

# 000 TOWARDS ROBUST OUT-OF-DISTRIBUTION GENERALIZATION 001 FOR DEEP NEURAL NETWORKS WITH TAI- 002 LORED DATA REGULARIZATION 003 004

005 **Anonymous authors**  
 006  
 007 Paper under double-blind review  
 008  
 009  
 010  
 011

## ABSTRACT

012 Out-of-Distribution (OOD) generalization remains both a fundamental challenge  
 013 and an often-overlooked aspect of modern machine learning—especially in the  
 014 context of Deep Neural Networks (DNNs), which are highly expressive yet prone  
 015 to overfitting under distributional stress. Classical learning theory highlights the  
 016 role of regularization in managing the bias-variance trade-off—particularly im-  
 017 portant for complex models with higher VC dimension. In this work, we explore  
 018 **stochastic data regularization techniques**—such as random transformations and  
 019 noise injection—applied not only as isolated strategies but also organized through  
 020 a Scheduling Policy framework using a Curriculum Learning-based approach.  
 021 By progressively increasing input difficulty during training, the scheduling aligns  
 022 model capacity with task complexity, promoting more **robust generalization**. We  
 023 also propose a novel statistical procedure to assess the consistency of performance  
 024 estimates across cross-validation folds, mitigating miscoverage effects in confi-  
 025 dence interval estimation. Altogether, our findings highlight the importance of a  
 026 **tailored data regularization, where the selection, combination, and schedul-**  
 027 **ing** of perturbations become key to achieving OOD robustness in DNNs.  
 028

## 029 1 INTRODUCTION

030 Robust generalization is the ability of models to maintain reliable performance under distribution  
 031 shifts—when test data deviate from the training distribution. It remains a significant challenge  
 032 for Deep Neural Networks (DNNs), which are highly susceptible to overfitting under distribution  
 033 shifts (Li et al., 2022; Hendrycks et al., 2021). Within the broader landscape of robust training  
 034 strategies, regularization techniques are commonly used to counteract this issue—but they often  
 035 fall short or even cause over-regularization, degrading model performance (Lin et al., 2024). These  
 036 limitations highlight the need for carefully tailored regularizers (Choi & Kim, 2024; Srivastava et al.,  
 037 2014) whose effectiveness depends heavily on both the task (e.g., classification) and the data domain  
 038 (e.g., vision, language). This motivates the need for domain-aware and task-sensitive regularization  
 039 approaches.  
 040

041 In this work, we adopt the OOD definitions of Farquhar & Gal (2022), focusing on the **transformed-**  
 042 **OOD** setting, where label-preserving corruptions are applied to in-distribution inputs. To assess the  
 043 impact of such shifts, we compute statistical distances—such as KL divergence—between clean and  
 044 corrupted latent representations, capturing the extent of deviation and the effect of regularization.

045 Building on this perspective, we explore whether stochastic data regularization—via random trans-  
 046 formations (Cubuk et al., 2020; Hendrycks et al., 2019) and input noise (Bishop, 1995; Yuan et al.,  
 047 2025; Filho et al., 2023)—can act as an implicit regularizer when applied dynamically. We show  
 048 that organizing these perturbations as a curriculum (Bengio et al., 2009; Lu & Lam, 2023), gradu-  
 049 ally increasing their intensity, is a promising yet underexplored strategy in computer vision (Choi &  
 050 Kim, 2024)—especially effective for compact models aiming at robust generalization.

051 A modular framework that systematizes data regularization is presented through three components:  
 052 a *Selection Policy* (e.g., choosing between noise types or augmentation pipelines), a *Combination*  
 053 *Policy* (e.g., composing augmentations and noise), and a *Scheduling Policy* (Curriculum Learning-  
 based approach).

054 This paper explores a core question in robust model design: **How can data regularization be dy-**  
 055 **namically adapted to a model’s capacity to improve robustness while mitigating overfitting**  
 056 **and underfitting, thereby enhancing out-of-distribution performance?** We hypothesize that  
 057 stochastic data regularization—whether applied uniformly or progressively—can drive consistent  
 058 gains in robustness across domains. When organized as a curriculum, aligning perturbation strength  
 059 with model maturity, such strategies can enhance learning stability. Moreover, even unstructured  
 060 randomness in augmentations and noise appears effective in reducing overfitting and promoting  
 061 generalization, particularly when carefully tuned to avoid early over-perturbation.

062 To evaluate our models, we assess both average performance and the reliability of performance es-  
 063 timates. We introduce a miscoverage-based analysis across cross-validation folds, inspired by Bates  
 064 et al. (2023), to quantify how well confidence intervals reflect true variability. Our findings indicate  
 065 that stronger data regularization reduces miscoverage—particularly in shallow architectures—by ad-  
 066 dressing the bias–variance trade-off and promoting more stable out-of-distribution generalization.

## 068 2 RELATED WORKS

070 **Out-of-Distribution Categorization** The term out-of-distribution is often used ambiguously in  
 071 the literature, leading to inconsistent interpretations and methodological practices. To clarify this,  
 072 recent work (Farquhar & Gal, 2022) categorizes the different distributions into four types: trans-  
 073 formed, related, complementary, and synthetic.

074 **Diversity in Data Regularization for Robust Learning** Recent work has shown that data regu-  
 075 larization plays a key role in improving both robustness and generalization (Li & Spratling, 2023).  
 076 However, simple transformations are often insufficient under distribution shifts. Increasing the diver-  
 077 sity of augmentations promotes better model performance. Diverse transformations help the model  
 078 to generalize to unseen data and improve stability under adversarial conditions. This highlights the  
 079 importance of carefully designed augmentation pipelines for robust learning.

080 **Miscoverage in Cross-Validation-Based Estimates** Standard  $K$ -fold cross-validation (CV) often  
 081 underestimates variance, resulting in confidence intervals with lower-than-nominal coverage (Ben-  
 082 gio & Grandvalet, 2004). This miscoverage is especially pronounced when folds are not indepen-  
 083 dent, as data points contribute to both training and evaluation. More recently, Bates et al. (2023)  
 084 showed that even in modern settings, such intervals can severely misrepresent model uncertainty.  
 085 They observed that stronger regularization mitigates this effect by providing CV estimator with  
 086 fresher data across folds. In our work, we capitalize on this known limitation by using leave-fold-  
 087 out replications to directly assess model consistency. This not only exposes the weaknesses of  
 088 standard cross-validation but also highlights the robustness gains achieved through our proposed  
 089 data regularization strategies.

