

# 1S-DAUG: ONE-SHOT DATA AUGMENTATION FOR ROBUST FEW-SHOT GENERALIZATION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Few-shot learning (FSL) demands generalization to novel classes based on just a few shots of labeled examples, a setting where traditional test-time augmentations fail to be effective. We introduce 1S-DAug, a one-shot generative augmentation operator that synthesizes diverse yet faithful variants from just one example image at test time. 1S-DAug couples traditional geometric perturbations with controlled noise injection and a denoising diffusion process conditioned on the original image. The generated images are then encoded and aggregated, alongside the original image, into a combined representation for more robust FSL predictions. Integrated as a training-free model-agnostic plugin, 1S-DAug consistently improves FSL across standard benchmarks of 4 different datasets without any model parameter update, including achieving 10% proportional accuracy improvement on the miniImagenet 5-way-1-shot benchmark. Codes will be released.

## 1 INTRODUCTION

Few-shot learning (FSL) is important for recognition systems deployed in the wild. While deep neural networks attain strong performance given abundant supervision, their accuracy degrades in rare-case generalization (Wang et al., 2020). Real-world data are long-tailed; rare categories with limited labels persist and cap overall system performance even as head classes continue to grow (Kang et al., 2020). Scarcity in the target domain induces a train-test gap that manifests as high generalization error on novel classes at test time (Wang et al., 2020).

FSL is a concrete instance of this long-tail problem. The FSL model must assign labels to previously unseen classes using only a handful of labeled examples per class (Vinyals et al., 2016; Snell et al., 2017). This low-data regime appears in practical settings such as medical imaging for rare diseases and autonomous driving with open-world, unpredictable events (Liu & Feng, 2024). The central challenge is to achieve robust generalization under strict label scarcity and distribution shift.

Data augmentation offers a natural handle on this challenge (Dvornik et al., 2019; Zhou, 2012). From an ensemble perspective, test-time augmentation aggregates predictions across multiple views of the same input and averages out error; in an idealized independence thought experiment, a single-view error rate  $\varepsilon$  would combine as  $\varepsilon^m$  for  $m$  views, while in practice correlation attenuates but does not eliminate the benefit (Zhou, 2012). From a margin perspective, classical generalization bounds can relate test error to data radius; augmentations that contract effective data radii can lower the Rademacher complexity and tighten such bounds (Bartlett & Mendelson, 2002; Bai et al., 2025). On the training side, augmentation increases observed data amount and can tighten the generalization bound further (Hariharan & Girshick, 2017b; Schwartz et al., 2018).

However, the augmentation is effective only with both diversity and accuracy (Zhou, 2012). Standard geometric or photometric transforms like cropping, resizing and scaling often add limited new information and may degrade image quality (Bai et al., 2025). In FSL, where each example carries high influence, degradations are especially harmful at test time, and the model must rely on precise visual cues. Achieving high diversity while preserving class-defining content is therefore central.

Generative data augmentation has potential, but deploying it under FSL constraints is nontrivial. Image-to-image translation with adversarial training (e.g., FUNIT, CycleGAN) can be effective in specific domains, yet it is prone to training instabilities and to inconsistent quality across dissimilar object categories and poses (Liu et al., 2019a; Zhu et al., 2017). Prior attempts to employ GANs

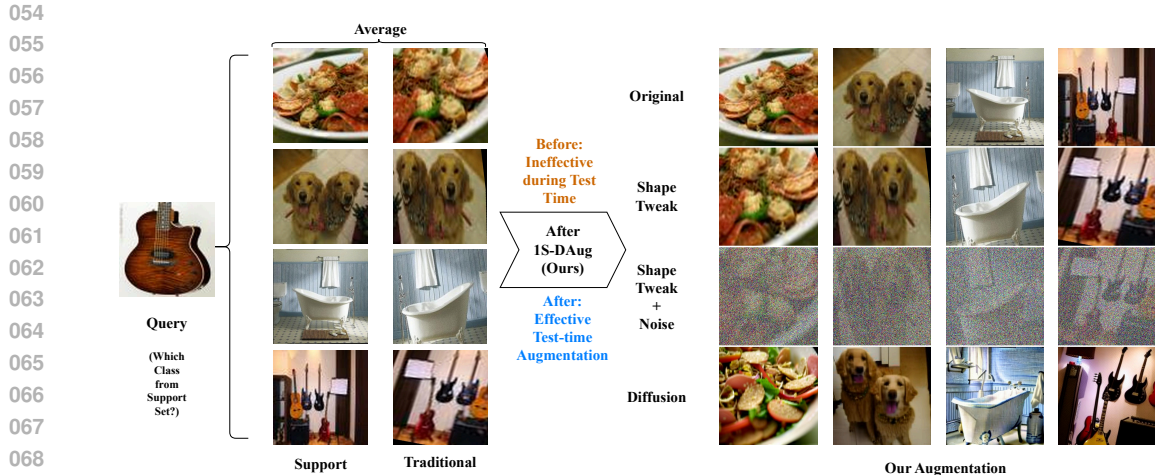


Figure 1: Pipeline of our one-shot, test-time augmentation. Given a query and supports, we apply a shape tweak and controlled noise, then perform attention-conditioned diffusion to synthesize class-faithful variants. Features from the original and the generated views are averaged before the few-shot head. For contrast, traditional test-time geometric transformations provide limited diversity.

for FSL augmentation typically restrict themselves to training-time augmentation (Hariharan & Girshick, 2017b; Schwartz et al., 2018; Hariharan & Girshick, 2017a). While generating data for seen classes is common, quality synthesis for unseen classes is challenging. To our best knowledge, Bai et al. (2025) is the only prior work explicitly using test-time generative augmentation for FSL (Liang et al., 2024), which relies on gan-based image-to-image translation comprising two samples and focuses on animal faces instead of more-diverse objects, limited in the broad dataset applicability (Liu et al., 2019a). A few recent work leverages diffusion-generated images for training data augmentation. However, the methods rely on fine-tuning using text prompts or additional target-class samples (Trabucco et al., 2024; He et al., 2023), thus not suitable for the challenging FSL test-time set-up.

We propose *IS-DAug* (One-Shot Data Augmentation), a one-shot generative operator that, given a single input image at test time, synthesizes a set of diverse yet faithful variants. The operator combines three ingredients: (i) class-preserving shape perturbations to broaden pose and layout coverage; (ii) controlled noise injection to hide non-class-defining and distorted details; and (iii) denoising guided by attention-based conditioning on the original image so that content-defining cues are preserved while appearance and pose vary in a controlled manner. During few-shot classification, the generated variants are processed by a fixed encoder and aggregated with the original into a single representation. The method is model-agnostic and requires no retraining. As large pretrained models proliferate and model sizes continue to explode, fine-tuning or model-side ensemble can be burdensome or infeasible due to compute or **restricted parameter access** (Zheng et al., 2023). A test-time, data-side plugin that treats the predictor in model-agnostic ways is thus attractive. In experiments on four FSL benchmarks (e.g., miniImagenet 5-way-1-shot), our method yields consistent gains without modifying any model parameters. Our contributions are as follows:

- We introduce *IS-DAug*, a one-shot data augmentation method. Given a single image, it leverages attention-guided diffusion with controlled noise to generate faithful yet diverse variants suitable for inference on unseen classes. Though synthetic data augmentation with target-class supervision is not new, to our best knowledge, we are the **first** to abide by the strict one-shot set-up for synthetic data augmentation, not relying on any label.
- Since our augmentation is strictly one-shot during inference, it can fulfill the challenging FSL test-time augmentation requirement. As FSL test-time augmentation is practically meaningful but **underexplored**, we analyse *IS-DAug*'s integration with FSL, where diversity and accuracy should be both achieved for performance optimality.
- We implement *IS-DAug* as a test-time plugin for trained FSL models, and we evaluate on benchmarks of four datasets, observing **consistent gains** (e.g., up to 20% relative improvement on 5-way-1-shot benchmarks) without any underlying parameter update to the off-the-shelf models.

## 2 RELATED WORK

For **few-shot learning**, augmentation-based approaches expand training diversity via feature hallucination (Hariharan & Girshick, 2017a) or GAN-driven synthesis (Wang et al., 2018), but they are largely *training-time* and rely on supervision from base classes. In contrast, *test-time* augmentation for FSL is harder; it must generate high-quality, class-faithful variants for unseen classes without any retraining or labels. The only explicit work in this direction is FSL-Rectifier (Bai et al., 2025), which uses FUNIT (Liu et al., 2019a) to combine the shape of one image with the class-defining style of another. This serves as a proof of concept, mainly targeting animal-face datasets with limited broad applicability. For **diffusion models**, attention-based conditioning adapters now inject external signals (e.g., an input image) into cross-attention, enabling faithful, controllable generation (Mou et al., 2023). SDEdit (Meng et al., 2022) adds controlled noise and then denoises with conditioning. However, lower noise level yields minimal change, and the higher sacrifices faithfulness (e.g., change object type), thus not suitable as data augmentation that requires both diversity and faithfulness. There are a few recent works using the diffusion-based synthetic images for training augmentation outside of few-shot learning, but they rely on fine-tuning based on text prompts or additional test-class samples, and are thus limited for test-time augmentation (He et al., 2023; Benigimim et al., 2023; Trabucco et al., 2024). More related works are available in Appendix Section §D.

## 3 DIFFUSION PRELIMINARIES

Let  $x_0 \in \mathcal{X}$  be an image and  $x_t$  its noised version at step  $t \in \{1, \dots, T\}$ . A variance schedule  $(\beta_t)_{t=1}^T$  defines  $\alpha_t = 1 - \beta_t$  and  $\bar{\alpha}_t = \prod_{\tau=1}^t \alpha_\tau$ . We use the variance-preserving (VP) forward process (Sohl-Dickstein et al., 2015; Ho et al., 2020; Nichol & Dhariwal, 2021; Meng et al., 2022)

$$q(x_t | x_0) = \mathcal{N}(\sqrt{\bar{\alpha}_t} x_0, (1 - \bar{\alpha}_t)\mathbf{I}), \quad (1)$$

so that sampling  $x_t$  can be written as  $x_t = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon$  with  $\epsilon \sim \mathcal{N}(0, \mathbf{I})$ .

A user noise level  $\eta \in [0, 1]$  is mapped to a discrete start time  $t_0 \in \{1, \dots, T\}$  by matching cumulative noise:

$$t_0 = \arg \min_{t \in \{1, \dots, T\}} |(1 - \bar{\alpha}_t) - \eta^2|. \quad (2)$$

Let  $\hat{\epsilon}_\varphi$  be a learned noise predictor with parameters  $\varphi$ , and let  $c_t$  be the conditioning signal at time  $t$ . The VP reverse update is

$$x_{t-1} = \mu_\varphi(x_t, t, c_t) + \sigma_t \epsilon, \quad (3)$$

with  $\epsilon \sim \mathcal{N}(0, \mathbf{I})$  and sampler noise level  $\sigma_t$  (we use  $\sigma_t=0$  in our experiments), where  $\mu_\varphi(x_t, t, c_t) = \alpha_t^{-1/2}(x_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \hat{\epsilon}_\varphi(x_t, t, c_t))$ . The same formulas apply in a latent space  $z_t$  via an encoder-decoder (Enc, Dec), with denoising output  $\tilde{x} = \text{Dec}(z_0)$  (Rombach et al., 2022b).

## 4 METHOD

**FSL Problem Setup.** We operate in standard inductive FSL classification, where a trained encoder  $\Phi_\theta : \mathcal{X} \rightarrow \mathbb{R}^d$  and a non-parametric classifier are evaluated on novel classes. Given a support set of labelled examples and an unlabelled query, the model must decide which support class the query belongs to. Typical FSL methods compute similarities between the query representation and the representations of support examples under a distance or similarity metric, and assign the query to the closest class prototype Wang et al. (2020); Sung et al. (2018); Ye et al. (2020); Snell et al. (2017). In what follows we denote the encoder by  $\Phi_\theta$  but, when  $\theta$  is fixed, we often simply write  $\Phi$  for brevity. Our goal is to wrap any such trained FSL model with a training-free test-time augmentation operator that, given a single image, synthesises faithful yet diverse variants and aggregates their representations for prediction.

**IS-DAug.** We produce variants of a single image by (1) applying traditional geometric changes, (2) injecting a controlled amount of noise to determine edit magnitude, and (3) denoising via a

diffusion process conditioned on the source image so that content-defining attributes are preserved while details and pose can vary. We write the resulting single-image augmentation operator as

$$\mathcal{A}(x; v) = \text{Den}_\varphi\left(\text{Noi}_\eta(T_\psi(x)), \lambda_{\text{img}}\right) \in \mathcal{X}, \quad (4)$$

where  $v = (\psi, \eta, \lambda_{\text{img}})$  collects the geometric, noising, and conditioning hyperparameters;  $T_\psi$  is a sampled geometric transform;  $\text{Noi}_\eta$  (add noise) applies the forward diffusion process up to a start time  $t_0$  determined by the noise level  $\eta$  (e.g. via equation 2); and  $\text{Den}_\varphi$  (denoise), parametrized by fixed diffusion parameters  $\varphi$ , runs the reverse process from  $t_0$  to 0 with image-conditioned attention. We next detail each stage of this operator.

**Stage 1: Shape Tweak.** Changes in pose/layout increase coverage of plausible views without altering class identity. Let  $\psi$  be shape-tweak parameters and

$$x_{\text{geom}} = T_\psi(x),$$

where  $T_\psi$  is drawn from a family of traditional image transformations, composed as rotations, anisotropic stretches, translations, perspective jitters, and horizontal flips.

**Stage 2: Controlled Noising.** The noise level sets the change magnitude during the diffusion denoising pipeline. Lower noise emphasises faithfulness; higher noise hides geometric distortion better and yields more diversity. We use the variance-preserving (VP) forward kernel (Sohl-Dickstein et al., 2015; Ho et al., 2020; Nichol & Dhariwal, 2021; Meng et al., 2022). Let  $t \in \{1, \dots, T\}$  index discrete diffusion steps, let  $\beta_t \in (0, 1)$  be a variance schedule, and define  $\alpha_t := 1 - \beta_t$  and  $\bar{\alpha}_t := \prod_{\tau=1}^t \alpha_\tau$ . Given a user noise level  $\eta \in [0, 1]$ , we choose a start time  $t_0 = t(\eta)$  (e.g., by matching cumulative noise as in equation 2), and sample

$$x_{t_0} \sim q(x_{t_0} | x_{\text{geom}}) = \mathcal{N}\left(\sqrt{\bar{\alpha}_{t_0}} x_{\text{geom}}, (1 - \bar{\alpha}_{t_0}) \mathbf{I}\right), \quad (5)$$

where  $x_{\text{geom}}$  is the geometrically perturbed input (Stage 1),  $\mathbf{I}$  is the identity covariance, and  $\mathcal{N}$  denotes a Gaussian distribution.

**Stage 3: Image-Conditioned Diffusion Denoising.** Let  $z_t$  be the latent at time  $t$ , and let  $f_{\text{img}}(x) \in \mathbb{R}^{L \times d_k}$  and  $f_{\text{txt}}(p) \in \mathbb{R}^{M \times d_k}$  be fixed-encoder outputs for the condition image  $x$  and optional text prompt  $p$ . For a U-Net (Ronneberger et al., 2015) block at time  $t$ , with queries  $Q_t = W_Q z_t \in \mathbb{R}^{N_q \times d_k}$  (here  $W_Q$  denotes the query weights and biases and  $N_q$  is the number of query tokens) and keys/values  $K_t, V_t \in \mathbb{R}^{(M+L) \times d_k}$ , the cross-attention (Vaswani et al., 2017) is

$$A_t(Q_t, K_t, V_t) = \text{softmax}\left(\frac{Q_t K_t^\top}{\sqrt{d_k}}\right) V_t, \quad (6)$$

and we concatenate text/image tokens with a scalar weight  $\lambda_{\text{img}} \geq 0$ :

$$K_t = [K_t^{\text{txt}}, \lambda_{\text{img}} K_t^{\text{img}}], \quad V_t = [V_t^{\text{txt}}, \lambda_{\text{img}} V_t^{\text{img}}], \quad (7)$$

with  $K_t^{\text{txt}}, V_t^{\text{txt}} = W_{\text{txt}, t} f_{\text{txt}}(p)$  and  $K_t^{\text{img}}, V_t^{\text{img}} = W_{\text{img}, t} f_{\text{img}}(x)$ . We set the conditioning variable for the reverse update to  $c_t := A_t(Q_t, K_t, V_t)$ . The VP reverse step equation 3 then reads

$$z_{t-1} = \mu_\varphi(z_t, t, c_t) + \sigma_t \epsilon, \quad \epsilon \sim \mathcal{N}(0, \mathbf{I}), \quad (8)$$

and rolling out  $t_0 \rightarrow 0$  produces  $\tilde{x} = \text{Dec}(z_0)$ .

**FSL Feature Aggregation.** Averaging features over faithful but non-identical views draws representations toward class-typical regions and improves robustness for non-parametric few-shot classifiers (Snell et al., 2017; Chen et al., 2019; Ye et al., 2020; Bai et al., 2025). For each image  $x$  we generate  $K_a^{\text{sup}} + 1$  support views  $\tilde{x}^{(k)} = \mathcal{A}(x; v^{(k)})$  with  $\tilde{x}^{(0)} \equiv x$ , where each  $v^{(k)}$  denotes a fresh sample of augmentation hyperparameters, and form an aggregated support embedding

$$\bar{z}_{\text{sup}}(x) = \sum_{k=0}^{K_a^{\text{sup}}} \alpha_k^{\text{sup}} \Phi(\tilde{x}^{(k)}), \quad \alpha_k^{\text{sup}} \geq 0, \quad \sum_{k=0}^{K_a^{\text{sup}}} \alpha_k^{\text{sup}} = 1. \quad (9)$$

Given a support set  $S = \{(s_i, y_i)\}_{i=1}^{NK}$  with  $N$  classes and  $K$  shots per class, class prototypes are  $p_c = \frac{1}{K} \sum_{i: y_i=c} \bar{z}_{\text{sup}}(s_i)$ . For each query  $q$  we generate  $K_a^{\text{qry}} + 1$  query views  $\tilde{q}^{(k)} = \mathcal{A}(q; v^{(k)})$ , encode them as  $z^{(k)}(q) = \Phi(\tilde{q}^{(k)})$ , and compute per-view logits  $\ell_c^{(k)}(q) = \kappa(z^{(k)}(q), p_c)$  for a chosen similarity  $\kappa$  (Euclidean or cosine). We then perform query-side logit averaging for FSL prediction:

$$\tilde{\ell}_c(q) := \sum_{k=0}^{K_a^{\text{qry}}} \alpha_k^{\text{qry}} \ell_c^{(k)}(q), \quad \alpha_k^{\text{qry}} \geq 0, \quad \sum_{k=0}^{K_a^{\text{qry}}} \alpha_k^{\text{qry}} = 1, \quad \hat{y}(q) = \arg \max_{c \in \{1, \dots, N\}} \tilde{\ell}_c(q). \quad (10)$$

## 5 THEORETICAL INSIGHTS

We analyse 1S-DAug in the standard episodic few-shot setting with a single trainable encoder and a fixed Euclidean nearest-prototype classifier. We show that (i) a simple risk decomposition separates the ensemble benefit into accuracy and diversity, and (ii) a margin-based generalisation bound for the encoder becomes strictly tighter after augmentation, via both empirical margin and feature-radius reduction. Lastly, based on the generalisation bound, we compare training and test-time augmentation in Appendix Section §H.7, highlighting the latter’s comparative advantage. All missing definitions and proofs are deferred to Appendix Sections §G, §H.2, §H.5, and §H.6.

### 5.1 EPISODIC EUCLIDEAN MODEL AND TEST-TIME AUGMENTATION

We consider the encoder-plus-Euclidean-prototype classifier of Section 4. In an  $N$ -way  $K$ -shot episode, class prototypes  $p_c$  are formed by averaging support embeddings, and a query  $q$  is assigned to the nearest prototype in squared Euclidean distance  $\|\Phi_\theta(q) - p_c\|_2^2$ . For the theory we reduce episodes to binary query–prototype pairs  $x = (q, p)$  with label  $y \in \{-1, 1\}$ , write

$$\Omega_\theta(x) := \Phi_\theta(q) - p \quad \text{and} \quad g_\theta(x) := -\|\Omega_\theta(x)\|_2^2,$$

and assume a uniform radius bound  $\|\Omega_\theta(x)\|_2 \leq r_0$ . Under test-time augmentation, multiple query views are combined by logit averaging; for squared Euclidean scores this is equivalent to using an averaged query embedding and hence an aggregated difference feature  $\Omega_\theta(x)$ . Full episodic details and the logit–feature equivalence are deferred to Appendix Section §H.1 and §H.3.

### 5.2 RISK DECOMPOSITION INTO ACCURACY AND DIVERSITY

We first quantify the ensemble effect at the pairwise level. For any real-valued predictor  $g$  on pairs  $x = (q, p)$  we use the scaled squared-loss risk

$$\mathcal{R}(g) := \frac{1}{4} \mathbb{E}_{(x,y) \sim D} [(g(x) - y)^2], \quad (11)$$

which coincides with 0–1 pairwise error when  $g(x) \in \{-1, 1\}$ . In particular, if  $g_\theta$  is a sign-valued score, the pairwise misclassification risk  $R_{\text{cls}}(\theta) := \mathbb{P}_{(x,y) \sim D}(y g_\theta(x) \leq 0)$  satisfies  $R_{\text{cls}}(\theta) = \mathcal{R}(g_\theta)$ . Let  $f(x)$  and  $f_A(x)$  be the sign predictors associated with the base and an augmented view, and define the (two-view) ensemble  $\tilde{f}(x) := \frac{1}{2}(f(x) + f_A(x)) \in \{-1, 0, 1\}$ . A direct calculation (Appendix Section §G) yields:

**Proposition 1** (Pairwise risk decomposition). *With  $\mathcal{R}(\cdot)$  as in equation 11,*

$$\mathcal{R}(\tilde{f}) - \mathcal{R}(f) = \underbrace{\frac{1}{4} (\mathbb{E}[f(x)y] - \mathbb{E}[f_A(x)y])}_{\text{accuracy gap}} + \underbrace{\frac{1}{8} (\mathbb{E}[f(x)f_A(x)] - 1)}_{\text{diversity term}}. \quad (12)$$

Thus improvements come from (i) maintaining or improving single-view accuracy, and (ii) making the augmented predictions sufficiently diverse on hard examples. This matches the empirical behaviour of 1S-DAug, which is designed to generate plausible but different-shape query views.

### 270 5.3 ENCODER GENERALISATION, RADIUS REDUCTION, AND 1S-DAUG

271  
272 With the encoder as the only learnable component, our method tightens a margin-based general-  
273 isation bound. Let  $\mathcal{G} := \{g_\theta : \theta \in \Theta\}$  be the score class induced by the encoder, and let  
274  $R_{\text{cls}}(\theta) := \mathbb{P}_{(x,y) \sim D}(y g_\theta(x) \leq 0)$  denote the pairwise misclassification risk. For a sample  
275  $S = \{(x_i, y_i)\}_{i=1}^m$  and margin parameter  $\rho > 0$ , let  $\widehat{R}_{S,\rho}(\theta)$  be the empirical  $\rho$ -margin loss for  $g_\theta$ ,  
276 and let  $\widehat{\mathfrak{R}}_S(\mathcal{G})$  denote the empirical Rademacher complexity of  $\mathcal{G}$ . Formal definitions are deferred  
277 to Appendix Section §H.2. A standard margin-based Rademacher argument gives the following.

278 **Theorem 1** (Encoder margin bound). *For any  $\rho > 0$  and  $\delta > 0$ , with probability at least  $1 - \delta$  over*  
279  *$S \sim D^m$ , every encoder  $\theta$  satisfies*

$$281 R_{\text{cls}}(\theta) \leq \widehat{R}_{S,\rho}(\theta) + \frac{2}{\rho} \widehat{\mathfrak{R}}_S(\mathcal{G}) + \sqrt{\frac{\log(1/\delta)}{2m}}. \quad (13)$$

284 Suppose  $\Phi_\theta$  is a feedforward network with 1-Lipschitz nonlinearities and layer spectral norms  
285  $\|W_\ell\|_2 \leq s_\ell$  such that  $\prod_{\ell=1}^L s_\ell \leq L_{\text{enc}}$ . Then  $g_\theta$  is Lipschitz in the difference feature with constant  
286 proportional to  $L_{\text{enc}} r_0$ , where  $r_0 := \sup_x \|\Omega_\theta(x)\|_2$  is the (pre-augmentation) feature radius from  
287 Section 5.1. Using standard spectral-norm bounds, Lemma 2 (Appendix Section §H.2) shows that  
288  $\widehat{\mathfrak{R}}_S(\mathcal{G}) \lesssim L_{\text{enc}} \frac{r_0}{\sqrt{m}}$ , so the complexity term in equation 13 scales linearly with the feature radius  
289  $r_0$ . Let  $\hat{r} := \sup_x \|\widehat{\Omega}_\theta(x)\|_2$  be the radius after augmentation and  $\widehat{R}_{\widehat{S},\rho}(\theta)$  the empirical margin loss  
290 computed with the aggregated sample  $\widehat{S}$ . By convexity, augmentation cannot increase the radius:  
291

292 **Lemma 1** (Radius contraction under augmentation). *If  $\|\Omega_\theta^{(k)}(x)\|_2 \leq r_0$  for all  $x$  and views  $k$ , then*  
293  *$\hat{r} := \sup_x \|\widehat{\Omega}_\theta(x)\|_2$  satisfies  $\hat{r} \leq r_0$ . Moreover, under a simple symmetric i.i.d. model for the views*  
294  *$\{\Omega_\theta^{(k)}(x)\}$ , the maximal empirical radius strictly decreases with non-zero probability.*

296 Combining Theorem 1, Lemma 2, and Lemma 1, the encoder-only risk before test-time augmenta-  
297 tion satisfies

$$299 R_{\text{cls}}(\theta) \leq \widehat{R}_{S,\rho}(\theta) + \frac{2C_{\text{enc}}L_{\text{enc}}}{\rho} \frac{r_0}{\sqrt{m}} + \sqrt{\frac{\log(1/\delta)}{2m}}, \quad (14)$$

301 for some architecture-dependent constant  $C_{\text{enc}} > 0$ . After test-time augmentation, the bound is  
302 the same except that  $\widehat{R}_{S,\rho}(\theta)$  is replaced with  $\widehat{R}_{\widehat{S},\rho}(\theta)$ , and  $r_0$  with  $\hat{r}$ . Since  $\hat{r} = \alpha r_0$  for some  
303  $\alpha \in (0, 1]$ , the complexity term contracts by factor  $\alpha$ , and Appendix Section §H.5 further shows that  
304 augmentation tends to move representations towards their class prototypes, increasing margins and  
305 thereby reducing the empirical term as well (i.e.,  $\widehat{R}_{\widehat{S},\rho}(\theta) < \widehat{R}_{S,\rho}(\theta)$ ). In summary, 1S-DAug tight-  
306 ens the encoder’s generalisation bound through both empirical margin and feature-radius reduction.  
307

## 308 6 EXPERIMENTS

### 309 6.1 SET-UP AND MAIN RESULTS

312 **Datasets and Pre-processing.** We follow the standard 5-way-1/5-shot episodic evaluation on  
313 miniImagenet and tieredImagenet (Russakovsky et al., 2015); additional experiments are conducted  
314 on CUB (fine-grained birds) (Wah et al., 2011) and an animal-face dataset used previously for test-  
315 time augmentation studies (Animals) (Liu et al., 2019a; Bai et al., 2025). We adhere to the conven-  
316 tional train/val/test splits for each benchmark. Besides, we set the aggregation weight of the original  
317 image as 0.5, and additional images as 0.5 altogether, so as to emphasize the original images. For  
318 diffusion noise addition, we set the noise level to 0.7. Details are available in Appendix Section §I.  
319

320 **Evaluation Protocol.** We evaluate three standard backbones used in FSL: a shallow 4-layer con-  
321 volutional neural network (ConvNet), a 12-layer residual network (Res12) (He et al., 2016), a  
322 small vision transformer backbone (ViTSmall) (Dosovitskiy et al., 2021) and a tiny swin trans-  
323 former backbone (SwinTiny) (Liu et al., 2021). Note that the encoders ViTSmall and SwinTiny are  
pretrained on the ImageNet-1k Russakovsky et al. (2015) dataset, while the Res12 and ConvNet

encoders on trained on the train-split of each FSL dataset. Classification is performed with a non-parametric Euclidean-distance/cosine-similarity prototype classifier in the episodic setting. These choices match common FSL practice, including ProtoNet and FEAT-style set-to-set variants that operate on support/query embeddings, ensuring comparability with prior work and our reproduced baselines (Snell et al., 2017; Ye et al., 2020). We evaluate off-the-shelf trained ProtoNet and FEAT models. We adopt Euclidean-distance-based classifier for miniImagenet and tieredImagenet, and cosine-similarity-based classifier for Animals and CUB. Our reproduced baselines closely match prior reports evaluation set-up (e.g., DeepEMD (Zhang et al., 2020), Meta-Baseline (Chen et al., 2021), MetaOptNet (Lee et al., 2019) on 5-way-1/5-shot inductive benchmarks, all with the Res12 backbone. Among the baselines, SLA-AG Lee et al. (2020) involves self-supervised label augmentation, and Meta-MaxUp Ni et al. (2020) involves training data augmentation, both being ensemble-based methods.). We sample 15,000 5-way-1/5-shot queries and report mean accuracy with 95% confidence intervals across episodes. This follows the standard protocol used in related FSL work (Ye et al., 2020; Snell et al., 2017).

Method (Res12)	5-Way-1-Shot (%)		5-Way-5-Shot (%)	
	miniImagenet	tieredImagenet	miniImagenet	tieredImagenet
DeepEMD Zhang et al. (2020)	65.91 ± 0.82	71.16 ± 0.87	82.41 ± 0.56	<b>86.03</b> ± 0.58
Meta-MaxUp Ni et al. (2020)	62.81 ± 0.34	-	79.38 ± 0.24	-
Meta-Baseline Chen et al. (2021)	63.17 ± 0.23	68.62 ± 0.27	79.26 ± 0.17	83.74 ± 0.18
MetaOptNet Lee et al. (2019)	62.64 ± 0.61	65.99 ± 0.72	78.63 ± 0.46	81.56 ± 0.53
SLA-AG Lee et al. (2020)	62.93 ± 0.63	-	79.63 ± 0.47	-
ProtoNet + TRAML Li et al. (2020)	60.31 ± 0.48	-	77.94 ± 0.57	-
ConstellationNet Xu et al. (2021)	64.89 ± 0.23	-	79.95 ± 0.17	-
Classifier-Baseline Chen et al. (2021)	58.91 ± 0.23	68.07 ± 0.26	77.76 ± 0.17	83.74 ± 0.18
DFR Cheng et al. (2023)	<b>67.74</b> ± 0.86	<b>71.31</b> ± 0.93	<b>82.49</b> ± 0.57	85.12 ± 0.64
ProtoNet-Res12 Snell et al. (2017)	62.39 ± 0.21	68.23 ± 0.23	80.53 ± 0.14	84.03 ± 0.16
ProtoNet-Res12 (re-impl.)	60.01 ± 0.65	65.28 ± 0.32	75.34 ± 0.49	81.13 ± 0.29
ProtoNet-Res12+1S-DAug-1 (Ours)	62.90 ± 0.66 (+2.89% ↑)	69.06 ± 0.32 (+3.78% ↑)	78.89 ± 0.48 (+3.55% ↑)	83.86 ± 0.27 (+2.73% ↑)
ProtoNet-Res12+1S-DAug-2 (Ours)	64.61 ± 0.66 (+4.60% ↑)	70.32 ± 0.32 (+5.04% ↑)	-	-
ProtoNet-Res12+1S-DAug-3 (Ours)	64.94 ± 0.66 (+4.93% ↑)	-	-	-
FEAT-Res12 Ye et al. (2020)	66.78 ± 0.20	70.80 ± 0.23	82.05 ± 0.14	84.79 ± 0.16
FEAT-Res12 (re-impl.)	63.31 ± 0.65	68.28 ± 0.28	77.90 ± 0.48	82.21 ± 0.28
FEAT-Res12+1S-DAug-1 (Ours)	67.08 ± 0.65 (+3.77% ↑)	71.85 ± 0.28 (+3.57% ↑)	81.96 ± 0.44 (+4.06% ↑)	84.82 ± 0.26 (+2.61% ↑)
FEAT-Res12+1S-DAug-2 (Ours)	69.04 ± 0.65 (+5.73% ↑)	<b>73.18</b> ± 0.27 (+4.90% ↑)	82.62 ± 0.45 (+4.72% ↑)	<b>85.55</b> ± 0.25 (+3.34% ↑)
FEAT-Res12+1S-DAug-3 (Ours)	<b>69.25</b> ± 0.65 (+5.94% ↑)	-	<b>83.38</b> ± 0.41 (+5.48% ↑)	-

Table 1: Inductive 5-way-1-shot and 5-way-5-shot accuracy (%) on miniImagenet and tieredImagenet with Res12 backbones. