish loss function for MobileNet-V1 as used in Wu et al. (2020), while employing 710 ReLU for the other two networks. Instead of using three fully-connected layers at the end of VGG-19, we use only one as done in Wu et al. (2020). Stochastic Gradient Descent (SGD) optimizer is used 711 with momentum 0.9 and weight decay 10^{-4} . As for the standard growth process, we use a Root 712 Mean Square Propagation (RMSprop) optimizer with momentum 0.9, alpha 0.1 and initial learning 713 rate of 9×10^{-5} . The number of workers is chosen as 4. For ERA, (W, D, J) = (1, 3, 4) is picked. 714 The augmentation transforms \mathcal{T} are same as that of AugMix Hendrycks et al. (2019) for all the 715 results except Table 5, where we pick the transforms from PRIME Modas et al. (2022). As specified 716 in AugMix, we also do not use any augmentations which are directly present in the corrupted test 717 dataset. 718

In case of OSG, the initial backbone f_0 is chosen as a network with $w_l = 45$ for all $l = \{1, ..., L\}$ and is thus extremely small. The number of randomly initialized neurons at each growth stage is 70. We ensure that \mathcal{E}_2 of GEARnn-1 and \mathcal{E}_r of GEARnn-2 are same for a fair comparison. \mathcal{E}_g is chosen as 1 based on Firefly Wu et al. (2020). The transforms used in AugMix are autocontrast, equalize, posterize, rotate, solarize, shear_x, shear_y, translate_x, translate_y.

Jetson Training: The two changes to the GEARnn algorithm when implementing on NVIDIA Jetson Xavier are - one we use j = 3 instead of j = 4, and two, we allow only 40 randomly initialized new neurons per layer in the growth step (as compared to 70 in Wu et al. (2020)). These measures are taken to stay within the memory constraints of the Edge device. We also reduce the batch size (and learning rate) appropriately in case the above measures are insufficient.

729 730 731

732

733

B ABLATION STUDIES

B.1 DIAGNOSTICS OF ROBUST AUGMENTATION METHODS

734 In this section, we investigate which aspects of the ro-735 bust augmentation framework described in Section 4.2 736 contribute most to the robustness while being training 737 efficient. Table 8 shows different modifications of 738 the stochastic chains obtained by varying (W, D, J)739 values. It can be observed that the basic version with 740 (W, D, J) = (1, 1, 0) (uses only standard cross en-741 tropy loss with the label and augmented data as input) has the least training time, but suffers a significant 742 drop in A_{rob} compared to standard AugMix. Cru-743 cially, we note that increase in D and J has more 744 impact on robustness at a lesser training cost com-745 pared to W. For ERA, we pick the modification with 746 (W, D, J) = (1, 3, 4) as it provides the highest ro-747 bustness while simultaneously reducing training time 748 over AugMix.

749 750

751

Experiment	W	D	J	$\mathcal{A}_{rob}(\%)$	$t_{\rm tr}({\rm min})$
Basic	1	1	0	77.74	10
+ width	3	1	0	78.51	16
+ depth	1	3	0	80.31	12
+ JSD-3	1	1	3	82.43	20
+ width + depth	3	3	0	80.47	21
+ width + JSD-3	3	1	3	82.27	32
+ depth + JSD-3	1	3	3	83.67	22
+ depth + JSD-2	1	3	2	82.41	13
+ depth + JSD-4	1	3	4	84.10	29
AugMix Hendrycks et al. (2019)	3	3	3	84.05	41

Table 8: Impact of training AugMixvariants on the robust accuracy and training time. Network f_2 from OSG is used as the starting network and $\mathcal{E}_r = 40$. All the methods are implemented for CIFAR-10 and 5% VGG-19 network on Quadro. W, D, J represent the width, depth and consistency samples used in the stochastic chains.

B.2 OSG VERSUS MULTI-SHOT GROWTH COMPARISONS

In this section, we first look at clean data growth comparisons on Jetson in Table 9. Then we look at robust data growth comparisons on Quadro in Table 14. When comparing various growth methods on clean data in Table 9, we also include the Small (\mathcal{D}_{in}) results to highlight the efficiency benefits of growth. We can see that OSG has comparable or better training efficiency than all the methods including Small (\mathcal{D}_{in}) . In case of clean accuracy, we observe that OSG has the highest among growth methods while being slightly lower than Small (\mathcal{D}_{in}). Looking at Table 14, we see that for both datasets OSG again provides comparable or best solution among all the growth methods across all metrics. Thus our choice of OSG over other Multi-Shot Growth methods is justified.

Table 9: Comparison of training complexities
and clean accuracy for different growth methods implemented using VGG-19 and clean
CIFAR-10 data on Jetson.

Table 10: OSG versus Multi-Shot Growth using ERA data, i.e. GEARnn-1 with Multi-Shot Growth. Results are shown for VGG-19 on Quadro.

Growth Steps	$ \mathcal{A}_{cln}(\%) $	$t_{\rm tr}~({\rm min})$	E (kJ)
Small (\mathcal{D}_{in})	91.96	267	73
1	90.80	210	51
2	90.49	209	53
3	90.31	231	59
4	90.08	275	65

Growth	C	IFAR-10	CIFAR-100						
Steps	\mathcal{A}_{rob} (%)	$t_{\rm tr}~({\rm min})$	$E\left(\mathrm{kJ}\right)$	$\left \mathcal{A}_{\mathrm{rob}} \left(\% \right) \right $	$t_{\rm tr}~({\rm min})$	E (kJ)			
1	82.86	75	449	52.68	84	426			
2	82.31	81	329	51.34	100	554			
3	81.94	82	506	50.52	103	517			
4	81.81	86	389	50.45	102	590			

B.3 BENEFITS OF ERA ON JETSON

In this section we will look at the benefits of using ERA over AugMix. Previously, we had looked at this comparison on Quadro using CIFAR-100 and VGG-19 in Table 7. Here, we will look at these results for CIFAR-10 and VGG-19 on Jetson when training a fixed-size Small (\mathcal{D}_{aug}) network for 160 epochs. Table 11 indicates that ERA is better than AugMix on all the metrics. Thus our choice of ERA over AugMix is justified.

Table 11: Comparison of ERA versus Aug-Mix on Jetson for VGG-19 and CIFAR-10 when trained for 160 epochs

Method	$ \mathcal{A}_{cln}(\%) $	$t_{\rm tr}~({\rm min})$	$E (\mathrm{kJ})$	
Small (AugMix)	93.36	85.73	1543	522
Small (ERA)	93.42	85.74	1542	486

B.4 GAUSSIAN AUGMENTATION

Table 12: Gaussian Augmentation comparison between the baselines and GEARnn on VGG-19 for CIFAR-10, CIFAR-100 on Quadro.

Architecture				CIFAR-1	0		CIFAR-100				
(full model size)	Method	Size (%)	Accu Accu Accu Mathematical Accu	$\mathcal{A}_{\rm rob}(\%)$	$\begin{vmatrix} Train \\ t_{tr}(min) \end{vmatrix}$	ing E(kJ)	Size (%)	$\begin{vmatrix} Accu \\ \mathcal{A}_{cln}(\%) \end{vmatrix}$	$\mathcal{A}_{rob}(\%)$	$\begin{vmatrix} \text{Train} \\ t_{tr}(\min) \end{vmatrix}$	ning E(kJ)
	Small (\mathcal{D}_{in})	5	92.69	70.57↓	31	241	9	68.07	41.24	38	335
VGG-19 (20M)	Small (\mathcal{D}_{aug}) GEARnn-1 GEARnn-2) 5 5 5	86.93 84.07 86.65	75.90 73.81 76.21	50 28 27	269 <u>93</u> 114	9 9 9	57.82 51.16 56.93	45.13 40.30 45.55	$\begin{vmatrix} 46 \\ \underline{30} \\ \overline{31} \end{vmatrix}$	298 <u>114</u> 203

B.5 RESULTS ON NVIDIA JETSON ORIN NANO

Table 13: Comparison of between the baselines and GEARnn for CIFAR-10 and CIFAR-100 using VGG-19 on Jetson Orin Nano.

		CIFAR-10							CIFAR-100					
Network	Method	$\begin{vmatrix} Accu \\ A_{cln}(\%) \end{vmatrix}$	$\mathcal{A}_{rob}(\%)$	Infe Size%	$\left \begin{array}{c} t_{inf}(ms) \end{array} \right $	Train t _{tr} (min)	ing E(kJ)	Acc $\mathcal{A}_{cln}(\%)$	$\mathcal{A}_{rob}(\%)$	Infe Size%	erence $t_{inf}(ms)$	$\begin{vmatrix} \text{Train} \\ t_{tr}(\min) \end{vmatrix}$	ning E(kJ)	
	Small (\mathcal{D}_{in})	92.84	70.03	5	0.3	149	49	67.07	40.29↓	9	0.3	175	59	
VGG-19	$\begin{array}{l} \text{Small} \left(\mathcal{D}_{aug} \right) \\ \text{GEARnn-1} \\ \text{GEARnn-2} \end{array}$	93.00 90.72 92.23	85.19 82.13 83.29	5 5 5	$\begin{array}{c c} \underline{0.3} \\ \underline{0.3} \\ \underline{0.3} \end{array}$	411 196 <u>174</u>	173 72 <u>57</u>	69.23 63.02 67.14	<u>55.72</u> 50.16 53.56	9 9 9	<u>0.3</u> 0.4 <u>0.3</u>	492 291 231	206 104 <u>82</u>	

C ACCURACY-ROBUSTNESS-EFFICIENCY TRADE-OFFS

C.1 TRAINING TIME VERSUS ACCURACIES

In Section 6 and Fig. 5 we observed that GEARnn-2 can achieve high robustness even when the robust training epochs are low. This is due to better initialization provided by OSG. We show the same results ablated for both VGG-19 and MobileNet-V1 for CIFAR-10 and CIFAR-100 in Fig. 9.



Figure 9: Plots (a)-(c) are implemented for VGG-19/CIFAR-10, (d)-(f) are for MobileNet-V1/CIFAR-10, (g)-(i) are for VGG-19/CIFAR-100, and (j)-(l) are for MobileNet-V1/CIFAR-100 on Quadro.
First two plots of each row indicates the robust accuracy as a function of epochs and training time respectively. The last plot in each row shows the clean accuracy as a function of training time. GEARnn-2 clearly achieves the best clean and robust accuracy at the same training cost.

864 C.2 MODEL SIZE VERSUS ACCURACIES 865

Fig. 10a and Fig. 10b show the results of L1-Unstructured pruning performed on GEARnn-2 final network. The global sparsity is varied from 10% to 90% in steps of 20%. Fig. 10c shows the impact of varying the growth ratio γ in GEARnn-2's OSG.



Figure 10: Plots (a) & (b) show the impact of L1-unstructured pruning done on the final network f_{2n}^* obtained from GEARnn-2 (1 million params). Plot (a) represents the sparsity-controlled pruning and Plot (b) shows the corresponding points plotted in parameter-space. Plot (c) indicates the impact of parameters on robust accuracy of GEARnn-2 when the growth ratio γ of OSG is varied. All experiments are conducted on CIFAR-10 data using VGG-19 network.

885 **PRIOR WORKS** D 886

FIREFLY **D**.1 887

888 In this section, we explain how the splitting and growing new neurons in Firefly Wu et al. (2020) 889 (and our growth technique \mathcal{G}) is implemented. We explain it in terms of fully-connected layers and neurons, but this can be easily extended to CNNs. Consider a multi-layered perception with two 890 neurons in the hidden layer as shown in Figure 1 of Wu et al. (2020). If x is the input to the neurons, 891 θ_i and 1 are the weights, σ is the activation function, then we can write the input to the final layer as 892 $\sigma(x, \theta_i)$. In case of splitting growth, we add a new incoming weight by perturbing the existing weight 893 $(\sigma(x, \theta_i - \varepsilon_i \delta_i))$ and adding the new perturbed weight $(\sigma(x, \theta_i + \varepsilon_i \delta_i))$. When adding a random 894 new-grown weight, we add a randomly initialized weight δ_i such that the input to the final layer is 895 $\varepsilon_i \sigma(x, \delta_i)$. We can write the function as follows: 896

897 898

899 900

901

902

903

909

866

867

868

870

871

872

873

874

875

876

877 878

879

880

881

882

883

884

$$f_{\boldsymbol{\varepsilon},\boldsymbol{\delta}} = \sum_{i=1}^{m} \frac{1}{2} \left(\sigma(x,\theta_i - \varepsilon_i \delta_i) + \sigma(x,\theta_i + \varepsilon_i \delta_i) \right) + \sum_{i=m+1}^{m+m'} \varepsilon_i \sigma(x,\delta_i)$$
$$\min_{\boldsymbol{\varepsilon},\boldsymbol{\delta}} \{ \mathcal{L}(f_{\boldsymbol{\varepsilon},\boldsymbol{\delta}}) \text{ s.t. } ||\boldsymbol{\varepsilon}||_0 \le \gamma \mathcal{C}(f_1), \, ||\boldsymbol{\varepsilon}||_{\infty} \le \epsilon, \, ||\boldsymbol{\delta}||_{2,\infty} \le 1 \}$$

where m denotes the number of split neurons and m' denotes the number of newly grown neurons. Solving the above minimization problem (denoted as \mathcal{G}) provides us the grown network.

D.2 AUGMIX 904

905 The working of AugMix Hendrycks et al. (2019) is similar to that of ERA. However AugMix uses 906 parallel concurrent transforms which makes it more inefficient than ERA. Below equations indicate 907 the working of AugMix. The notation is same as ERA and W denotes the width of the block of 908 transforms.

$$A_j^w(\mathbf{x}) = a_1 \circ a_2 \circ ... \circ a_{d_j}(\mathbf{x})$$

- $A_j(\mathbf{x}) = \sum_{w=1}^W \alpha_w A_j^w(\mathbf{x})$ $\mathbf{x}_j^{\text{aug}} = p\mathbf{x} + (1-p)A_j(\mathbf{x})$ 912
- 913

- $\mathcal{R}((\mathbf{x},y)|\mathcal{T}) = (\{\mathbf{x}_1^{\mathrm{aug}}, ..., \mathbf{x}_{J-1}^{\mathrm{aug}}, \mathbf{x}\}, y) \implies \mathcal{D}_{\mathrm{aug}} := \mathcal{R}(\mathcal{D}_{\mathrm{in}}|\mathcal{T})$ 915
- 916 where $a_i \sim \text{Unif}(\mathcal{T}), p \sim \beta(1,1), d_i \sim \text{Unif}(\{1,...,D\}), j \in \{1,...,J-1\}, \alpha \sim \text{Dirichlet}(W)$ 917 (3)

D.3 COMPARISON WITH PRIOR GROWTH WORKS

Table 14: Comparing GEARnn with SOTA growth methods

Growth		CIFA	R-10			CIFA		
Method	$\left \mathcal{A}_{cln} \left(\% \right) \right.$	$\mathcal{A}_{rob}~(\%)$	Size (M)	$t_{\rm tr}~({\rm min})$	$ \mathcal{A}_{cln}(\%) $	$\mathcal{A}_{rob}~(\%)$	Size (M)	$t_{\rm tr}~({\rm min})$
Splitting (Wu et al. (2019))	93.43	70.91	1.11	145	70.78	42.20	1.81	194
Firefly (Wu et al. (2020))	93.59	73.00	1.38	169	69.35	41.90	2.39	204
GEARnn-2 (ours)	92.18	83.77	1.06	53	68.44	54.31	1.81	65

Here we compare GEARnn-2 with the direct implementations of state-of-the-art growth methods. For implementation purposes the Splitting method has 1 randomly initialized neuron. The networks from Splitting and Firefly with the closest number of parameters to GEARnn-2 are picked for comparison. We find that GEARnn-2 has better robustness, training and inference efficiency compared to the SOTA methods.

E DISCUSSION (CONTINUED)

937 938 939

940

918

919 920 921

930

931

932

933

934 935 936

E.1 MOTIVATION FOR ON-DEVICE EDGE TRAINING

GEARnn performs isolated on-device Edge training. However, there are two other scenarios (fine tuning and federated learning) which can be used interchangeably with GEARnn depending on the
 application and constraints. But all three methods are important and none is a replacement for the
 other. Below we highlight the scenarios where on-device Edge training is necessary:

945 Fine-tuning: There are scenarios where existing pre-trained networks are not useful for the task at 946 hand due to lack overlap in data domains (eg: sensitive medical data, geographical data from other planets). Apart from this, fine-tuning of large models can lead to over-fitting when limited data is 947 available. This is countered by our method since growth involves training extremely small models 948 that gradually increase in size (the problem of over-parametrization is avoided and thus over-fitting). 949 Along with that, our growth-based approach allows us to design custom parameter-efficient models 950 based on the dataset (Section 8.1). This is not possible in existing methods since they either fit large 951 models on the Edge, or compress the existing models on the Cloud without the local data available 952 at the Edge. Lastly, even if robust fine-tuning is used instead of GEARnn, our method provides a 953 lesson to fine-tune on clean data before moving to augmented data (2-Phase approach) for efficient 954 fine-tuning. 955

Federated Learning: It was proposed to preserve privacy by transmitting the weights to the Cloud instead of the data directly. However, recent works in security have shown that training data can be extracted from the weights of these models Carlini et al. (2021) thus putting the privacy of Federated Learning in jeopardy. Practical Federated Learning also suffers from convergence issues due to non-IID data Tang et al. (2018). Lastly, there are several medical or defense applications where the user does not want to share any data or weights and would prefer on-device training.

With regard to the concern about lack of labeled data available on a single Edge device, there are several ways such as active learning Bengar et al. (2021); Shan et al. (2024), sel-supervision Qin et al. (2024), semi-supervision Nukavarapu & Nadeem (2021) and quality sampling Liu et al. (2019) to convert the abundant unlabeled data from the Edge device sensor into trainable data.

966 967

968

E.2 IMPLEMENTATION FOR TRANSFORMERS

For Edge devices (which is our focus) with limited compute, parameters and data, CNNs work as well
as transformers Pan et al. (2022). Hence, we focus our attention on growing CNNs for the Edge and
demonstrate our results on VGG-19, MobileNet-V1 and ResNet-18. Robust growth on transformers for other applications is a good direction for future work.

972 E.3 CHOICE OF NETWORK SIZE

974 Our choice of network size is determined by the growth ratio γ . The growth ratio is chosen such 975 that the model lies within the device memory constraints while achieving the desired accuracy. 976 Analytically estimating the memory consumption during training is challenging due to the dynamic 977 computation graphs, gradient calculations, data movement and hardware optimizations implemented. 978 Hence we resort to an empirical thumb rule - to ensure the nominal batch size used for this task 979 (128 in this case) is maintained at the Edge without any reduction. The growth ratio γ is chosen 980 accordingly for each network.