090 **Curriculum Learning** Curriculum Learning (CL) is a training paradigm inspired by the human  
 091 learning process, where models are exposed to increasingly difficult examples (Bengio et al., 2009).  
 092 It has shown benefits for convergence and generalization across domains. A recent causal analy-  
 093 sis (Li et al., 2024) highlights that CL is more effective when early tasks reinforce decision patterns  
 094 that remain valid throughout training. While originally studied in reinforcement learning, the under-  
 095 lying principle—aligning task difficulty with model capacity—can be extended to other learning  
 096 settings. For instance, in computer vision, curriculum-based augmentation strategies that gradually  
 097 increase corruption severity during training have gained attention (Lu & Lam, 2023; Choi & Kim,  
 098 2024).

## 100 3 METHODS

### 101 3.1 STOCHASTIC DATA REGULARIZATION

102 We explore data regularization through stochastic transformations applied during training, struc-  
 103 tured along three core dimensions: *selection policies* (e.g., noise injection, random augmentation  
 104 pipelines), *combination policies* (e.g., jointly applying multiple data regularization strategies), and  
 105 *scheduling policy* (e.g., curriculum learning-based approach).

108 **Selection Policies** Selection policies define stochastic mechanisms for perturbing training inputs,  
 109 thereby inducing regularization without altering the underlying task. In this work, we consider two  
 110 complementary strategies under this paradigm: direct **Noise Injection** and **Random Transforma-**  
 111 **tions** drawn from augmentation sets.

112 *Noise Injection* applies input-space corruption by sampling corruption parameters dynamically at  
 113 each step. We define a noise operator  $\nu(\cdot)$  such that:

$$116 \quad \tilde{x} = \nu(x; \theta_t), \quad \theta_t \sim \mathcal{P}(\theta_{\min}, \theta_{\max}), \quad (1)$$

119 where  $\mathcal{P}$  is a generic sampling distribution (e.g., Uniform, Alpha-Stable (Yuan et al., 2025)), and  
 120  $\theta_t$  modulates the corruption strength (e.g., standard deviation for Gaussian noise or the corruption  
 121 factor for Salt-and-Pepper noise). This technique is theoretically equivalent to Tikhonov regularization  
 122 (Bishop, 1995) and is applied exclusively during training.

123 *Random Transformations*, in turn, select stochastic augmentations from a candidate set  $\mathcal{T} =$   
 124  $\{\tau_1, \tau_2, \dots, \tau_k\}$ . At each training step, a random subset  $\mathcal{T}^* \subseteq \mathcal{T}$  is sampled and applied to the  
 125 input with randomized parameters:

$$128 \quad \tilde{x} = \tau(x), \quad \tau \in \mathcal{T}^*, \quad \text{with parameters from predefined ranges.} \quad (2)$$

131 **Combination Policies** In practice, multiple selection policies may be combined—e.g., applying  
 132 both  $\tau$  and  $\nu$  sequentially—to form a compound perturbation strategy. Such combination policies  
 133 can unify coarse (e.g., geometric, shuffling) and fine-grained (e.g., noisy) transformations, enabling  
 134 richer training signals and broader robust generalization capabilities.

136 **Scheduling Policy** We implement a Curriculum Learning-based approach by progressively train-  
 137 ing across stages of increasing regularization, each governed by Early Stopping, as detailed in Al-  
 138 gorithm 1. Each stage  $s \in \{1, \dots, S\}$  introduces a transformation  $\Phi_s(x)$  selected from a predefined  
 139 scheduling sequence  $\mathcal{S} = (\Phi_1, \dots, \Phi_S)$ , which may include selection policies (e.g. Eq. 2, Eq. 1) or  
 140 their combinations.

---

**Algorithm 1** Scheduling Policy with Early Stopping
 

---

142 **Input:** Training set  $\{X^{train}, Y^{train}\}$ , Validation set  $\{X^{val}, Y^{val}\}$   
 143 **Input:** Scheduling sequence  $(\Phi_1, \dots, \Phi_S)$ , Early Stopping patience values  $(p_1, \dots, p_S)$   
 144 **Input:** Hypothesis Space  $\mathcal{H}$ , Optimizer  $\mathcal{U}$ , Loss  $\mathcal{L}$   
 145 **Output:** Final hypothesis  $h_S \in \mathcal{H}$

146 1:  $h_1 \leftarrow h \in \mathcal{H}$ ,  $\epsilon \leftarrow \infty$   
 147 2: **for**  $s = 1$  to  $S$  **do**  
 148 3:    $p \leftarrow 0$ ,  $\tilde{h} \leftarrow h_s$   
 149 4:    $X_s^{train} \leftarrow \Phi_s(X^{train})$   
 150 5:   **while**  $p > p_S$  **do**  
 151 6:      $\tilde{h} \leftarrow \mathcal{U}(X_s^{train}, Y^{train}, \tilde{h})$   
 152 7:      $\hat{Y}^{val} \leftarrow \tilde{h}(X^{val})$   
 153 8:      $\tilde{\epsilon} \leftarrow \mathcal{L}(\hat{Y}^{val}, Y^{val})$  ▷ Validation set used without transformations  
 154 9:     **if**  $\tilde{\epsilon} < \epsilon$  **then**  
 155 10:        $\epsilon \leftarrow \tilde{\epsilon}$ ,  $p \leftarrow 0$ ,  $h_s \leftarrow \tilde{h}$   
 156 11:     **else**  
 157 12:        $p \leftarrow p + 1$   
 158 13:     **end if**  
 159 14:   **end while**  
 160 15:    $h_{s+1} \leftarrow h_s$ ,  $s \leftarrow s + 1$   
 161 16: **end for**

---

162 3.2 CHARACTERIZING OUT-OF-DISTRIBUTION  
163

164 We characterize out-of-distribution (OOD) datasets by estimating their divergence from the in-  
165 distribution data in a shared latent representation space (Algorithm 2). This approach avoids direct  
166 comparisons in the input space, which may be sensitive to superficial or non-semantic differences.  
167 We train an autoencoder ( $\mathbf{h} \equiv (\mathbf{h}_f, \mathbf{h}_g)$ ) on the in-distribution training set and use its encoder ( $\mathbf{h}_f$ ) to  
168 extract latent representations for both the clean test set ( $\mathbf{Z}_{\text{in}}$ ) and each OOD variant ( $\mathbf{Z}_{\text{out}}$ ). To quanti-  
169 fy the shift between these distributions, we employ Kullback–Leibler (KL) divergence, although  
170 other statistical distances (e.g., Wasserstein) are compatible with our framework.

