



**HAL**  
open science

# SLACK: Stable Learning of Augmentations with Cold-start and KL regularization

Juliette Marrie, Michael Arbel, Diane Larlus, Julien Mairal

► **To cite this version:**

Juliette Marrie, Michael Arbel, Diane Larlus, Julien Mairal. SLACK: Stable Learning of Augmentations with Cold-start and KL regularization. CVPR 2023 - IEEE/CVF Conference on Computer Vision and Pattern Recognition, IEEE, Jun 2023, Vancouver, Canada. pp.1-17, 10.1109/cvpr52729.2023.02328 . hal-04135386

**HAL Id: hal-04135386**

**<https://inria.hal.science/hal-04135386>**

Submitted on 20 Jun 2023

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

# SLACK: Stable Learning of Augmentations with Cold-start and KL regularization

Juliette Marrie<sup>1,2</sup>

Michael Arbel<sup>1</sup>

Diane Larlus<sup>2</sup>

Julien Mairal<sup>1</sup>

<sup>1</sup> Univ. Grenoble Alpes, Inria, CNRS, Grenoble INP, LJK

<sup>2</sup> NAVER LABS Europe

## Abstract

*Data augmentation is known to improve the generalization capabilities of neural networks, provided that the set of transformations is chosen with care, a selection often performed manually. Automatic data augmentation aims at automating this process. However, most recent approaches still rely on some prior information; they start from a small pool of manually-selected default transformations that are either used to pretrain the network or forced to be part of the policy learned by the automatic data augmentation algorithm. In this paper, we propose to directly learn the augmentation policy without leveraging such prior knowledge. The resulting bilevel optimization problem becomes more challenging due to the larger search space and the inherent instability of bilevel optimization algorithms. To mitigate these issues (i) we follow a successive cold-start strategy with a Kullback-Leibler regularization, and (ii) we parameterize magnitudes as continuous distributions. Our approach leads to competitive results on standard benchmarks despite a more challenging setting, and generalizes beyond natural images.<sup>1</sup>*

## 1. Introduction

Data augmentation, which encourages predictions to be stable with respect to particular image transformations, has become an essential component in visual recognition systems. While the data augmentation process is conceptually simple, choosing the optimal set of image transformations for a given task or dataset is challenging. For instance, designing a good set for ImageNet [4] or even CIFAR-10/100 [13] has been the result of a long-standing research effort. Whereas data augmentation strategies that have been chosen by hand for ImageNet have been used successfully for many recognition tasks involving natural images, they may fail to generalize to other domains such as medical imaging, remote sensing or hyperspectral imaging.

This has motivated automating the design of data augmentation strategies [10, 12, 14–16, 19, 27, 31]. Those are often represented as a stochastic *policy* that randomly draws a combination of transformations along with their magnitudes from a large predefined set, each time an image is sampled. The goal becomes to learn strategies that effectively compose multiple transformations, which is a challenging task given the large search space of augmentations.

A natural framework for learning the parameters of this policy is that of bilevel optimization. Intuitively, one looks for the best possible policy such that a neural network trained with this policy on a training set (inner problem) generalizes well on a distinct validation set (outer problem). Optimizing the resulting formulation is challenging as the outer problem depends on the solution of the inner problem. Classical techniques for solving this bilevel problem, such as unrolled optimization, can become highly unstable as the network weights become progressively suboptimal for the current policy during the learning process.

Moreover, augmentations are often non-differentiable in the parameters of the policy, thus requiring techniques other than direct differentiation, such as Bayesian optimization [15], gradient approximations (*e.g.* RELAX [7]), or the score method / REINFORCE [28] algorithm. While these techniques bypass the differentiability issues, they can suffer from large bias or variance. As a result, learning augmentation policies is a difficult problem whose challenges are exacerbated by the inherent instability of the optimization techniques developed to solve bilevel problems, such as unrolled optimization [1].

A standard way to improve stability and make the automatic data augmentation problem simpler is to reduce the search space. This is often achieved by learning the policy on top of “default” transformations such as Cutout [5], random cropping and resizing, or color jittering, all known to be well-suited to natural images which compose standard benchmarks such as CIFAR or ImageNet, or by discarding transformations known to be harmful such as Invert. Fixing some of the transformations and removing others mitigate the challenges inherent to learning a compo-

<sup>1</sup>Project page: <https://europe.naverlabs.com/slack>

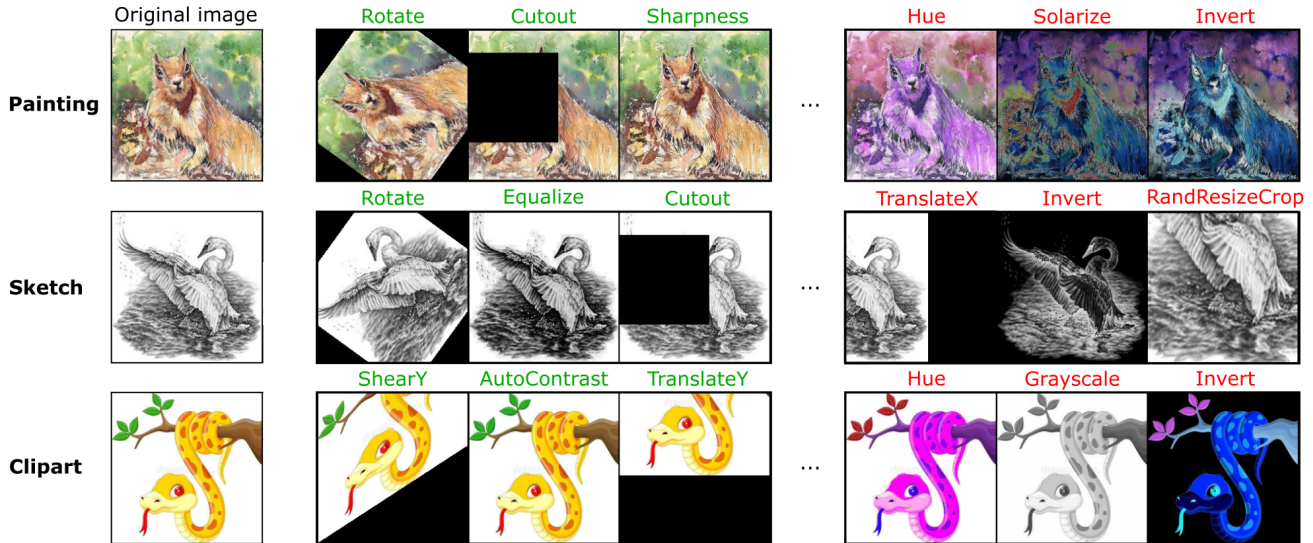


Figure 1. For different domains of the DomainNet dataset [21] (one per line), we show an image from that domain (left) and that image transformed using the three most likely (middle) and the three least likely (right) augmentations for that domain, as estimated by SLACK.

sition of transformations. TrivialAugment [19] also shows that state-of-the-art results can be achieved on these previous benchmarks simply by directly applying the policy classically used for initializing auto-augmentation models, up to minor modifications. Moreover, all methods rely on carefully chosen ranges that constraint the transformation’s magnitudes. Despite its effectiveness, manually selecting default transformations and magnitude ranges restricts the applicability of such policies to natural images and prevents generalisation to other domains.

In this paper, our goal is to choose augmentation strategies without relying on default transformations nor on hand-selected magnitude ranges known to suit common benchmarks. To achieve this objective, we first introduce a simple interpretable model for the augmentation policies which allows learning both the frequency by which a given augmentation is selected and the magnitude by which it is applied. Then, we propose a method for learning these augmentation policies by solving a bilevel optimization problem. Our method relies on the REINFORCE technique for computing the gradient of the policy and on unrolled optimization for learning the policy, both of which can result in instabilities and yield high variance estimates.

To address these issues, we introduce an efficient multi-stage algorithm with a cold-start strategy and a Kullback-Leibler (KL) regularization that are designed to improve the stability of the process for learning the data augmentation policy. More precisely, the algorithm first pre-trains a network with a data augmentation policy uniformly sampling over all transformations. Then, each stage uses a “cold-start” strategy by restarting from the pre-trained network and performs incremental updates of the current policy.

This multi-stage approach with cold start prevents the network from becoming progressively suboptimal as the policy is updated using unrolled optimization. The KL regularization defines a trust region for the policy to compensate for the possibly high variance of gradient estimates obtained using the REINFORCE technique and encourages exploration during training, preventing collapse to trivial solutions. This regularization is inspired by proximal point algorithms in convex optimization [22], which have also been successful in reinforcement learning tasks [23].

By combining the regularized multi-stage approach with our interpretable model of the augmentation policies, we obtained the proposed SLACK method, which stands for *Stable Learning of Augmentations with Cold-start and Kullback-Leibler regularization*. SLACK is an efficient data augmentation learning method that is able to address the challenging bilevel optimization problem of learning a stochastic data augmentation policy without relying strongly on prior knowledge. Figure 1 illustrates the transformations found by SLACK to be most important / detrimental on a dataset of different domains including non-natural images.

To summarize, our contribution is threefold. (i) We propose a simple and interpretable model of the policies which allows learning both frequency and magnitudes of the augmentations. (ii) We propose a regularized multi-stage strategy to improve the stability of the bilevel optimization algorithm used for solving the data augmentation learning problem. (iii) We evaluate our method on challenging experimental settings, and show that it finds competitive augmentation strategies on natural images without resorting to prior information and generalizes to other domains.

## 2. Related Work

The choice of image transformation (also known as data augmentation) has become central in the design of computer vision pipelines. To remove the burden of manual selection, automatic data augmentation strategies have been proposed [20, 26]. AutoAugment [2], one of the earliest methods, uses for instance a recurrent neural network for designing the augmentation policy. Because such an approach requires retraining a prediction model at each iteration, it is prohibitively slow, and more efficient alternatives have been proposed. They aim at reducing the training cost using, *e.g.*, population-based training [12], Bayesian optimization [15], and more recently, gradient-based approaches based on bilevel optimization [9, 10, 14, 15, 18, 27, 30], relying on various gradient estimation techniques such as RELAX [7] or the Score method [28]. While the former is inherently biased, the latter is theoretically exact, but has a high variance when approximated in the context of stochastic optimization. Therefore, these approximations may lead to diverging gradient updates. Our method alleviates this by introducing a KL regularization that defines a trust-region for the policy.

**Automatic augmentation using prior knowledge.** Most previous works learn augmentations using a small network learned on a subset of the dataset of interest, before retraining the prediction model on a larger network using the full (augmented) data. This choice is appealing to recent gradient-based methods [9, 14] as the search phase for an augmentation policy is often reduced to minutes. Nevertheless, [10, 30] have observed that policies found with such a reduced setup may be suboptimal compared to approaches exploiting full datasets for training both the augmentation policy and the prediction model. This observation was confirmed in [3], which shows that a naive grid search could actually yield state-of-the-art results when directly training on the full-size network and the full data. These results are however obtained at the expense of using strong prior knowledge: augmentation policies are applied on top of default transformations that are manually and independently chosen for each benchmark. Lately, [19] has shown that with a few additional careful choices regarding the augmentation policies, applying a single random transformation on top of the default ones could lead to state-of-the-art results.

To avoid relying on default augmentations, DeepAA [31] has recently proposed a greedy approach that is able to learn these transformations. Yet, learning is performed after a “pre-training” phase leveraging the usual default transformations. Moreover, while such a greedy approach simplifies the search procedure and reduces its stochasticity, the resulting computational cost is high. Instead, our approach improves stability and allows directly learning the joint probability of sampling multiple transformations, reducing the search time twofold compared to DeepAA.

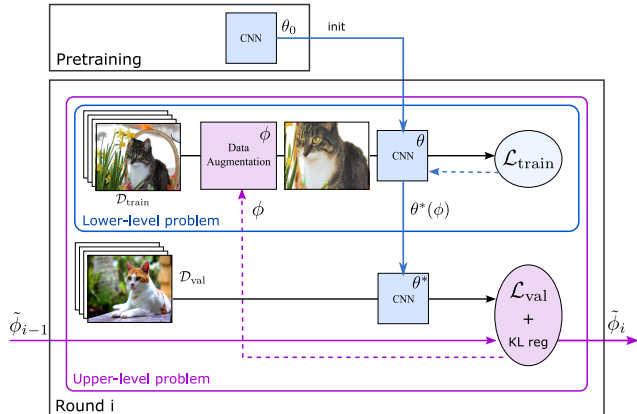


Figure 2. **Overview of our proposed SLACK method.** We learn a data augmentation policy parameterized by  $\phi$  using bilevel optimization. The inner loop finds the optimal network parameter  $\theta^*$  on images from  $\mathcal{D}_{\text{train}}$ . The outer loop trains on a disjoint set of images  $\mathcal{D}_{\text{val}}$  using this network and finds the optimal transformation parameters  $\phi$ . The method is enhanced with i) a cold-start strategy that structures the learning into rounds which share  $\phi$  but restart the network from the pretrained one, and ii) a KL regularization.

## 3. Method

Our method, SLACK, defines an *augmentation policy*, which is a probabilistic model for generating data augmentations. The goal is to learn the parameters of this augmentation policy so as to improve the performance of the trained classifier on a held-out dataset. We first describe the augmentation policy (Sec. 3.1) and then formalize the problem of learning data augmentations with bilevel optimization (Sec. 3.2). We then describe our approach for solving such a bilevel optimization problem, which aims at stabilizing the optimization (Sec. 3.3).

### 3.1. Stochastic data augmentations policy

An augmentation function  $\tau$  transforms an image  $x$  into another *augmented* image  $\tau(x)$  of the same dimensions. We consider composite augmentations obtained by combining simpler augmentations selected from a finite set  $\mathcal{S} = \{s_1, \dots, s_N\}$  of  $N$  candidate *elementary transformations*, such as rotations, translations, shearing, etc. Each elementary transformation depends on a *magnitude* parameter  $m$  that controls the strength of the transformations, for instance, the angle by which an image is rotated. Magnitudes are normalized to be in the unit interval  $[0, 1]$ .

**Augmentation policy.** We define the augmentation policy as a probabilistic model  $p_\phi$  that generates composite augmentations given some parameter  $\phi$  to be learned. The model generates an augmentation in three steps: (1) it samples  $K$  elementary transformations  $t_1, \dots, t_K$  from  $\mathcal{S}$  according to a categorical distribution  $p_\pi$  of parameter  $\pi$ , (2)

it samples values for the magnitudes  $m_1, \dots, m_K$  for each of the selected elementary transformations  $t_k$  according to a smoothed uniform distribution  $p_\mu$  of parameter  $\mu$  and (3) it composes the  $K$  elementary transformations to obtain the composite augmentation, with each  $t_k$  applied using its corresponding magnitude  $m_k$ . Therefore, the augmentation policy  $p_\phi(\tau)$  takes the form

$$p_\phi(\tau) = \prod_{i=1}^K p_\pi(t_i) p_\mu(m_i | t_i), \quad (1)$$

where the parameters  $\phi = (\pi, \mu)$  are learned jointly. Next, we describe the sampling of the transformations and of their magnitudes.

**Sampling transformations.** We sample elementary transformations  $t_k$  with replacement from a categorical distribution  $\text{Cat}_{\pi_k}$  of dimension  $N$  parameterized by a logit vector  $\pi_k := (\pi_{k,n})_{1 \leq n \leq N}$ . The probability  $p_\pi(t_1, \dots, t_K)$  of sampling the  $K$  transformations is given by:

$$p_\pi(t_1, \dots, t_K) = \prod_{k=1}^K \text{Cat}_{\pi_k}(t_k), \quad (2)$$

where we collect all logits to form a parameter matrix  $\pi$  of size  $K \times N$ . These parameters are learned.

**Sampling magnitudes.** The magnitudes of each elementary transformation  $s_i$  in  $S$  are sampled from a smoothed uniform distribution between  $[0, \mu_i]$  whose upper-bound  $\mu_i$  is learned. More precisely, the distribution’s density is defined as

$$p_{\mu_i}(m_i) = \frac{1}{\mu_i} \int_0^{\mu_i} \mathcal{N}(m_i, \sigma)(u) du,$$

where  $\mathcal{N}(m_i, \sigma)$  is the Gaussian distribution of mean  $m_i$  and deviation  $\sigma$ . The density  $p_{\mu_i}(m_i)$  approximates the uniform distribution  $\frac{1}{\mu_i} \mathbf{1}_{[0, \mu_i]}$  as the deviation  $\sigma$  approaches 0. In practice, we set  $\sigma = 0.1$  as we found it to achieve a good trade-off between smoothing and approximation.

**Why a uniform distribution?** We ran some ablations on previous methods (see supplementary) which suggest that a uniform sampling works on par with more elaborate sampling strategies, and that the magnitude range has more impact on the results.

### 3.2. Bilevel formulation for policy search

We consider a prediction task, such as predicting the class  $y$  of some natural image  $x$  using a model  $f_\theta(x)$  with parameter  $\theta$ . We are interested in finding the best policy parameter  $\phi$  so that the prediction model  $f_\theta$ , when trained using such policy on a training set  $\mathcal{D}$  of input/output pairs  $(x, y)$ , generalizes well on the test set  $\mathcal{D}_{\text{test}}$ . The problem naturally decomposes in two phases. During the *search phase*, the optimal augmentation policy  $p_\phi$  is learned on  $\mathcal{D}$ .

During the *evaluation phase*, the model is re-trained on  $\mathcal{D}$  using  $p_\phi$  and is then evaluated on  $\mathcal{D}_{\text{test}}$ . The *evaluation phase* is performed using standard optimization methods. However, the *search phase* requires solving a complex optimization problem that we describe next.

**Search phase as a bilevel problem.** The *search phase* naturally writes as a bilevel problem involving two interdependent losses: a lower-level loss  $\mathcal{L}_{\text{train}}(\theta, \phi)$  for learning an optimal model parameter  $\theta^*(\phi)$  obtained using the augmentation policy  $p_\phi$  and an upper-level loss  $\mathcal{F}(\phi)$  for learning the policy parameter  $\phi$  by evaluating the optimal model with parameter  $\theta^*(\phi)$ . Each of these objectives is evaluated on two separate splits of the available data  $\mathcal{D}$ : a training split  $\mathcal{D}_{\text{train}}$  for the lower-level loss and a validation split  $\mathcal{D}_{\text{val}}$  for the upper-level loss. Below, we describe both losses.

**Lower-level loss.** We first introduce the training loss  $\ell_{\text{train}}(\theta, \tau)$  when only a fixed augmentation  $\tau$  is used:

$$\ell_{\text{train}}(\theta, \tau) := \mathbb{E}_{(x,y) \sim \mathcal{D}_{\text{train}}} [\ell(y, f_\theta(\tau(x)))],$$

where  $(x, y)$  is an (image,label) pair drawn from  $\mathcal{D}_{\text{train}}$  and  $\ell$  is a pointwise prediction loss (e.g. cross-entropy). We then define the training loss  $\mathcal{L}_{\text{train}}(\theta, \phi)$  for an augmentation policy  $p_\phi$  by taking the expectation of  $\ell_{\text{train}}(\theta, \tau)$  over augmentations  $\tau$  sampled according to the policy  $p_\phi$ :

$$\mathcal{L}_{\text{train}}(\theta, \phi) := \mathbb{E}_{\tau \sim p_\phi} [\ell_{\text{train}}(\theta, \tau)].$$

Hence, for a given policy  $p_\phi$ , the goal is to learn the optimal model parameter  $\theta^*(\phi)$  by minimizing  $\mathcal{L}_{\text{train}}(\theta, \phi)$  over  $\theta$ .

**Upper-level loss.** We first denote by  $\mathcal{L}_{\text{val}}(\theta)$  the validation loss for a given model of parameter  $\theta$ :

$$\mathcal{L}_{\text{val}}(\theta) := \mathbb{E}_{(x,y) \sim \mathcal{D}_{\text{val}}} [\ell(y, f_\theta(x))].$$

The validation loss  $\mathcal{L}_{\text{val}}(\theta)$  is computed over the validation set  $\mathcal{D}_{\text{val}}$  *without* applying any augmentation and thus provides a proxy for the performance on the test dataset. We then define the upper-level loss to be the validation loss of an optimal model  $\theta^*(\phi)$  learned using a policy  $p_\phi$ :

$$\mathcal{F}(\phi) := \mathcal{L}_{\text{val}}(\theta^*(\phi)). \quad (3)$$

While optimizing the lower-level loss is relatively standard, minimizing the upper-level loss  $\mathcal{F}$  is more challenging due to the complex dependence of the optimal model parameter  $\theta^*(\phi)$  on the policy. Next, we describe our proposed algorithm for solving the bilevel problem.

### 3.3. SLACK algorithm

We propose Algorithm 1 for learning the optimal policy during the search phase. SLACK first pre-trains the prediction model using the objective  $\mathcal{L}_{\text{train}}(\theta, \phi_{\text{uniform}})$  for an

initial policy parametrized by  $\phi_{\text{uniform}}$  which samples uniformly among all elementary transformations. It then performs  $n_{\text{rounds}}$  rounds to update the parameters  $\theta$  and  $\phi$  jointly using a bilevel optimization algorithm. This approach is reminiscent of the one of AutoAugment (AA) [2] that fully re-trains the network for each policy update. Yet, it is several orders of magnitude faster than AA, as it benefits from pre-training and from our bilevel optimization.

SLACK relies on two strategies to ensure the stability of the parameter updates during each round: a *cold-start strategy* for the prediction model and an *anchoring strategy* for the policy. The *cold-start* initializes the prediction model at the beginning of each round using a pre-trained model  $\theta_0$ . *Anchoring* is achieved by encouraging the current policy to remain close to some *anchor policy*  $p_{\tilde{\phi}}$ . We set  $\tilde{\phi}$  to the current policy parameter  $\phi$  at the beginning of each round. During the first  $n_{\text{retrain}}$  steps of each round, the algorithm only updates the model parameter using a stochastic estimate  $\hat{g}_\theta$  of  $\nabla_\theta \mathcal{L}_{\text{train}}(\theta, \phi)$  while maintaining the policy fixed. Then for the last  $n_{\text{total}} - n_{\text{retrain}}$  steps, the algorithm alternates between model updates and policy updates. The policy updates aim to minimize the sum of the upper-level objective  $\mathcal{F}$  and an anchoring  $d(p_\phi, p_{\tilde{\phi}}) := \text{KL}(p_\pi, p_{\tilde{\pi}})$  encouraging the policy  $p_\phi$  to remain close to the *anchor policy*  $p_{\tilde{\phi}}$ . These updates are obtained using a stochastic estimate  $\hat{G}_\phi$  along with the exact gradient of the KL regularization which admits a closed-form expression. Next we explain how we estimate the gradients  $\hat{g}_\theta$  and  $\hat{G}_\phi$  and discuss the effect of cold-start and KL-regularization.

---

#### Algorithm 1 SLACK

---

```

1: Initialize policy parameter  $\phi \leftarrow \phi_{\text{uniform}}$ .
2: Pre-training:  $\theta_0 \leftarrow \text{optimize}(\mathcal{L}_{\text{train}}(\theta, \phi))$ .
3: for  $i \in \{1, \dots, n_{\text{rounds}}\}$  do
4:   Cold-start:  $\theta \leftarrow \theta_0$ .
5:   Update anchor policy:  $\tilde{\phi} \leftarrow \phi$ .
6:   for  $j \in \{1, \dots, n_{\text{total}}\}$  do
7:     Compute stochastic gradient  $\hat{g}_\theta \approx \nabla_\theta \mathcal{L}_{\text{train}}(\theta, \phi)$ 
8:     Update  $\theta$ :  $\theta \leftarrow \theta - \eta \hat{g}_\theta$ .
9:     if  $j > n_{\text{retrain}}$  then
10:      Compute stochastic gradient  $\hat{G}_\phi \approx \nabla_\phi \mathcal{F}(\phi)$ 
11:      Update  $\phi$ :  $\phi \leftarrow \phi - \alpha(\hat{G}_\phi + \lambda \nabla_\phi d(p_\phi, p_{\tilde{\phi}}))$ .
12:     end if
13:   end for
14: end for

```

---

**Gradient estimation.** Algorithm 1 requires estimating the gradient of  $\mathcal{F}(\phi)$ , which is challenging given the complex dependence of the upper-level loss on the policy  $p_\phi$  through the optimal model parameter  $\theta^*(\phi)$  learned using such a policy. In line with previous works [9, 14, 18, 27], we approximate the optimal model parameter  $\theta^*(\phi)$  with a simpler function  $\hat{\theta}(\phi)$  that is easier to compute:

$$\hat{\theta}(\phi) := \theta - \eta \nabla_\theta \mathcal{L}_{\text{train}}(\theta, \phi). \quad (4)$$

Eq. (4) corresponds to one gradient step to optimize the lower-level loss starting from the current parameter  $\theta$  and  $\phi$  and using step-size  $\eta > 0$ . By keeping track of the dependence in  $\phi$  and exploiting the fact that the augmentation policy  $p_\phi$  has a score  $\nabla_\phi \log p_\phi(\tau)$  that can be computed explicitly using Eq. (1), we can use the REINFORCE/Score method [6] to derive a closed-form expression for  $\nabla_\phi \hat{\theta}(\phi)$  which will serve for approximating the gradient of  $\mathcal{F}$ :

$$\nabla_\phi \hat{\theta}(\phi) = -\eta \mathbb{E}_{\tau \sim p_\phi} [\nabla_\theta \ell_{\text{train}}(\theta, \tau) \nabla_\phi \log p_\phi(\tau)^\top].$$

Then, we approximate the upper-level loss  $\mathcal{F}(\phi)$  with a simpler function  $\hat{\mathcal{F}}(\phi) := \mathcal{L}_{\text{val}}(\hat{\theta}(\phi))$  and the gradient  $\nabla_\phi \mathcal{F}(\phi)$  with  $\nabla_\phi \hat{\mathcal{F}}(\phi)$  which is obtained using the chain rule:

$$\nabla_\phi \mathcal{F}(\phi) \approx \nabla_\phi \hat{\mathcal{F}}(\phi) = \nabla_\theta \mathcal{L}_{\text{val}}(\hat{\theta}(\phi))^\top \nabla_\phi \hat{\theta}(\phi). \quad (5)$$

The above expression requires only first-order derivatives and matrix-vector products, which is amenable to efficient implementation using automatic differentiation softwares.

**Stochastic gradient estimates.** In practice, we replace all expectations by estimates on a batch of data and sampled augmentations. More precisely, to compute the approximation  $\hat{g}_\theta$  to  $\nabla_\theta \mathcal{L}_{\text{train}}(\theta, \phi)$ , we sample  $B_{\text{aug}}$  augmentations from  $p_\phi$  and then apply each of them to a batch of training data  $B_{\text{train}}$  from  $\mathcal{D}_{\text{train}}$ . Using the same batch of data and augmentation, we approximate  $\hat{\theta}(\phi)$  and  $\nabla_\phi \hat{\theta}(\phi)$  appearing in Eq. (5). Finally, we use a batch  $B_{\text{val}}$  of data from  $\mathcal{D}_{\text{val}}$  to estimate  $\nabla_\theta \mathcal{L}_{\text{val}}(\hat{\theta}(\phi))$  and compute  $\hat{G}_\phi$ , which is a stochastic estimate of  $\nabla_\phi \hat{\mathcal{F}}(\phi)$  in Eq. (5).

**Cold-start.** The cold-start strategy allows to re-train the model at each round with the current augmentation policy starting from the pre-trained model. This approach is closer to the original bilevel formulation which implies finding an optimal prediction model for each policy. Initializing with a pre-trained model yields computational gain as fewer iterations are needed to optimize the model. We could instead use a *warm-start* strategy which initializes the model at each round with the learned model at the previous round. Yet we experimentally observe that such approach progressively leads to overfitting and degrades the quality of the learned policies (see supplementary).

**Anchoring using KL regularization.** We experimentally found that adding an anchoring  $d(p_\phi, p_{\tilde{\phi}}) := \text{KL}(p_\pi, p_{\tilde{\pi}})$  with strength parameter  $\lambda$  when updating the policy prevents the algorithm from collapsing towards trivial policies. The anchoring affects only the categorical distribution  $p_\pi$ . For the magnitudes  $p_\mu$ , we did not use anchoring as it is ill-defined for a uniform distribution. Instead, we simply used smaller step-sizes. Our approach takes inspiration from Proximal Policy Optimization [23] used in the context of reinforcement learning which is known to improve policy search.

## 4. Experiments

In this section, we first briefly describe our experimental setup (Sec 4.1). Then we evaluate our approach on several standard benchmarks composed of natural images (Sec 4.2) as well as on a benchmark with other domains (Sec 4.3). We finally report some ablation studies (Sec 4.4).

### 4.1. Experimental setup

**Benchmarks.** We first evaluate our model on three standard benchmarks, CIFAR10 [13], CIFAR100 [13] and ImageNet-100 [25], all composed of natural images. To study how well our method generalizes beyond natural images, we also evaluate on the DomainNet dataset [21], which contains 345 classes for 6 different domains. To ensure our protocol uses a similar number of training images for each domain, we use a reduced set of 50,000 training images for the two largest domains (real, quickdraw) and leave the remaining images for testing. For the other domains, we isolate 20% of the data for testing.

**Architectures.** CIFAR10/100 are evaluated with two architectures that are standard for automatic data augmentation: WideResNet-40x2 and WideResNet-28x10 [29]. Unlike previous works whose search phase is only conducted with the smaller WideResNet-40x2, we search and evaluate with the same architecture, as we found it to be better (see Sec. 4.4). ImageNet-100 and DomainNet are evaluated with a ResNet-18 [11] architecture.

**Transformation space.** Our data augmentation search space is composed of the standard pool of 15 transformations: *Identity*, *ShearX*, *ShearY*, *TranslateX*, *TranslateY*, *Rotate*, *AutoContrast*, *Equalize*, *Invert*, *Solarize*, *Posterize*, *Contrast*, *Brightness*, *Sharpness*, *Color*. We add to this pool the transformations that previous methods usually apply by default: Cutout and RandomCrop for CIFAR, RandomResizeCrop for ImageNet, Grayscale for DomainNet. Following standard practice, when RandomResizeCrop is sampled, it is always applied first. We learn the range of its scale parameter. We do not add ColorJitter that is also applied by default in prior work for ImageNet, as it is already a mix of Brightness, Contrast and Color. However we add Hue, which is one component of ColorJitter and never applied by default. Following prior work, the magnitudes are mapped to [0,1]. After mapping,  $\mu$  is initialized at 0.75 to favour exploration (see details in supplementary). We also uniformly sample magnitudes for Cutout and RandomCrop, whereas their value is hand-picked in prior work. Since the datasets are horizontally symmetric, we follow common practice and apply flip by default.

**Policy search.** We apply a train/val split of 0.5/0.5, meaning that half of the data is used to train the model parameters while the other half is used to learn the augmentation policy. Pre-training is done in the same setting as the evaluation

(see next paragraph), except that we train only with the train data in the train/val split of the search phase. We use SGD with momentum for the optimization of the validation and training losses. For the latter, we use the same weight decay as for the final policy evaluation. We sample 8 different augmentations for computing the expectation that is needed for the stochastic gradient estimate, as detailed in Sec. 3.3.

**Policy evaluation.** We evaluate our models following the framework of TrivialAugment [19]. The corresponding hyperparameters can be found in the supplementary. We evaluate each policy with 4 independent runs, meaning that our results are averaged over a total of  $4 \times 4 = 16$  evaluations. Our Uniform policy (corresponding to SLACK’s initialization) and our reported results on TrivialAugment are evaluated with 8 independent runs. We also report a confidence interval which contains the true mean with probability  $p = 95\%$ , under the assumption of normally distributed accuracies.

### 4.2. Comparison with the state of the art

We compare our method with a Uniform augmentation policy as well as many previous approaches for data augmentation, including AutoAugment (AA) [2], Fast AutoAugment (FastAA) [15], Differentiable Automatic Data Augmentation (DADA) [14], RandAugment (RA) [3], Teach Augment [24], UniformAugment [20], TrivialAugment (TA) [19], and Deep AutoAugment (DeepAA) [31].

For each method, we indicate the total number of composed transformations, and the number of hard-coded transformations among those (Tables 1 and 2). For SLACK, we evaluate the policies obtained from 4 independent search runs (each with 4 different train/val splits) to assess the robustness of our approach. We follow the same process when reproducing DeepAA on CIFAR10/100. Note that all previous methods use a single run for search, before evaluating the policy with one or multiple runs. We report 95% confidence intervals for those evaluating with multiple runs.

The supplementary provides qualitative results showing the evolution of the probability distributions over the transformations and the final estimated policies for all datasets.

**CIFAR.** In Table 1, we observe that, despite not hard-coding Cutout and RandomCrop in our policy, our method is competitive on both CIFAR10 and CIFAR100.

We found that, in general, Cutout and Rotate are selected with a high probability, while the Invert transformation is systematically discarded (see supplementary). This is consistent with the choices made in practice by prior work of adding/removing these transformations manually.

We observe a mismatch between DeepAA’s reported results [31], and those we obtain when evaluating their approach on multiple search runs, using the author’s code and following their recommendations. This is likely due to the stochasticity of the search procedure.

	# Augmentations		CIFAR10		CIFAR100	
	Total	Hard-coded	WRN-40-2	WRN-28-10	WRN-40-2	WRN-28-10
AA [2]	4	2	96.3	97.4	79.3	82.9
FastAA [15]	4	2	96.4	97.3	79.4	82.7
DADA [14]	4	2	96.4	97.3	79.1	82.5
RA [3]	4	2	-	97.3	-	83.3
TeachA [24]	4	2	-	97.5	-	83.2
UniformAugment [17]	4	2	96.25	97.33	79.01	82.82
TA (Wide) [19]	3	2	96.32 ± .05	97.46 ± .06	79.86 ± .19	84.33 ± .17
Uniform policy	3	0	96.12 ± .08	97.26 ± .07	78.79 ± .25	82.82 ± .24
DeepAA [31]	6**	0*	-	97.56 ± .14	-	84.02 ± .18
DeepAA (reproduced)†	6**	0*	96.25 ± .11	97.27 ± .11	79.26 ± .35	83.38 ± .33
<b>SLACK (Ours)</b>	3	0	96.29 ± .08	97.46 ± .06	79.87 ± .11	84.08 ± .16

Table 1. Test accuracies on CIFAR10 and CIFAR100. For SLACK and DeepAA (reproduced) we conduct 4 independent searches, and evaluate each policy with 4 evaluation runs, meaning that we report averages over 16 evaluations. TA and DeepAA are also evaluated with multiple evaluation runs. Results for the remaining methods are reported from the corresponding papers and based on a single run.

\*: DeepAA uses hard-coded transformations for pre-training. \*\*: DeepAA learns random flipping unlike other baselines. † We evaluate the policies found from 4 independent search runs as we do for SLACK, using the code from the authors and following their recommendations.

	# Augmentations		ImageNet-100
	Total	Hard-coded	ResNet18
TA (RA) [19]†	5	4	85.87 ± .30
TA (Wide) [19]†	5	4	86.39 ± .18
Uniform policy	3	0	85.78 ± .32
<b>SLACK</b>	3	0	86.06 ± .11

Table 2. Test accuracies on ImageNet-100.

**ImageNet-100.** Results for ImageNet-100 are reported in Table 2. We compare SLACK to our Uniform policy and to TrivialAugment (RA) and (Wide) variants, the latter using larger magnitude ranges for its random transformation. SLACK’s results lie in between both variants and improve over our Uniform policy.

Interestingly, for ImageNet-100, we found that RandomResizeCrop is not favoured during the search phase (see supplementary), suggesting that it is not critical for ImageNet-100. Instead the performance gap between TA (Wide) and TA (RA) suggest that harder transformations are key to a better performance for this dataset.

### 4.3. Beyond natural images

For the DomainNet dataset, we compare SLACK to a Uniform policy, to the augmentations used by DomainBed [8] for domain generalization, and to the TrivialAugment (RA) and (Wide) methods with their ImageNet and CIFAR default settings. Results can be found in Table 3.

DomainBed uses the same default transformations as TA ImageNet together with Grayscale, but with smaller magnitudes and unlike TA, does not add a random transformation. Yet it strongly overfits and performs much lower than TA.

This suggests that augmentations well suited for domain generalization do not perform well on the individual tasks. TA (Wide) ImageNet consistently outperforms all other TA flavors. This further justifies the need to learn the magnitude range and to eliminate any manual range selection process.

SLACK is a close second, yet it learns the policy end-to-end. The learned policies are illustrated as pie charts in Fig. 3. The slices represent the probability  $\pi$  over the different transformations while their radius represent the corresponding magnitudes. They differ from a domain to another, suggesting that the gain compared to the initialization (*i.e.* Uniform policy) results from SLACK’s ability to learn and adapt to each domain.

### 4.4. Ablation study

In this section, we evaluate our contributions and main design choices: the network architecture used for search, the KL regularization, and the benefits of learning  $\pi$  and  $\mu$ . More ablations can be found in the supplementary. Note that hyperparameters are adjusted to each baseline included in the comparison, to make them as competitive as possible.

**Network architecture for search.** In prior works, the search phase (when there is one) is conducted on the smaller WideResNet-40x2 architecture for CIFAR10 and CIFAR100, and the learned policy is evaluated for both WideResNet-40x2 and WideResNet-28x10. Table 4 shows that for SLACK, searching directly with WideResNet-28x10 gives the best results for that architecture.

**KL regularization.** We compare SLACK with a flavor that does not apply KL-regularization. For the latter, we reduce the outer learning rate so that the augmentation policies with and without regularization evolve at similar speeds. Results in Table 5 show that our regularization is beneficial.



	# Augmentations		Real-50k	Quickdraw-50k	Infograph	Sketch	Painting	Clipart	Average
	Total	Hard-coded							
DomainBed <sup>†</sup>	5	5	62.54 ± .15	66.54 ± .91	26.76 ± .36	59.54 ± .37	58.31 ± .25	66.23 ± .10	57.23 ± .18
TA (RA) ImageNet <sup>†</sup>	5	4	70.85 ± .13	67.85 ± .07	<b>35.24 ± .19</b>	65.63 ± .11	64.75 ± .18	70.29 ± .18	62.43 ± .05
TA (Wide) ImageNet <sup>†</sup>	5	4	<b>71.56 ± .07</b>	68.60 ± .05	<b>35.44 ± .33</b>	<b>66.21 ± .16</b>	<b>65.15 ± .20</b>	71.19 ± .19	<b>63.03 ± .07</b>
TA (RA) CIFAR <sup>†</sup>	3	2	70.28 ± .08	68.35 ± .07	33.85 ± .21	64.13 ± .12	64.73 ± .17	70.33 ± .21	61.94 ± .05
TA (Wide) CIFAR <sup>†</sup>	3	2	71.12 ± .10	<b>69.29 ± .05</b>	34.21 ± .29	65.52 ± .25	64.81 ± .14	71.01 ± .21	62.66 ± .07
Uniform policy	3	0	70.37 ± .08	68.27 ± .06	34.11 ± .21	65.22 ± .17	63.97 ± .24	72.26 ± .14	62.37 ± .06
SLACK (ours)	3	0	71.00 ± .13	68.14 ± .11	34.78 ± .18	65.41 ± .16	64.83 ± .12	<b>72.65 ± .20</b>	62.80 ± .06

Table 3. Test accuracies on DomainNet.

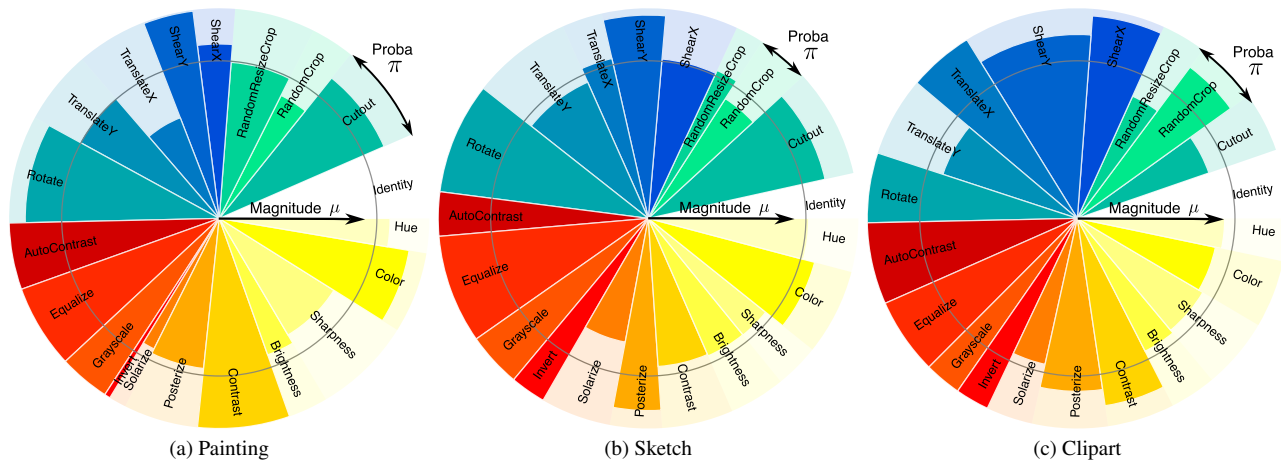


Figure 3. Policies found on DomainNet for the best search split. Gray circle: initial magnitude upper-bound. Radius of each pie: learned upper-bound. Size of each pie: probability of each transformation, averaged over the three composite distributions. Transformations which are parameter-free, *AutoContrast*, *Equalize*, *Grayscale*, and *Invert*, are displayed with maximal magnitude upper-bound.

Search architecture	CIFAR10	CIFAR100
WRN-40-2	97.43 ± .04	83.94 ± .20
WRN-28-10 (ours)	97.46 ± .06	84.08 ± .16

Table 4. CIFAR10/100 accuracy evaluated with WRN-28-10: impact of using a smaller architecture for the search phase.

SLACK variant	CIFAR10		CIFAR100	
	WRN-40-2	WRN-28-10	WRN-40-2	WRN-28-10
Uniform policy	96.22 ± .10	97.38 ± .05	79.07 ± .24	83.26 ± .17
$\mu$ only	96.20 ± .08	97.42 ± .05	79.22 ± .17	83.57 ± .18
$\pi$ only	96.22 ± .09	97.35 ± .04	79.36 ± .11	83.45 ± .15
$\pi$ and $\mu$ (ours)	96.29 ± .08	97.46 ± .06	79.87 ± .11	84.08 ± .16

Table 6. CIFAR10/100 accuracy when only learning part of the policy parameters

SLACK variant	CIFAR10		CIFAR100	
	WRN-40-2	WRN-28-10	WRN-40-2	WRN-28-10
without KL	96.27 ± .05	97.06 ± .11	79.61 ± .13	83.79 ± .19
with KL (ours)	96.29 ± .08	97.46 ± .06	79.87 ± .11	84.08 ± .16

Table 5. CIFAR10/100 accuracy with/without KL regularization.

**Joint learning of  $\pi$  and  $\mu$ .** Lastly, we study how beneficial jointly learning our augmentation parameters is compared to the initial Uniform policy and to a setting where only  $\pi$  or  $\mu$  is learned. Results can be found in Table 6.

## 5. Conclusion

In this paper, we address the task of automatic data augmentation. Considering the more challenging bilevel optimization problem that arises when the search space is not

reduced with default transformations, our proposed SLACK method tackles the resulting stability issues thanks to a multi-stage approach based on cold-start, coupled with a KL-regularization. Combined, they allow to reduce the variance of the gradient estimate and to better control the optimization process. We have experimentally observed that our method performs on par with recent approaches leveraging prior knowledge. It has also proved versatile enough to select domain-specific transformations when confronted to non-natural images.

**Acknowledgments.** This work was supported by ANR 3IA MIAI@Grenoble Alpes (ANR-19-P3IA-0003) and performed using HPC resources from GENCI-IDRIS (Grant 2022-AD011013343).

## References

- [1] Michael Arbel and Julien Mairal. Non-convex bilevel games with critical point selection maps. *preprint arXiv:2207.04888*, 2022. **1**
- [2] Ekin D. Cubuk, Barret Zoph, Dandelion Mané, Vijay Vasudevan, and Quoc V. Le. Autoaugment: Learning augmentation strategies from data. In *Proc. CVPR*, 2019. **3, 5, 6, 7, 10**
- [3] Ekin Dogus Cubuk, Barret Zoph, Jon Shlens, and Quoc Le. RandAugment: Practical automated data augmentation with a reduced search space. In *Proc. NeurIPS*, 2020. **3, 6, 7, 10**
- [4] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale hierarchical image database. In *Proc. CVPR*, 2009. **1**
- [5] Terrance DeVries and Graham W Taylor. Improved regularization of convolutional neural networks with cutout. *arXiv preprint arXiv:1708.04552*, 2017. **1**
- [6] Michael C Fu. Gradient estimation. *Handbooks in operations research and management science*, 13:575–616, 2006. **5**
- [7] Will Grathwohl, Dami Choi, Yuhuai Wu, Geoffrey Roeder, and David Duvenaud. Backpropagation through the void: Optimizing control variates for black-box gradient estimation. In *Proc. ICLR*, 2018. **1, 3**
- [8] Ishaan Gulrajani and David Lopez-Paz. In search of lost domain generalization. In *Proc. ICLR*, 2021. **7, 11**
- [9] Ryuichiro Hataya, Jan Zdenek, Kazuki Yoshizoe, and Hideki Nakayama. Faster AutoAugment: Learning augmentation strategies using backpropagation. In *Proc. ECCV*, 2020. **3, 5, 10**
- [10] Ryuichiro Hataya, Jan Zdenek, Kazuki Yoshizoe, and Hideki Nakayama. Meta approach to data augmentation optimization. In *Proc. WACV*, 2022. **1, 3**
- [11] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proc. CVPR*, 2016. **6**
- [12] Daniel Ho, Eric Liang, Xi Chen, Ion Stoica, and Pieter Abbeel. Population based augmentation: Efficient learning of augmentation policy schedules. In *Proc. ICML*, pages 2731–2741, 2019. **1, 3, 10**
- [13] Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. 2009. **1, 6**
- [14] Yonggang Li, Guosheng Hu, Yongtao Wang, Timothy M. Hospedales, Neil Martin Robertson, and Yongxin Yang. DADA: differentiable automatic data augmentation. In *Proc. ECCV*, 2020. **1, 3, 5, 6, 7, 10, 11**
- [15] Sungbin Lim, Ildoo Kim, Taesup Kim, Chiheon Kim, and Sungwoong Kim. Fast AutoAugment. In *Proc. NeurIPS*, 2019. **1, 3, 6, 7, 10, 11**
- [16] C. Lin, M. Guo, C. Li, X. Yuan, W. Wu, J. Yan, D. Lin, and W. Ouyang. Online hyper-parameter learning for auto-augmentation strategy. In *Proc. ICCV*, 2019. **1**
- [17] Tom Ching LingChen, Ava Khonsari, Amirreza Lashkari, Mina Rafi Nazari, Jaspreet Singh Sambee, and Mario A Nascimento. Uniformaugment: A search-free probabilistic data augmentation approach. *arXiv preprint arXiv:2003.14348*, 2020. **7**
- [18] Aoming Liu, Zehao Huang, Zhiwu Huang, and Naiyan Wang. Direct differentiable augmentation search. In *Proc. ICCV*, 2021. **3, 5**
- [19] Samuel G. Müller and Frank Hutter. TrivialAugment: tuning-free yet state-of-the-art data augmentation. In *Proc. ICCV*, 2021. **1, 2, 3, 6, 7, 10, 11, 12**
- [20] Mattis Paulin, Jérôme Revaud, Zaid Harchaoui, Florent Perronnin, and Cordelia Schmid. Transformation pursuit for image classification. In *Proc. CVPR*, 2014. **3, 6**
- [21] Xingchao Peng, Qinxun Bai, Xide Xia, Zijun Huang, Kate Saenko, and Bo Wang. Moment matching for multi-source domain adaptation. In *Proc. ICCV*, pages 1406–1415, 2019. **2, 6**
- [22] R Tyrrell Rockafellar. Monotone operators and the proximal point algorithm. *SIAM journal on control and optimization*, 14(5):877–898, 1976. **2**
- [23] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017. **2, 5**
- [24] Teppei Suzuki. Techaugment: Data augmentation optimization using teacher knowledge. In *Proc. CVPR*, 2022. **6, 7**
- [25] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive multiview coding. *arXiv preprint arXiv:1906.05849*, 2019. **6**
- [26] Riccardo Volpi and Vittorio Murino. Addressing model vulnerability to distributional shifts over image transformation sets. In *Proc. ICCV*, pages 7980–7989, 2019. **3**
- [27] Xiaoxing Wang, Xiangxiang Chu, Junchi Yan, and Xiaokang Yang. DAAS: Differentiable architecture and augmentation policy search. *arXiv preprint arXiv:2109.15273*, 2021. **1, 3, 5**
- [28] Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3):229–256, 1992. **1, 3**
- [29] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. In Edwin R. Hancock Richard C. Wilson and William A. P. Smith, editors, *Proc. BMVC*, pages 87.1–87.12. BMVA Press, September 2016. **6**
- [30] Xinyu Zhang, Qiang Wang, Jian Zhang, and Zhao Zhong. Adversarial AutoAugment. In *Proc. ICLR*, 2020. **3**
- [31] Yu Zheng, Zhi Zhang, Shen Yan, and Mi Zhang. Deep AutoAugmentation. In *Proc. ICLR*, 2022. **1, 3, 6, 7, 10, 11**

In this supplementary, we first give a short overview of prior work’s key aspects (Section A), then we provide additional details about the experimental protocol that was used to obtain our results (Section B) and analyse in more details the results reported in the paper and the impact of our approach on the stability of the search process (Section C).

## Contents

<b>1. Introduction</b>	<b>1</b>
<b>2. Related Work</b>	<b>3</b>
<b>3. Method</b>	<b>3</b>
3.1. Stochastic data augmentations policy . . . . .	3
3.2. Bilevel formulation for policy search . . . . .	4
3.3. SLACK algorithm . . . . .	4
<b>4. Experiments</b>	<b>6</b>
4.1. Experimental setup . . . . .	6
4.2. Comparison with the state of the art . . . . .	6
4.3. Beyond natural images . . . . .	7
4.4. Ablation study . . . . .	7
<b>5. Conclusion</b>	<b>8</b>
<b>A Summary of prior art</b>	<b>10</b>
<b>B Experimental setup</b>	<b>10</b>
B.1. Gradient estimation . . . . .	10
B.2. Magnitude ranges. . . . .	11
B.3. Image pre-processing . . . . .	11
B.4. Policy search . . . . .	11
B.5. Policy evaluation . . . . .	11
<b>C Extended analysis</b>	<b>11</b>
C.1. Uniform distribution: ablations on prior work	11
C.2. Visualization of the learned policies . . . . .	13
C.3. Avoiding instabilities with SLACK . . . . .	13
C.4. An ensembling approach . . . . .	16
C.5. Warm-start vs cold-start . . . . .	16

## A. Summary of prior art

Overall, methods for augmentation search mostly differ in key design choices that are highlighted in Table 7.

## B. Experimental setup

In this section, we first describe the practical implementation of our gradient estimation, then we compare our magnitude ranges to those of other methods, and finally we describe in more detail our search and evaluation protocols.

Method	Optimization	Prior-free*	Full set	Search time
AA [2]	Reinfor. learning	✗	✗	5000
PBA [12]	Population-based	✗	✗	5
Fast AA [15]	Bayesian optim.	✗	✗	3.5
Faster AA [9]	RELAX	✗	✗	0.23
DADA [14]	RELAX	✗	✗	0.1
RA [3]	Exhaustive search	✗	✓	25
TrivialA [19]	No optimization	✗	✓	NA
DeepAA [31]	Greedy algorithm	✓**	✓***	9
<b>SLACK (Ours)</b>	REINFORCE	✓	✓	4

\* Prior-free refers to methods that do not use any default transformations.

\*\* DeepAA does not use default transformations during search but during pretraining.

\*\*\* DeepAA pretrains on a reduced dataset.

Table 7. **Key aspects of automatic data augmentation methods.**

Our approach, SLACK, tackles the corresponding bilevel optimization problem using the REINFORCE gradient estimator. It is prior-free (*i.e.* does not rely on default transformations) and can use the full training set while maintaining a reasonable search time (indicated in GPU hours, for search on CIFAR with WRN-40-2).

## B.1. Gradient estimation

In this section, we describe the practical implementation of the gradient estimates derived in our method section.

To optimize the augmentation policies, we minimize an approximation to the upper-level objective  $\mathcal{F}(\phi) := \mathcal{L}_{\text{val}}(\theta^*(\phi))$  defined as  $\hat{\mathcal{F}}(\phi) := \mathcal{L}_{\text{val}}(\hat{\theta}(\phi))$ , where we replaced the intractable lower-level solution  $\theta^*(\phi)$  by an approximate solution  $\hat{\theta}(\phi)$ . Such approximate solution is obtained by performing one gradient step to optimize the lower-level objective starting from the current parameter  $\theta$ , *i.e.*  $\hat{\theta}(\phi) := \theta - \eta \nabla_{\theta} \mathcal{L}_{\text{train}}(\theta, \phi)$ . The gradient  $\nabla_{\phi} \mathcal{F}$  is then naturally approximated by  $\nabla_{\phi} \hat{\mathcal{F}}$  which is computed by applying the chain rule:

$$\nabla_{\phi} \mathcal{F} \approx \nabla_{\phi} \hat{\mathcal{F}} = \nabla_{\theta} \mathcal{L}_{\text{val}}(\hat{\theta}(\phi))^{\top} \nabla_{\phi} \hat{\theta}(\phi).$$

The Jacobian  $\nabla_{\phi} \hat{\mathcal{F}}$  can be computed explicitly using the Score method which yields:

$$\nabla_{\phi} \hat{\theta}(\phi) = -\eta \mathbb{E}_{\tau \sim p_{\phi}} [\nabla_{\theta} \ell_{\text{train}}(\theta, \tau) \nabla_{\phi} \log p_{\phi}(\tau)^{\top}].$$

In practice, expectations over the data and augmentation policies are estimated with batches. At a given iteration, we sample  $B_{\text{aug}}$  augmentations from  $p_{\phi}$  and then apply each of them to a batch of training data  $B_{\text{train}}$  from  $\mathcal{D}_{\text{train}}$  to approximate  $\nabla_{\phi} \hat{\theta}(\phi)$ . Finally, we use a batch  $B_{\text{val}}$  of data from  $\mathcal{D}_{\text{val}}$  to estimate the validation loss. Denoting  $N_a, N_t, N_v$  the size of the augmentation, training and validation batches respectively, and

$$\hat{l}_{\text{val}}(\theta) := \frac{1}{N_v} \sum_{(x,y) \in B_{\text{val}}} [\ell(y, f_{\theta}(x))],$$

$$\hat{l}_{\text{train}}(\theta, \tau) := \frac{1}{N_t} \sum_{(x,y) \in B_{\text{train}}} [\ell(y, f_{\theta}(\tau(x)))],$$

our gradient estimate can be expressed as

$$\begin{aligned} \nabla_{\phi} \mathcal{F} &\approx -\frac{\eta}{N_a} \nabla_{\hat{\theta}} \hat{l}_{\text{val}}(\theta) \left( \sum_{\tau \in B_{\text{aug}}} \nabla_{\theta} \hat{l}_{\text{train}}(\theta, \tau) \nabla_{\phi} \log p_{\phi}(\tau)^T \right) \\ &= -\frac{\eta}{N_a} \sum_{\tau \in B_{\text{aug}}} \left( \nabla_{\theta} \hat{l}_{\text{val}}(\theta)^T \nabla_{\theta} \hat{l}_{\text{train}}(\theta, \tau) \right) \nabla_{\phi} \log p_{\phi}(\tau) \end{aligned}$$

In other words, the upper-level gradient is a weighted sum of the scores  $\nabla_{\phi} \log p_{\phi}(\tau)$ , with the weights representing the alignment between the gradients of i) the loss on the training data transformed with  $\tau$  (evaluated at  $\theta$ ), and ii) the loss on the validation data (evaluated at  $\hat{\theta}$ , *i.e.* one step ahead).

In practice, the lower-level learning rate decreases with a cosine schedule. As we do not want our upper-level gradient updates to shrink, we set  $\eta$  to the *initial* value of the lower-level learning rate instead of its current value.

## B.2. Magnitude ranges.

The ranges used for mapping the magnitudes to  $[0,1]$  vary across methods; we indicate this mapping for each method in Table 8. For transformations with respect to which the datasets naturally exhibit symmetries (Shear, Translate, Rotate, Enhance), once we have sampled a magnitude, we randomly select a direction. Note that SLACK’s ranges are larger than the usual ones (*i.e.* those of TA (RA)), which gives more flexibility during the optimization of our magnitude upper-bounds  $\mu$ . The latter is initialized at 0.75. Experimentally, we noted that this initialization should be high enough to favour exploration and avoid over-fitting during pre-training. We observed that any initialization in the  $[0.75, 0.9]$  range consistently works well across datasets.

## B.3. Image pre-processing

Table 9 indicates the image pre-processing choices on ImageNet-100 and DomainNet for TrivialAugment [19], DomainBed [8] and SLACK.

ImageNet-100 and DomainNet images have variable original sizes. In the literature, training images are commonly resized with RandomResizeCrop. For testing, TrivialAugment uses `Resize(256)+CenterCrop((224,224))`, preserving the aspect ratio, while DomainBed directly applies `Resize((224,224))`, degrading the aspect ratio but preserving the image content. For each method, we stick to the authors’ choices, as we experimentally noted that they yield the best results (*e.g.* using TrivialAugment’s pre-processing for DomainBed degrades the performance, and vice-versa).

For SLACK, which does not apply RandomResizeCrop by default, we preprocess the training data and validation/testing data in the same way. For training, random cropping is applied instead of center cropping to fully exploit the data. For ImageNet-100, we use TrivialAug-

ment’s pre-processing. For DomainNet, we select the pre-processing strategy by cross-validation after pre-training.

## B.4. Policy search

Hyperparameters used for policy search are indicated in Table 10. They are chosen to satisfy two criteria that we found to be useful for obtaining a successful policy search: i) the validation loss after re-training should be similar (experimentally, slightly lower) to the one obtained after pre-training, and ii) the probability distributions should vary at the same speed for all datasets. Our learning rate is 4 times larger for re-training on CIFAR10 than on CIFAR100. We observed that gradients on CIFAR10 are 4 times smaller in norm than those on CIFAR100, and that re-scaling the updates allows satisfying i) and ii) empirically. For DomainNet, we adapt the number of re-training steps to the dataset size. A fixed lower-level learning rate for all datasets experimentally satisfies i). We observed that the lower-level gradients differ in scale for each dataset. Satisfying ii) requires re-scaling the KL regularisation and accordingly changing the upper-level learning rate (so that  $\text{KL weight} \times \text{upper-level lr}$  is constant).

The upper-level learning rate indicated in the Tables is the one used for updating  $\pi$ . We divide it by 40 for the optimization of  $\mu$  to ensure slower updates for the magnitude parameter which we found to be sensitive to variations (or by 10 for ablations removing the KL regularization).

## B.5. Policy evaluation

The hyperparameters used for the evaluation phase are indicated in Table 11. For CIFAR10 and CIFAR100, we use the same hyperparameters as prior work.

## C. Extended analysis

In this section, we further analyse the results of our search algorithm. We first illustrate the policies learned for all datasets. We then study in more detail the impact of our multi-stage approach with KL regularization, comparing it with single-stage (unrolled) and unregularized approaches and illustrating their instability.

### C.1. Uniform distribution: ablations on prior work

In this section, we motivate our choice of a uniform magnitude distribution, showing that it globally outperforms optimized magnitude models in prior work. To this end, we directly evaluate the policies provided by the authors without re-running their search procedure. We compare their learned magnitude model with a simpler one that consists in sampling the magnitudes uniformly on their  $[0, 1]$ -mapped ranges. We study three baselines: DADA [14], FastAA [15] and DeepAA [31]. Note that their policy results from a search on CIFAR10, that they also use when evaluating on CIFAR100. Results are reported in Table 12.

Application	Transformation	Method			
		TA (RA)	TA (Wide)	DomainBed	Ours
Sampled	ShearX/Y	[0, 0.3]	[0, 0.99]	-	[0, 1]
	Translate X/Y	[0, 0.45]*	[0, 32px]**	-	[0, 0.75]
	Rotate	[0, 30]	[0, 135]	-	[0, 90]
	Posterize	[4, 8]	[2, 8]	-	[2, 8]
	Solarize	[0, 255]	[0, 255]	-	[0, 255]
	Enhance***	[0, 0.9]	[0, 0.99]	-	[0, 0.99]
	Cutout	[0, 0.2]	[0, 0.6]	-	[0, 1]
	RandCrop	-	-	-	[0, 0.5]
	RandResizeCrop	-	-	-	[0.05, 1]
Default	ColorJitter****	[0, 0.4]	[0, 0.4]	[0, 0.3]	-
ImageNet/DomainNet	RandResizeCrop	[0.08, 1]	[0.08, 1]	[0.7, 1]	-
Default CIFAR	Cutout	0.5	0.5	NA	-
	RandCrop	0.125	0.125	NA	-

Table 8. Our magnitude ranges compared to those used by other methods.

\* TrivialAugment [19] uses  $[-0.31, 0.31]$ ,

\*\* TrivialAugment [19] sets the upper-bounds in pixels, not in proportion

\*\*\* Color, Contrast, Brightness, Sharpness,

\*\*\*\* Color, Contrast, Brightness

Dataset	Model	Train	Test
ImageNet-100	TrivialAugment	RandResizeCrop((224,224))	Resize(256)+CenterCrop((224,224))
	SLACK	Resize(256)+RandomCro((224,224))	Resize(256)+CenterCrop((224,224))
DomainNet	TrivialAugment ImageNet	RandResizeCrop(224,224)	Resize(256)+CenterCrop(224,224)
	TrivialAugment CIFAR	Resize(256)+RandomCrop((224,224),padding=28)	Resize(256)+CenterCrop((224,224))
	DomainBed	RandResizeCrop((224,224))	Resize((224,224))
	SLACK (Clipart, Sketch, Quickdraw)	Resize((224,224))	Resize((224,224))
	SLACK (Painting, Infograph, Real)	Resize(256)+RandomCrop((224,224))	Resize(256)+CenterCrop((224,224))

Table 9. Image pre-processing on ImageNet-100 and DomainNet.

Dataset	Network	Re-train iter	Unrolled iter	Batch size	Lower lr	Upper lr	KL weight $\times$ Upper lr
CIFAR10/100	WRN-40-2/WRN-28-10	1000	400	$8 \times 128$	0.4 / 0.1	1	0.02
ImageNet-100	ResNet-18	2000	800	$8 \times 256$	0.1	0.5	0.005
DomainNet	ResNet-18	800 - 1200	400	$8 \times 128$	0.1	0.625 - 1.25	0.01

Table 10. Hyperparameters used for search on CIFAR, ImageNet-100 and DomainNet.

Dataset	Network	Epochs	Batch size	Learning rate	Weight decay
CIFAR10/100	WRN-40-2, WRN-28x10	200	128	0.1	0.0005*
ImageNet-100	ResNet-18	270	256	0.1	0.001
DomainNet	ResNet-18	200	128	0.1	0.001

Table 11. Hyperparameters used for training on CIFAR, ImageNet-100 and DomainNet. \*: 0.0002 for CIFAR10 on WRN-40-2

**Parametrization.** DADA and FastAA directly optimize a probability distribution over the set of all possible composite transformations (*sub-policies*) and learn a single magnitude value for each transformation in a sub-policy. They keep the top- $k$  sub-policies for evaluation. DeepAA learns to compose transformations in a greedy manner and discretizes the magnitude ranges, learning a probability for each magnitude. We compare these learned magnitude values (FastAA, DADA) or learned probabilities (DeepAA) to our approach based on a uniform sampling.

**DADA/FastAA.** Our approach compares favorably to DADA’s and FastAA’s optimized models. We also compare both approaches on their initial policy (equal probabilities for all sub-policies, magnitudes set at mid-range). With uniform magnitude sampling, their initial policy (sampling among all possible sub-policies) performs similarly if not better than their optimized one (sampling among their top- $k$  sub-policies).

**DeepAA.** Results on the policy provided by DeepAA are more nuanced: using uniform sampling improves results

Model	Magnitude model	CIFAR10		CIFAR100	
		WRN-40-2	WRN-28-10	WRN-40-2	WRN-28-10
FastAA/DADA initialization	Theirs	96.22	97.08	78.26	82.17
	Uniform	96.37	97.25	79.10	82.80
FastAA, reported	Original	96.4	97.3	79.3	82.7
FastAA, reproduced (evaluation only)	Original	96.4	97.22	79.11	82.82
	Uniform	96.37	97.30	79.15	82.84
DADA, reported	Original	96.4	97.3	79.1	82.5
DADA, reproduced (evaluation only)	Original	96.33	97.19	79.07	82.05
	Uniform	96.37	97.35	78.97	82.57
DeepAA, reported	Original	-	97.56	-	84.02
DeepAA, reproduced (evaluation only)	Original	96.46	97.48	79.62	83.85
	Uniform	96.55	97.47	78.89	83.62

Table 12. **Why learning the magnitude range?** CIFAR10/100 accuracies for DADA, FastAA and DeepAA, when using their original magnitude model, or a simpler one which samples in a uniform manner in their ranges.

on CIFAR10 (on which their search was conducted) and degrades them on CIFAR100.

## C.2. Visualization of the learned policies

**CIFAR.** The evolution of probability distributions for CIFAR10 and CIFAR100 and pie charts of the final policies are illustrated in Fig. 4. It can be noted that Invert and Solarize, known to be detrimental, are systematically discarded. The policies learned are quite diverse, with different leading transformations for each distribution but global predominance of some transformations such as Cutout or Rotate. It can also be noted that magnitudes upper-bounds are in average higher for the larger WideResNet-28x10 networks: a larger learning capacity benefits more from harder transformations.

**ImageNet-100.** The best policy found for ImageNet-100 is illustrated in Fig. 5. Interestingly, RandomResizeCrop is ranked quite low, yet our policy yields results comparable to TrivialAugment’s (with a 86.18 average accuracy on this split), suggesting that other geometrical transformations such as Cutout, Rotate and ShearY are equivalently beneficial for training on ImageNet-100. We can note rather high magnitudes upper-bounds for the color jittering transformation TrivialAugment also applies by default (Color, Contrast, Brightness), which is consistent with the higher performance of TrivialAugment’s (Wide) version compared to (RA).

**DomainNet.** The policies found by SLACK on DomainNet are illustrated in the pie charts Fig. 6 for all domains. Some similarities with policies found on CIFAR and ImageNet can be noted. In particular, Invert and Solarize (that only inverts part of the pixels) are systematically discarded for all domains except Quickdraw. Invert is manually removed from TrivialAugment’s baseline as it is known to be

detrimental, and this seems to generalize to other domains. Also, Rotate and Cutout are globally favoured, similarly to the policies found on CIFAR and ImageNet-100.

However some differences mark specificities to each domain: i) on the strength of the transformations: for example, geometrical transformations are given high magnitudes on Clipart and lower ones on Real, ii) on their probabilities: color jittering transformations used for real images are globally assigned a high probability for Real, Painting and Infograph domains, and a much lower one for Clipart, which suggests that changes in color, contrast or brightness are less meaningful for this domain.

## C.3. Avoiding instabilities with SLACK

In this section, we study how our augmentation policies evolve when removing the KL regularization or when using a single optimization stage instead of multiple ones, which corresponds to standard unrolled optimization. Evaluations in both settings are reported in Table 13.

We first show that unrolled optimization is globally unstable and easily collapses, justifying the need for a regularization. We illustrate how entropy regularization prevents collapse and yields competitive results, but at the cost of high ‘local’ instability. These instabilities make the final performance highly dependent on the choice of some hyperparameters, such as the learning rate. The necessity to overcome these instabilities motivates our multi-stage procedure with an adaptive anchoring for the regularization. Lastly, we show that unregularized multi-stage optimization, while more stable than unregularized unrolled optimization, does not yield competitive results, confirming again the benefits of our KL regularization.

**Unregularized unrolled optimization.** Unrolled optimization is subject to two sources of instability: first, the approx-

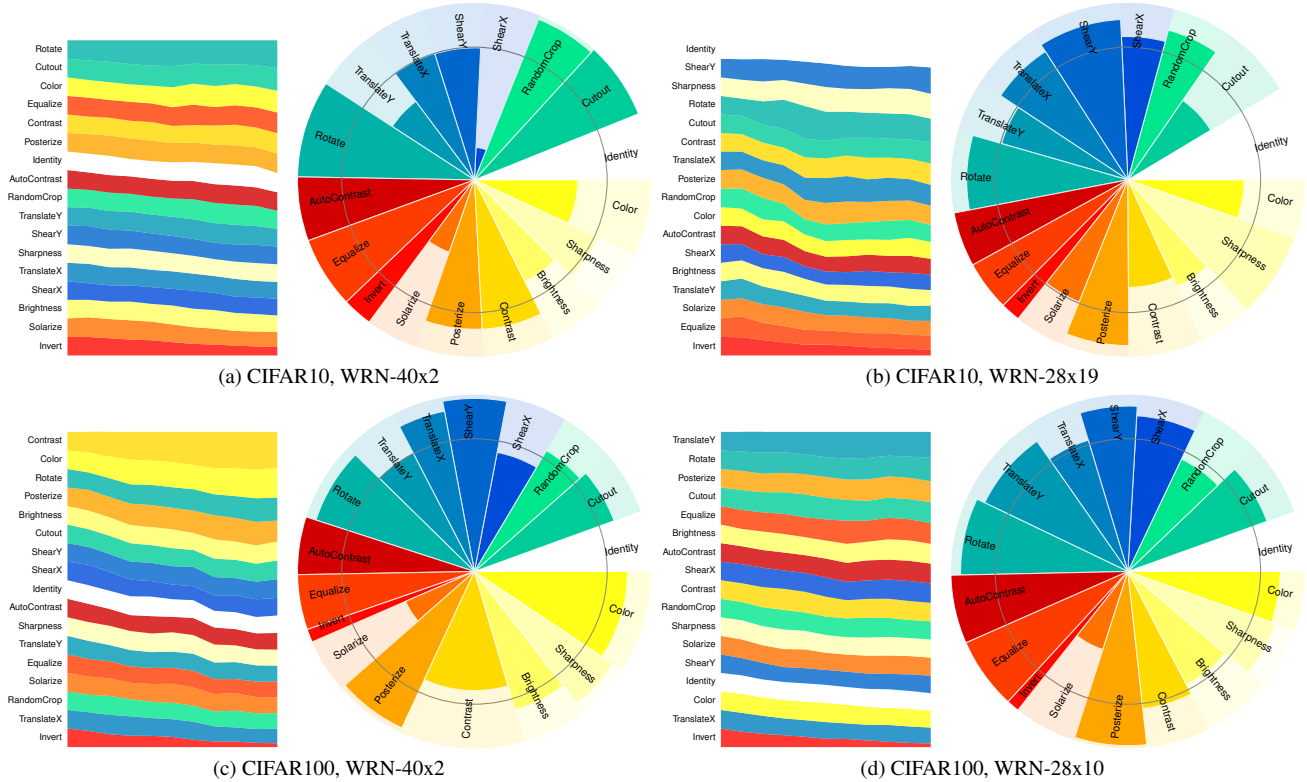


Figure 4. Illustrations of best policies found for CIFAR10 and CIFAR100 on WideResNet-40x2 and WideResNet-28x10 architectures. For each dataset and architecture, we show the evolution of the probability distribution  $\pi$  as training progresses (left) and the final learned policy as a pie chart (right), where slice widths represent  $\pi$  and slice radii represent  $\mu$ .

SLACK variant	Upper-level iterations	Upper-level lr	KL weight	CIFAR10		CIFAR100	
				WRN-40-2	WRN-28-10	WRN-40-2	WRN-28-10
Unrolled w/ KL (Fig. 8)	10000	0.25	0.005	96.30 ± .08	97.43 ± .04	79.54 ± .20	84.11 ± .13
SLACK w/o KL (Fig. 9)	10 × 400	0.25	0	96.27 ± .05	97.06 ± .11	79.61 ± .13	83.79 ± .19
SLACK (Fig. 4)	10 × 400	1	0.02	96.29 ± .08	97.46 ± .06	79.87 ± .11	84.08 ± .16

Table 13. CIFAR10/100 accuracy with unregularized and single-stage approaches.

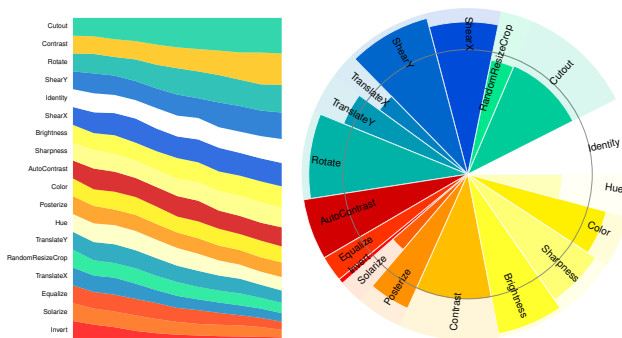


Figure 5. ImageNet-100 policy on the best search split

imation  $\theta^*(\phi) = \hat{\theta}(\phi)$  with a single gradient step inherently leads to wrong gradient updates; second, the REINFORCE gradient estimation is theoretically exact but has a high variance in practice when approximated in the context of stochastic optimization. Fig. 7 illustrates these instabilities: blindly following wrong gradient directions exacerbated by an oversampling of the dominant transformation leads to a progressive collapse of the policy.

**Unrolled optimization with entropy regularization.** In the case of a single-stage unrolled optimization, the KL regularization uses a uniform distribution as an anchor, which corresponds to an entropy regularization. By maximizing the entropy, the algorithm encourages exploration of the augmentation policies and prevents the divergence phenomenon observed above. While this regularization leads

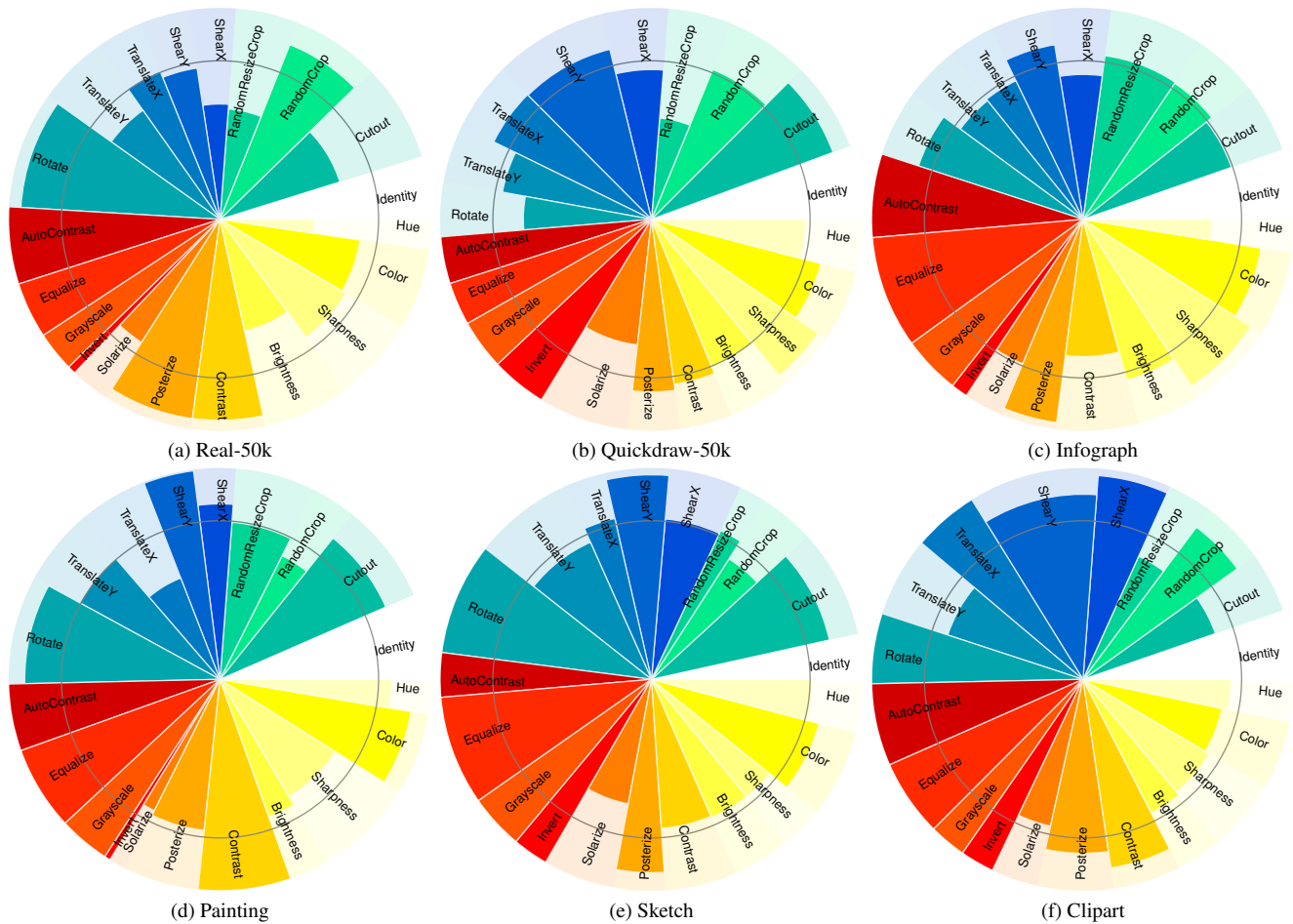


Figure 6. Policies found on DomainNet. The three distributions from  $\pi$  forming the composite transformation are averaged.

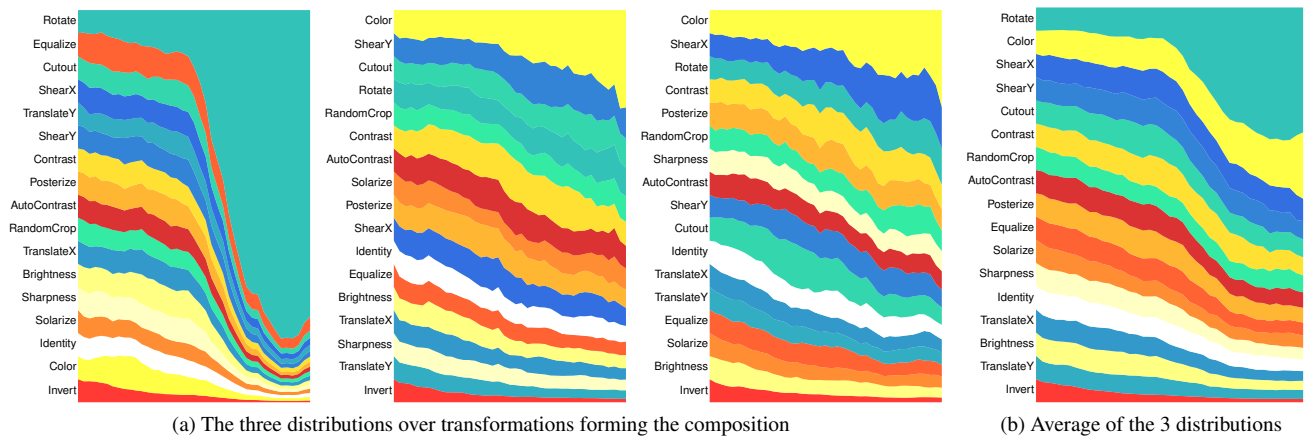


Figure 7. Evolution of the probability distributions  $\pi$  for CIFAR100 with unregularized unrolled optimization in a case of collapse.

to competitive results as reported in Table 13, it does not mitigate the inherent instability of the gradient updates. On the other hand, the multi-stage algorithm we proposed in

SLACK yields more stable gradient updates.

**Multi-stage optimization without KL regularization.** In our multi-stage approach,  $\theta^*(\check{\phi})$  is well-approximated at the



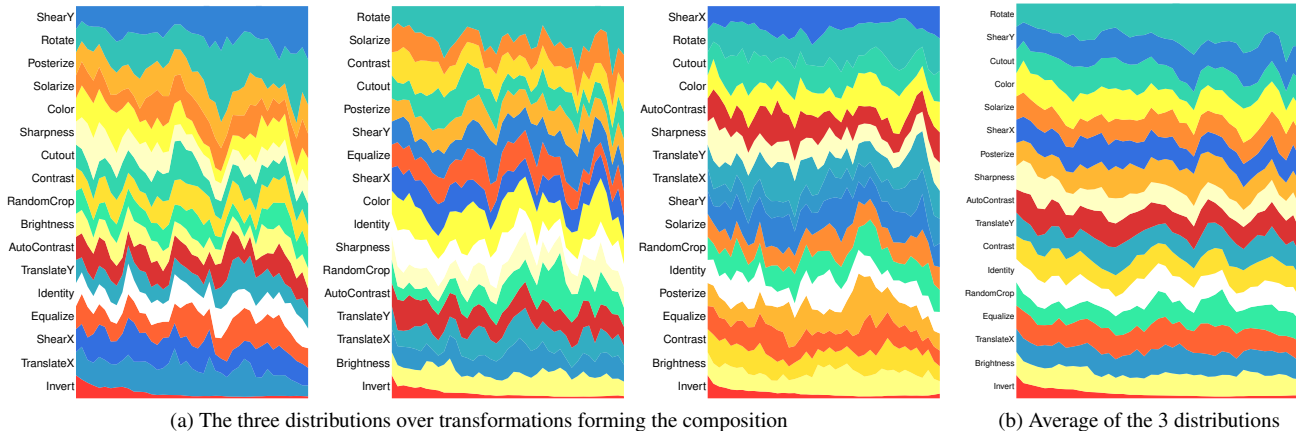


Figure 8. Evolution of the distributions  $\pi$  for CIFAR100 with entropy-regularized unrolled optimization, on one of the search splits.

beginning of each stage, as the model is re-trained with the current policy  $\tilde{\phi}$ . Gradient updates close to this policy are ‘trusted’ since our current  $\theta$  after re-training stays close to  $\theta^*(\tilde{\phi})$ , meaning that we strongly mitigate the approximation inherent to unrolled optimization. Our KL regularization encourages the policy to stay in this *trust region* and without it, the stochasticity of the optimization combined with the high variance from REINFORCE may drive the policy away. Fig. 9 shows the evolution of our probability distributions under an unregularized multi-stage search with two different learning rates, one twice larger than the other. This evolution is smoother than with single-stage unrolled optimization and also quite stable when using a small learning rate, but this slows down convergence, yielding a sub-optimal policy (see Table 13). The larger one leads again to a progressive divergence: the policy is driven too far and the  $\theta^*(\tilde{\phi})$  obtained after re-training becomes sub-optimal for the current  $\phi$ . In other words, the KL regularization allows making large updates in the parameter space, while remaining close to a reference/anchor policy.

#### C.4. An ensembling approach

In this section, we investigate the effect of an ensembling strategy for SLACK to reduce the variance of gradient updates in the search phase. More precisely, the strategy consists in independently training multiple models on the lower-level loss while averaging their contributions to the upper-level gradient. Each model is initialized (and subsequently re-initialized at each stage) based on a pre-training with a different seed. This ensembling strategy was implemented using multiple GPUs, where each GPU trains one copy of the model and only the upper-levels gradients are communicated and averaged across GPUs.

Results on CIFAR10/100 are reported in Table 14. While there is a small improvement in most cases, the method still

has a strong computational overhead. Yet it might be a relevant line of research for datasets for which the training procedure has a higher variance, *e.g.* smaller datasets where the additional cost of ensembling is not a significant overhead.

#### C.5. Warm-start vs cold-start

In this section, we study the model behaviour when searching with warm-start instead of cold-start. By warm-start, we mean that re-training is performed starting from the current network’s weights at the beginning of each stage instead of re-initializing it to its pre-trained weights. We experimentally observe that warm-start with the same hyper-parameters as for cold-start leads to a progressive overfitting of the network. Increasing the lower-level learning rate mitigates this phenomenon, but still yields sub-optimal results as reported in Table 15. This suggests that re-training from  $\theta_0^*(\phi_0)$  gives a better estimate of  $\theta^*(\phi_i)$  at stage  $i$  than re-training from the biased state close to  $\theta^*(\phi_{i-1})$ .

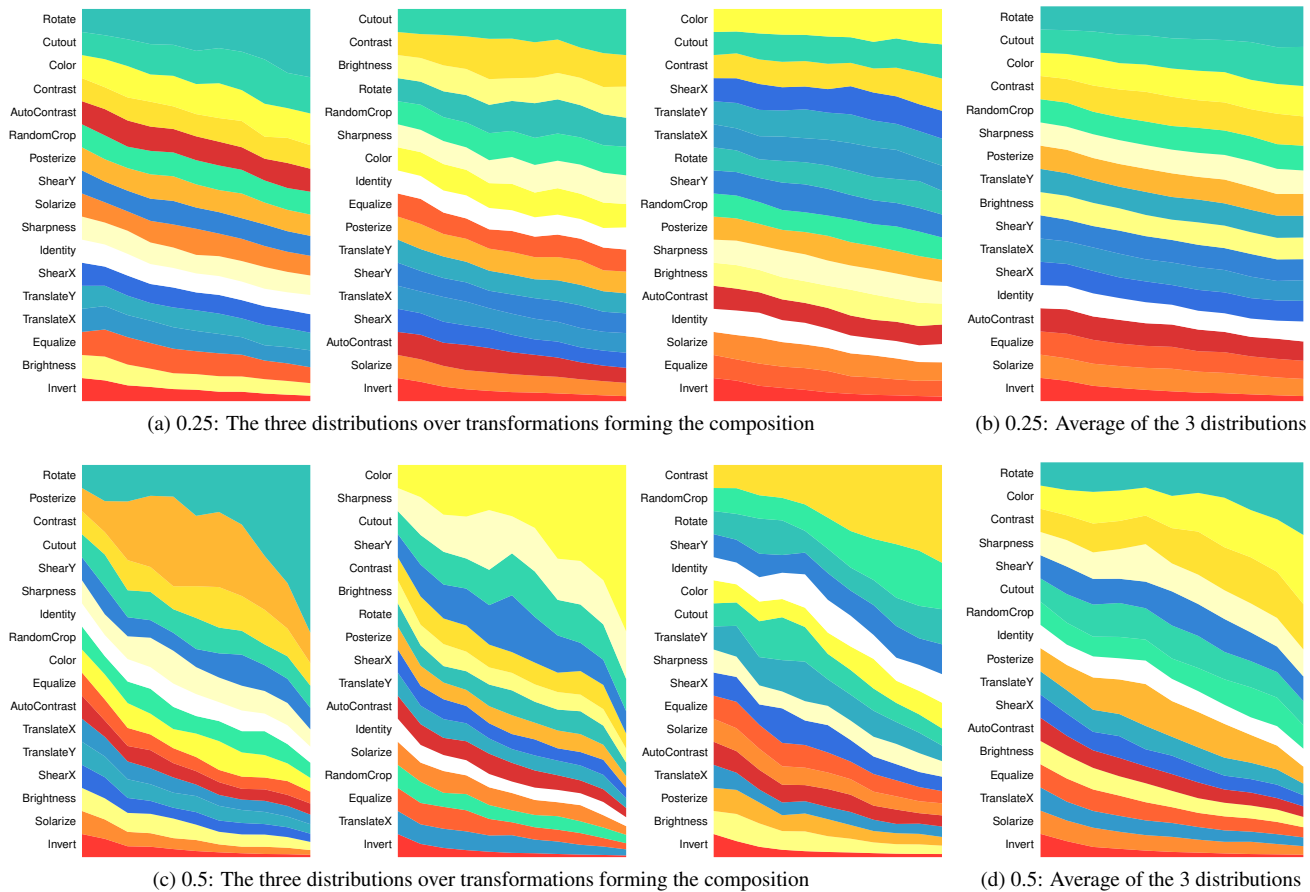


Figure 9. Evolution of the probability distributions  $\pi$  for CIFAR100 with unregularized multi-stage optimization using upper-level learning rates of 0.25 and 0.5.

	CIFAR10		CIFAR100	
	WRN-40-2	WRN-28-10	WRN-40-2	WRN-28-10
SLACK	96.29 ± .08	97.46 ± .06	79.87 ± .11	84.08 ± .16
Ensembling of SLACK (4 GPUs)	96.33 ± .08	97.48 ± .06	79.94 ± .13	84.01 ± .14

Table 14. CIFAR10/100 accuracy with ensembling strategy.

SLACK variant	CIFAR10		CIFAR100	
	WRN-40-2	WRN-28-10	WRN-40-2	WRN-28-10
Warm-start	96.27 ± .09	97.05 ± .15	79.70 ± .11	83.90 ± .10
Cold-start (ours)	96.29 ± .08	97.46 ± .06	79.87 ± .11	84.08 ± .16

Table 15. CIFAR10/100 accuracy with cold start and warm start.