171

**Algorithm 2** Characterizing Out-of-Distribution Data using Latent Representations

---

172 **Input:** In-distribution dataset  $X^{\text{in}} = \{x_i\}_{i=1}^{n_{\text{in}}}$ ,  $x_i \in \mathbb{R}^d$ , with training and test splits:  $X^{\text{train}}, X^{\text{test}}$   
173 **Input:** Out-of-distribution dataset  $X^{\text{out}} = \{x_j\}_{j=1}^{n_{\text{out}}}$ ,  $x_j \in \mathbb{R}^d$   
174 **Input:** Learning Algorithm  $\mathcal{A} : \mathbb{R} \rightarrow (\mathbf{h}_f, \mathbf{h}_g)$  where  $\mathbf{h}_f : \mathbb{R}^d \rightarrow \mathbb{R}^p$  (encoder) and  $\mathbf{h}_g : \mathbb{R}^p \rightarrow \mathbb{R}^d$   
175 (decoder)  
176 **Output:** Out-of-Distribution metric  $M = \text{KL}(\mathbf{Z}_{\text{in}} \parallel \mathbf{Z}_{\text{out}})$

177 1:  $\mathbf{h} \leftarrow \mathcal{A}(X^{\text{train}})$   
178 2:  $\varepsilon \leftarrow 10^{-10}$  ▷ Define a small positive constant  
179 3:  $\mathbf{Z}_{\text{in}} \leftarrow \mathbf{h}_f(X^{\text{test}})$  ▷ Encode in-distribution test dataset  $\mathbf{Z}_{\text{in}} \in \mathbb{R}^{n_{\text{in}} \times p}$   
180 4:  $\mathbf{Z}_{\text{in}} \leftarrow \text{flatten}(\mathbf{Z}_{\text{in}}) + \varepsilon$  ▷ Encode in-distribution test dataset  $\mathbf{Z}_{\text{in}} \in \mathbb{R}^{n_{\text{in}} \times p} \times 1$   
181 5:  $\mathbf{Z}_{\text{in}} \leftarrow \frac{\mathbf{Z}_{\text{in}}}{\mathbf{1}^T \mathbf{Z}_{\text{in}}}$   
182 6:  $\mathbf{Z}_{\text{out}} \leftarrow \mathbf{h}_f(X^{\text{out}})$  ▷ Encode out-of-distribution dataset  $\mathbf{Z}_{\text{out}} \in \mathbb{R}^{n_{\text{out}} \times p}$   
183 7:  $\mathbf{Z}_{\text{out}} \leftarrow \text{flatten}(\mathbf{Z}_{\text{out}}) + \varepsilon$  ▷ Encode in-distribution test dataset  $\mathbf{Z}_{\text{out}} \in \mathbb{R}^{n_{\text{out}} \times p} \times 1$   
184 8:  $\mathbf{Z}_{\text{out}} \leftarrow \frac{\mathbf{Z}_{\text{out}}}{\mathbf{1}^T \mathbf{Z}_{\text{out}}}$   
185 9:  $M = \mathbf{Z}_{\text{in}}^T \cdot \log \mathbf{Z}_{\text{in}} - \mathbf{Z}_{\text{in}}^T \log \mathbf{Z}_{\text{out}}$  ▷ Compute KL Divergence

---

186

187

188 3.3 MISCOVERAGE STATISTICAL ANALYSIS  
189

190 We formalize our miscoverage evaluation via leave-fold-out analysis, detailed in Algorithm 3.

191

**Algorithm 3** Leave-Folds-Out Miscoverage Analysis

---

192 **Input:** Set of  $K$ -Fold Estimates  $\mathcal{F} = \{\mathbf{F}_i : \mathbf{F}_i \in \mathbb{R}^B\}_{i=1}^K$   
193 **Input:** Number  $L$  of folds to leave out  
194 **Output:** Set of Tuples  $\mathcal{R} = \{(\tilde{\mu}_{\mathbf{R}_j}, \sigma_{\mathbf{R}_j})\}_{j=1}^J$

195 1:  $J \leftarrow K - L$   
196 2:  $\mu_{\mathbf{F}} \leftarrow \frac{1}{N} \mathbf{1}^T \mathbf{F}$  ▷ Compute mean  
197 3:  $\mathcal{R} \leftarrow \emptyset$   
198 4: **for**  $j = 1$  to  $J$  **do**  
199 5:    $\mathbf{R}_j \leftarrow \mathbf{F} - \{\mathbf{F}_j, \dots, \mathbf{F}_{j+L-1}\}$  ▷  $\mathbf{R}_j \in \mathbb{R}^{B \times J}$   
200 6:    $\tilde{\mu}_{\mathbf{R}_j} = \mu_{\mathbf{R}_j} - \mu_{\mathbf{F}}$  ▷ Mean-centered  
201 7:    $\sigma_{\mathbf{R}_j} \leftarrow \text{bootstrap}(\mathbf{R}_j)$   
202 8:    $\mathcal{R} \leftarrow \mathcal{R} \cup (\tilde{\mu}_{\mathbf{R}_j}, \sigma_{\mathbf{R}_j})$   
203 9: **end for**

---

204

205

206

207

208 4 EXPERIMENTAL SETUP  
209

210

211

212

213

214

215

216 Our experimental setup is designed to evaluate the impact of data regularization strategies on in-  
217 distribution, out-of-sample and out-of-distribution scenarios, as seen in Figure 3. All results reported  
218 in the main paper refer to the CIFAR-10 (Krizhevsky et al., 2009) dataset and its corrupted variant,  
219 CIFAR-10-C (Hendrycks & Dietterich, 2019). The CIFAR-10-C benchmark includes 19 corruption  
220 types (e.g. JPEG compression, contrast, brightness) across 5 severity levels, resulting in a total  
221 of 95 out-of-distribution datasets. We use F1-score as the evaluation estimator  $\theta$  throughout all  
222 in-distribution and out-of-distribution assessments.

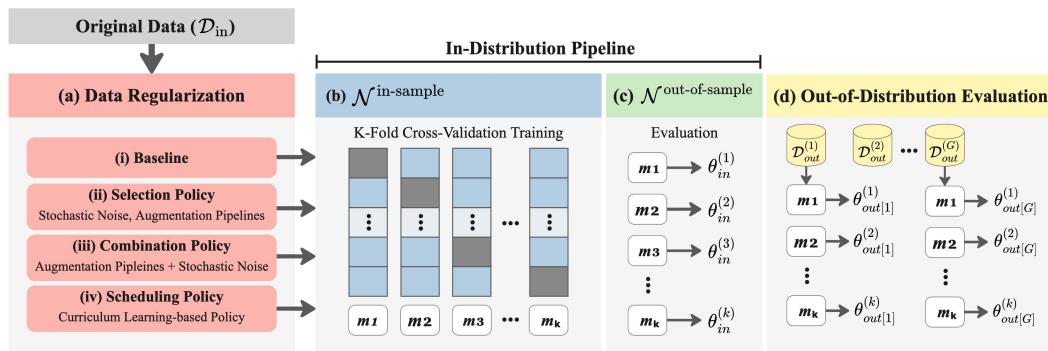


Figure 1: Overview of our evaluation pipeline. **(a)** We apply distinct stochastic data regularization strategies to the original in-sample dataset  $\mathcal{D}_{in}$  for training, including no regularization (Baseline) and our modular framework consisting in 3 different policies—*selection*, *combination*, and *scheduling*—to regularize our data. **(b)** For each strategy, models are trained using  $K$ -fold cross-validation over  $\mathcal{D}_{in}$ . **(c)** Trained models are evaluated on  $\mathcal{D}_{in}$  to obtain  $\theta_{in}^{(i)}$  estimators. **(d)** Each model is then evaluated on a collection of  $G$  corrupted datasets  $\{\mathcal{D}_{out}^{(g)}\}_{g=1}^G$  to compute out-of-distribution estimators  $\theta_{out[g]}^{(i)}$ , enabling robust generalization and miscoverage analysis under domain shift.

#### 4.1 ARCHITECTURES

We evaluate three representative architectures to assess the generalization effects of data regularization. **ResNet-20** (He et al., 2016) serves as a compact and shallow baseline with approximately 280K parameters. **WideResNet-28-10** (Zagoruyko & Komodakis, 2016) is a deeper and significantly wider CNN, totaling over 36M parameters, representing a high-capacity architecture. Finally, **CCT** (Compact Convolutional Transformer) (Hassani et al., 2021) introduces a hybrid transformer-based model with convolutional tokenization and positional encoding, comprising around 930K parameters. This setup allows us to contrast different architectural families—shallow CNNs, wide CNNs, and attention-based models—under a unified training protocol. All models are trained from scratch, without architectural-level regularization (e.g., Dropout or LayerNorm), to isolate the effects of data regularization alone. Inputs are 32×32×3 and training uses a batch size of 128.

#### 4.2 IMPLEMENTATION DETAILS

Early Stopping was applied consistently across all training routines—including both standard and curriculum-based strategies—to prevent overfitting and stabilize convergence. Although this promotes fair evaluations, baseline models often converge prematurely, leading to lower final performance that may differ from typical state-of-the-art results. This is intentional, as our focus lies in understanding the robust generalization capabilities of models rather than maximizing absolute performance.

Standard training strategies used a fixed patience of 10 epochs, while Curriculum Learning stages followed a progressive patience schedule tailored to the difficulty of each stage. The same Early Stopping protocol was also applied to the Autoencoder used in the KL divergence characterization (see Section 3.2), ensuring consistent training dynamics across all components of the experimental pipeline.

#### 4.3 DATA REGULARIZATION STRATEGIES

All stochastic data regularization strategies evaluated in this study are organized within our modular framework of **Selection Policies**, **Combination Policies**, and **Scheduling Policy**. Table 1 summarizes the configuration details, including the maximum Salt & Pepper factor, Gaussian noise standard deviation, RandAugment parameters, Curriculum Learning stage schedule, and Early Stopping values.

We apply **RandAugment** (Cubuk et al., 2020)—composed of both transformations (e.g. color jitter, Gaussian blur, and saturation adjustments) and standard augmentations like random cropping—as a representative *Selection Policy*, using three transformations per image with a fixed magnitude of 0.3. As a *Combination Policy*, we compose RandAugment with additive noise—either Gaussian or Salt & Pepper—to enhance perturbation diversity. At each training step, the noise intensity is dynamically sampled from a uniform distribution:  $\sigma \sim \mathcal{U}(0, \sigma_{\max})$  for Gaussian noise, and  $\lambda \sim \mathcal{U}(0, \lambda_{\max})$  for Salt & Pepper.

Finally, we implement a curriculum-based *Scheduling Policy*, where the regularization severity increases across training stages. All parameter values used in these strategies were selected through a lightweight parameter search. As a result, the Scheduling Policy focuses on the most effective configurations found—namely, RandAugment followed by RandAugment combined with Gaussian noise.

For a visual illustration of the applied perturbations under each strategy, see Appendix A. The GitHub repository will be made publicly available for full reproducibility.

Table 1: Configurations for CIFAR-10 training strategies, including data regularization types, noise levels, and Early Stopping (ES) patience.

Policy	Strategy	Max S&P Factor ( $\alpha$ )	Max Gaussian StdDev ( $\sigma$ )	ES Patience
Baseline	None	—	—	10
Selection Policy	RandAugment	—	—	10
Combination Policy	RandAugment + S&P	0.3	—	10
Combination Policy	RandAugment + GN	—	0.2	10
Scheduling Policy	<b>Stage 1:</b> RandAugment	—	—	3
	<b>Stage 2:</b> RandAugment + GN	—	0.1	5
	<b>Stage 3:</b> RandAugment + GN	—	0.2	8

## 5 RESULTS

To enable severity-aware comparisons, we characterize each of the 95 CIFAR-10-C corruptions by measuring their divergence from the clean CIFAR-10 distribution in latent space and calculating the KL Divergence. This is performed using Algorithm 2. Based on the resulting metric vector, we sort all corruptions and divide them into three severity bands according to their percentile rank: **Lowest** (0–33rd), **Mid-Range** (34–66th), and **Highest** (67–100th). Figure 2 shows this categorization, with average KL values and bootstrapped confidence intervals per group. This severity stratification underpins all subsequent robustness analyses presented in this study.

To evaluate the effectiveness of data regularization strategies, we report F1-scores under both in-distribution (out-of-sample) and out-of-distribution (OOD) scenarios. As shown in Figure 5 and Table 2, augmenting training with stochastic regularization significantly improves performance across all models and severity levels. In particular, we observe that applying RandAugment improves OOD generalization compared to the baseline. Furthermore, combining RandAugment with noise yields even stronger improvements under high-severity corruptions. Finally, the Scheduling Policy consistently achieve superior F1-scores under stronger corruptions, especially for lightweight architectures like ResNet20.

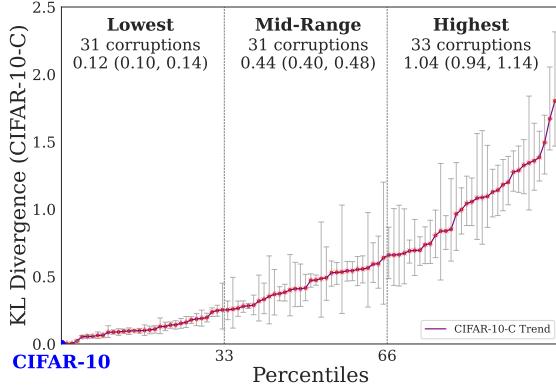


Figure 2: KL-based categorization of CIFAR-10-C corruptions.

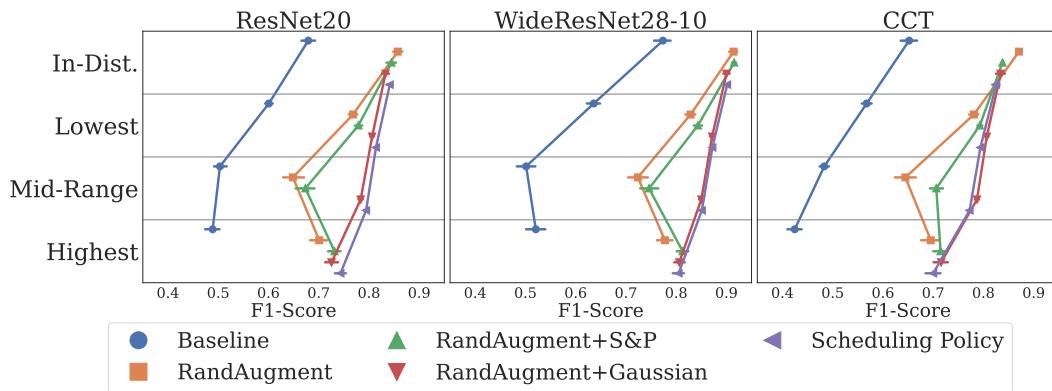


Figure 3: F1-scores for each model under In-Distribution and Out-of-Distribution (OOD) scenarios, with corruptions grouped by severity (Lowest, Mid-Range, Highest).

Table 2: Out-of-Distribution (OOD) characterization for CIFAR-10-C grouped by severity. Values are F1-score (95% CI). Epochs are average values. Best results per column are in bold.

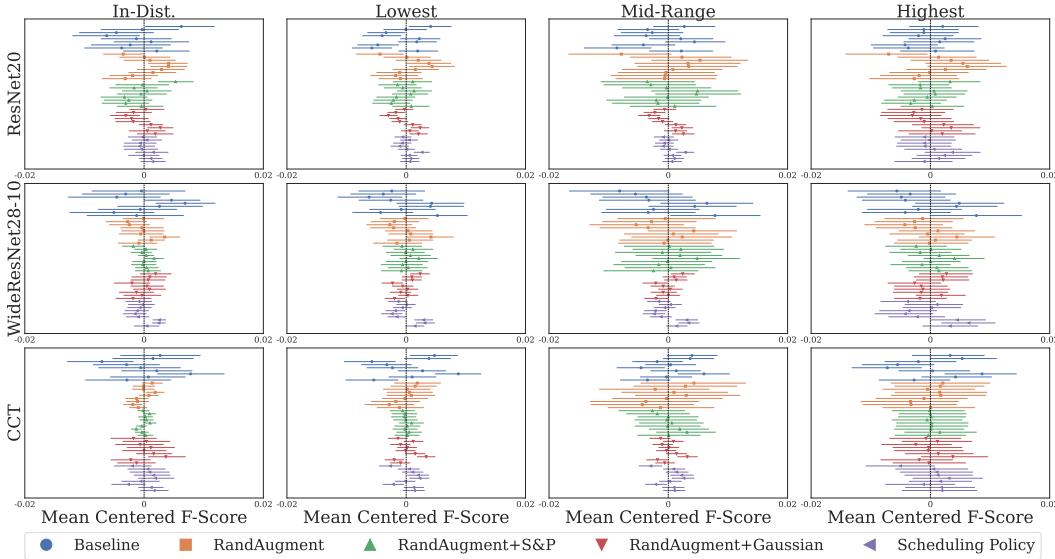
ResNet20	In-Dist.	Lowest	Mid-Range	Highest	Epochs (avg)
Baseline	67.9 (66.6, 69.4)	60.1 (59.3, 60.9)	50.4 (49.1, 51.6)	48.9 (47.5, 50.3)	15.4
RandAugment	<b>85.7 (84.8, 86.7)</b>	76.8 (75.9, 77.8)	64.9 (62.8, 66.8)	70.1 (68.5, 71.8)	51.4
RandAugment+S&P	84.5 (83.6, 85.5)	77.9 (77.2, 78.7)	67.4 (65.6, 69.2)	73.2 (72.0, 74.4)	65.3
RandAugment+Gaussian	83.3 (82.6, 84.0)	80.6 (80.3, 81.0)	78.4 (77.9, 78.8)	72.6 (71.4, 73.9)	63.0
Scheduling Policy	84.2 (83.6, 84.8)	<b>81.4 (81.1, 81.8)</b>	<b>79.4 (79.1, 79.8)</b>	<b>74.3 (73.3, 75.5)</b>	91.6
WideResNet-28-10					
Baseline	77.4 (75.7, 79.2)	63.6 (62.4, 64.9)	50.1 (48.2, 52.0)	52.1 (50.3, 53.8)	19.0
RandAugment	91.4 (90.7, 92.3)	82.9 (81.9, 83.8)	72.4 (70.5, 74.4)	77.8 (76.3, 79.3)	53.6
RandAugment+S&P	<b>91.5 (91.1, 92.0)</b>	84.3 (83.5, 85.1)	74.6 (72.9, 76.4)	<b>81.3 (80.1, 82.6)</b>	56.6
RandAugment+Gaussian	90.0 (89.4, 90.7)	<b>87.2 (86.8, 87.6)</b>	<b>85.0 (84.6, 85.5)</b>	80.6 (79.6, 81.7)	59.9
Scheduling Policy	90.2 (89.7, 90.7)	<b>87.2 (86.8, 87.6)</b>	<b>85.1 (84.7, 85.5)</b>	80.5 (79.4, 81.8)	62.1
CCT					
Baseline	65.2 (63.6, 66.8)	56.8 (55.8, 57.8)	48.3 (47.3, 49.4)	42.4 (41.1, 43.8)	15.6
RandAugment	<b>87.1 (86.7, 87.5)</b>	78.1 (77.2, 79.1)	64.5 (62.4, 66.7)	69.5 (67.7, 71.5)	95.0
RandAugment+S&P	83.8 (83.6, 84.0)	79.3 (78.8, 79.8)	70.6 (69.4, 71.8)	71.5 (70.1, 73.0)	100.0
RandAugment+Gaussian	83.5 (82.6, 84.4)	<b>80.7 (80.3, 81.1)</b>	<b>78.7 (78.3, 79.1)</b>	<b>71.6 (70.3, 73.1)</b>	99.1
Scheduling Policy	82.5 (81.8, 83.3)	79.4 (79.0, 79.9)	77.3 (76.8, 77.7)	70.0 (68.5, 71.6)	77.7

378

379  
380  
381  
Table 3: Standard deviation (95% CI) of F1-scores across CIFAR-10-C severity ranges. Lower  
values indicate more stable performance across folds. Best results per column are in bold.

382	ResNet20	In-Dist.	Lowest	Mid-Range	Highest
383	Baseline	.0237 (.0229, .0244)	.0718 (.0707, .0728)	.1144 (.1134, .1154)	.1337 (.1321, .1354)
384	RandAugment	.0139 (.0135, .0143)	.0643 (.0636, .0650)	.1418 (.1409, .1427)	.1292 (.1278, .1307)
385	RandAugment+S&P	.0157 (.0149, .0164)	.0670 (.0660, .0681)	.1603 (.1590, .1616)	.1162 (.1140, .1184)
386	RandAugment+Gaussian	.0106 (.0101, .0112)	.0313 (.0305, .0319)	.0389 (.0385, .0394)	.1173 (.1140, .1206)
387	Scheduling Policy	<b>.0098 (.0094, .0103)</b>	<b>.0306 (.0299, .0313)</b>	<b>.0338 (.0334, .0341)</b>	<b>.1048 (.1016, .1081)</b>
388	<b>WideResNet-28-10</b>				
389	Baseline	.0273 (.0255, .0293)	.1167 (.1150, .1186)	.1658 (.1646, .1671)	.1695 (.1679, .1710)
390	RandAugment	.0106 (.0103, .0110)	.0674 (.0668, .0681)	.1394 (.1386, .1403)	.1192 (.1180, .1206)
391	RandAugment+S&P	<b>.0072 (.0068, .0076)</b>	.0736 (.0725, .0747)	.1635 (.1622, .1647)	.1128 (.1108, .1148)
392	RandAugment+Gaussian	.0111 (.0107, .0116)	.0356 (.0348, .0365)	.0427 (.0422, .0432)	<b>.1002 (.0979, .1030)</b>
393	Scheduling Policy	.0075 (.0068, .0081)	<b>.0343 (.0335, .0353)</b>	<b>.0399 (.0395, .0404)</b>	.1039 (.1006, .1068)
394	<b>CCT</b>				
395	Baseline	.0250 (.0237, .0261)	.0881 (.0866, .0895)	.0897 (.0889, .0906)	.1317 (.1303, .1330)
396	RandAugment	.0094 (.0090, .0099)	.0588 (.0581, .0595)	.1316 (.1306, .1325)	.1490 (.1472, .1507)
397	RandAugment+S&P	<b>.0041 (.0039, .0043)</b>	.0439 (.0432, .0447)	.1136 (.1124, .1147)	<b>.1319 (.1287, .1350)</b>
398	RandAugment+Gaussian	.0143 (.0136, .0150)	<b>.0349 (.0342, .0357)</b>	<b>.0363 (.0359, .0367)</b>	.1332 (.1299, .1365)
399	Scheduling Policy	.0113 (.0107, .0119)	.0377 (.0369, .0385)	.0400 (.0395, .0405)	.1340 (.1307, .1373)

400 To complement average performance results, we apply our proposed Miscoverage Analysis (Algo-  
401 rithm 3) to assess the stability of model predictions across cross-validation folds. Figure 5 visualizes  
402 miscoverage behaviors, while Table 3 reports the standard deviation and 95% confidence intervals  
403 of F1-scores across severity categories, highlighting the consistency gains from data regularization  
404 strategies.



424  
425  
426  
427  
428  
429  
Figure 4: Mean-Centered F-Score distributions for each data regularization strategy across three ar-  
430  
431  
chitectures (ResNet20, WideResNet-28-10, and CCT) and four evaluation domains: In-Distribution and OOD corruptions grouped by severity levels. Each point represents a mean-centered F1-score from one replication, and horizontal lines denote 95% confidence intervals derived via bootstrap resampling. These results are generated using the Leave-Folds-Out Miscoverage Analysis (Algo-  
432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485  
486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539  
540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593  
594  
595  
596  
597  
598  
599  
600  
601  
602  
603  
604  
605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647  
648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701  
702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755  
756  
757  
758  
759  
750  
751  
752  
753  
754  
755  
756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809  
810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863  
864  
865  
866  
867  
868  
869  
860  
861  
862  
863  
864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917  
918  
919  
910  
911  
912  
913  
914  
915  
916  
917  
918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971  
972  
973  
974  
975  
976  
977  
978  
979  
970  
971  
972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025  
1026  
1027  
1028  
1029  
1020  
1021  
1022  
1023  
1024  
1025  
1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079  
1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133  
1134  
1135  
1136  
1137  
1138  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133  
1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133  
1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133  
1134  
1135  
1136  
1137  
1138  
1139  
1130  
1131  
1132  
1133  
1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187  
1188  
1189  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187  
1188  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241  
1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1240  
1241  
1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295  
1296  
1297  
1298  
1299  
1290  
1291  
1292  
1293  
1294  
1295  
1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320  
1321  
1322

432 

## 6 DISCUSSION

434 Our findings underscore that robust generalization is strongly influenced by how data regularization  
 435 is designed and scheduled during training. Rather than relying solely on isolated techniques, **our**  
 436 **modular framework offers a structured lens for understanding and improving robust gener-**  
 437 **alization**, especially under out-of-distribution shifts.

438 Scheduling Policies (Curriculum Learning-based) consistently with state-of-the-art proves effec-  
 439 tive in enhancing robustness across architectures and corruption severities. It reliably reduces per-  
 440 formance variability (i.e., lower standard deviations) and often outperforms isolated regularization  
 441 strategies. These findings emphasize the importance of the training data presentation order for gen-  
 442 eralization under distribution shift.

443 Combining stochastic data regularization techniques also yields measurable benefits. In particular,  
 444 pairing RandAugment with Gaussian noise improves robustness for WideResNet and CCT, whereas  
 445 Salt & Pepper noise produces less consistent gains. This suggests that the interaction between model  
 446 architecture and corruption type is nontrivial and deserves further attention.

447 Notably, improvements are not solely driven by the type of regularization applied, but also by in-  
 448 creased training exposure. Curriculum-based strategies tend to train longer before Early Stopping is  
 449 triggered. This extended exposure appears necessary for learning robust representations, indicating  
 450 that vision models benefit from prolonged stochastic regularization.

451 Some of the applied data regularization strategies may partially overlap with corruptions present  
 452 in the CIFAR-10-C dataset. To understand the effects under this lens, we include a discussion  
 453 in Appendix B comparing overall performance on the full CIFAR-10-C benchmark and a filtered  
 454 version that excludes overlapping corruptions—specifically, cases where similar transformations  
 455 (e.g., Gaussian noise or contrast adjustments introduced by RandAugment) are used during training  
 456 but are also present in the evaluated corrupted test sets.

457 In summary, our results highlight that the structure and progression of regularization—not just its  
 458 presence—play a critical role in enabling robust generalization, particularly for compact models. A  
 459 tailored and modular approach to data regularization, as proposed in this study, offers a promising  
 460 direction for building more reliable machine learning systems under distributional stress.

463 

## 7 CONCLUSION

464 This work explored stochastic data regularization strategies to improve model robustness under dis-  
 465 tribution shift. We find that organizing these transformations into a curriculum—progressively in-  
 466 creasing complexity during training—consistently leads to more stable and generalizable models.  
 467 While curriculum-based approaches demonstrate strong regularization capabilities, especially un-  
 468 der challenging conditions, we also observe that simpler strategies combining transformations with  
 469 noise injection offer competitive trade-offs in terms of effectiveness and efficiency. These results  
 470 highlight that the structure and dynamics of data exposure can be as important as the regularization  
 471 technique itself. For future work, we aim to explore a Variational Autoencoder to extract a structured  
 472 latent space, thereby enhancing our quantification of distribution shifts. Additionally, we intend to  
 473 apply the methods proposed in this work to other data modalities, such as acoustic signals and text.

475 

## ACKNOWLEDGMENTS

476 Not applicable.

480 

## REFERENCES

481 Stephen Bates, Trevor Hastie, and Robert Tibshirani. Cross-validation: what does it estimate and  
 482 how well does it do it? *Journal of the American Statistical Association*, pp. 1–12, 2023.

483 Yoshua Bengio and Yves Grandvalet. No unbiased estimator of the variance of k-fold cross-  
 484 validation. *Journal of Machine Learning Research*, 5:1089–1105, 2004.

486 Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. Curriculum learning.  
 487 In *Proceedings of the 26th Annual International Conference on Machine Learning*, ICML  
 488 '09, pp. 41–48, New York, NY, USA, 2009. Association for Computing Machinery. ISBN  
 489 9781605585161. doi: 10.1145/1553374.1553380. URL <https://doi.org/10.1145/1553374.1553380>.

491 Chris M Bishop. Training with noise is equivalent to tikhonov regularization. *Neural computation*,  
 492 7(1):108–116, 1995.

494 Juhwan Choi and YoungBin Kim. Colorful cutout: Enhancing image data augmentation with cur-  
 495 riculum learning. *arXiv preprint arXiv:2403.20012*, 2024.

497 Ekin D Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V Le. Randaugment: Practical automated  
 498 data augmentation with a reduced search space. In *Proceedings of the IEEE/CVF conference on*  
 499 *computer vision and pattern recognition workshops*, pp. 702–703, 2020.

500 Sebastian Farquhar and Yarin Gal. What ‘out-of-distribution’ is and is not. In *NeurIPS ML Safety*  
 501 *Workshop*, 2022.

503 Umberto Tenório de Barros Filho, Paulo Rocha, Marcos Oliveira, Andrea Maria Nogueira Caval-  
 504 canti Ribeiro, Rodrigo de Paula Monteiro, and Diego Pinheiro. Regularizing neural networks  
 505 with noise injection for classification of brain tumor in magnetic resonance imaging. In *2023*  
 506 *IEEE Latin American Conference on Computational Intelligence (LA-CCI)*, pp. 1–6, 2023. doi:  
 507 10.1109/LA-CCI58595.2023.10409397.

508 Ali Hassani, Steven Walton, Nikhil Shah, Abulikemu Abuduweili, Jiachen Li, and Humphrey Shi.  
 509 Escaping the big data paradigm with compact transformers. *arXiv preprint arXiv:2104.05704*,  
 510 2021.

512 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-  
 513 nition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp.  
 514 770–778, 2016.

515 Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common cor-  
 516 ruptions and perturbations. *Proceedings of the International Conference on Learning Represen-  
 517 tations*, 2019.

519 Dan Hendrycks, Norman Mu, Ekin D Cubuk, Barret Zoph, Justin Gilmer, and Balaji Lakshmi-  
 520 narayanan. Augmix: A simple data processing method to improve robustness and uncertainty.  
 521 *arXiv preprint arXiv:1912.02781*, 2019.

522 Dan Hendrycks, Nicholas Carlini, John Schulman, and Jacob Steinhardt. Unsolved problems in ml  
 523 safety. *arXiv preprint arXiv:2109.13916*, 2021.

525 Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.  
 526 *Technical Report TR-2009*, 2009.

527 Binghui Li, Jikai Jin, Han Zhong, John Hopcroft, and Liwei Wang. Why robust generalization  
 528 in deep learning is difficult: Perspective of expressive power. *Advances in Neural Information  
 529 Processing Systems*, 35:4370–4384, 2022.

531 Lin Li and Michael W. Spratling. Data augmentation alone can improve adversarial training. In  
 532 *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=y4uc4NtTWaq>.

534 Mingxuan Li, Junzhe Zhang, and Elias Bareinboim. Causally aligned curriculum learning. In  
 535 *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=hp4y0jhwTs>.

538 Chi-Heng Lin, Chiraag Kaushik, Eva L Dyer, and Vidya Muthukumar. The good, the bad and  
 539 the ugly sides of data augmentation: An implicit spectral regularization perspective. *Journal of  
 Machine Learning Research*, 25(91):1–85, 2024.

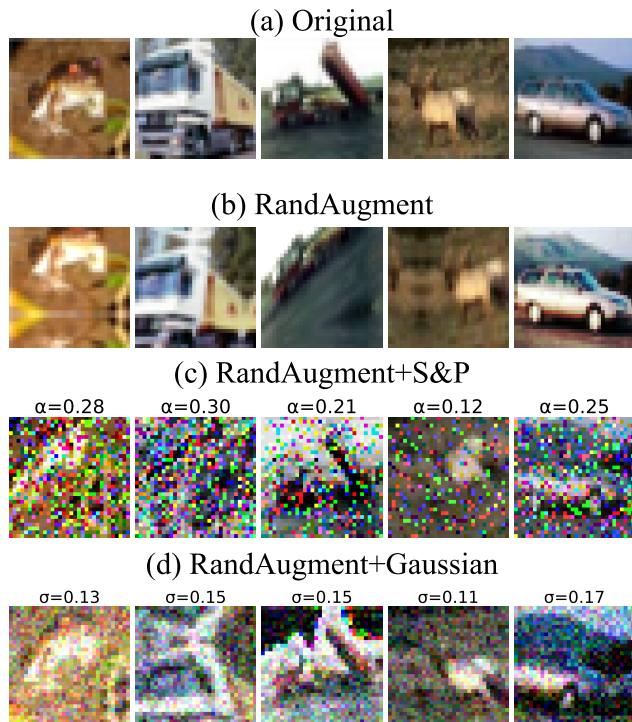
540 Hongyuan Lu and Wai Lam. PCC: Paraphrasing with bottom-k sampling and cyclic learning for  
 541 curriculum data augmentation. In Andreas Vlachos and Isabelle Augenstein (eds.), *Proceed-  
 542 ings of the 17th Conference of the European Chapter of the Association for Computational Lin-  
 543 guistics*, pp. 68–82, Dubrovnik, Croatia, May 2023. Association for Computational Linguis-  
 544 tics. doi: 10.18653/v1/2023.eacl-main.5. URL <https://aclanthology.org/2023.eacl-main.5/>.

545 Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov.  
 546 Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine  
 547 Learning Research*, 15(56):1929–1958, 2014. URL <http://jmlr.org/papers/v15/srivastava14a.html>.

550 Xueqiong Yuan, Jipeng Li, and Ercan Engin Kuruoglu. Robustness enhancement in neural networks  
 551 with alpha-stable training noise. *Digital Signal Processing*, 156:104778, 2025.

553 Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *ArXiv*, abs/1605.07146, 2016.  
 554 URL <https://api.semanticscholar.org/CorpusID:15276198>.

## 556 A DATA REGULARIZATION TRANSFORMATIONS



583 Figure 5: Visual comparison of stochastic data regularization strategies applied to CIFAR-10 sam-  
 584 ples. (a) Original images. (b) RandAugment with 3 transformations per image and magnitude 0.3.  
 585 (c) RandAugment combined with Salt & Pepper noise (*Combination Policy*), where the noise factor  
 586  $\alpha$  is sampled uniformly from  $\mathcal{U}(0, \alpha_{\max} = 0.3)$ . (d) RandAugment combined with Gaussian noise,  
 587 with standard deviation  $\sigma \sim \mathcal{U}(0, \sigma_{\max} = 0.2)$ . The values of  $\alpha$  and  $\sigma$  shown below each image  
 588 indicate the sampled noise intensity for that example.

## 591 B OVERALL RESULTS ACROSS ALL CORRUPTIONS

593 Some of the data regularization strategies we adopt—such as Gaussian noise, Salt & Pepper noise,  
 and RandAugment—introduce transformations that can be thought of as partially overlapping with

specific corruptions in the CIFAR-10-C dataset (e.g., *Gaussian Noise*, *Shot Noise*, *Speckle Noise*, *Impulse Noise*, *Contrast*, *Brightness*). We conducted a sensitivity analysis comparing a dataset with all corruptions included and a dataset (*w/o Overlap*) that excludes these potentially overlapping corruptions. The (*w/o Overlap*) dataset is, when compared to the dataset containing all corruptions, OOD to a greater extent. The results remain consistent, and the exclusion of potentially overlapping corruptions has even improved in some cases.

Table 4: General performance for three models using different augmentation strategies. Values are F1-score (95% CI). Best results per column are in bold.

ResNet20	In-Dist.	All Corruptions	w/o Overlap
Baseline	67.9 (66.5, 69.4)	53.0 (52.3, 53.9)	55.2 (54.5, 56.0)
RandAugment	<b>85.7 (84.8, 86.7)</b>	70.6 (69.7, 71.5)	73.3 (72.3, 74.2)
RandAugment+S&P	84.5 (83.6, 85.5)	72.8 (72.0, 73.6)	75.0 (74.4, 75.7)
RandAugment+Gaussian	83.3 (82.6, 84.0)	77.1 (76.6, 77.6)	76.8 (76.4, 77.2)
Scheduling Policy	84.2 (83.6, 84.8)	<b>78.3 (77.9, 78.8)</b>	<b>78.1 (77.7, 78.5)</b>
<b>WideResNet-28-10</b>			
Baseline	77.4 (75.7, 79.4)	55.2 (54.2, 56.3)	57.3 (56.2, 58.4)
RandAugment	<b>91.4 (90.7, 92.3)</b>	77.7 (76.8, 78.6)	79.6 (78.6, 80.6)
RandAugment+S&P	<b>91.5 (91.1, 92.0)</b>	80.1 (79.3, 80.9)	81.1 (80.2, 82.0)
RandAugment+Gaussian	90.0 (89.3, 90.7)	<b>84.2 (83.7, 84.7)</b>	<b>83.5 (83.1, 84.0)</b>
Scheduling Policy	90.2 (89.8, 90.7)	<b>84.2 (83.7, 84.7)</b>	<b>83.6 (83.2, 84.1)</b>
<b>CCT</b>			
Baseline	65.2 (63.7, 66.8)	49.0 (48.3, 49.8)	50.1 (49.3, 50.9)
RandAugment	<b>87.1 (86.6, 87.5)</b>	70.7 (69.6, 71.7)	75.3 (74.5, 76.2)
RandAugment+S&P	83.8 (83.5, 84.0)	73.7 (73.0, 74.5)	75.4 (74.9, 76.0)
RandAugment+Gaussian	83.5 (82.6, 84.4)	<b>76.9 (76.3, 77.5)</b>	<b>76.8 (76.4, 77.2)</b>
Scheduling Policy	82.5 (81.9, 83.2)	75.4 (74.8, 76.1)	75.3 (74.8, 75.8)

## C ACKNOWLEDGMENT OF LLM USE

We acknowledge the use of large language models to aid in polishing the writing and, primarily, to help build and check how tables and plots could be presented in the best way. The models were not used to generate original research content, experiments, or results.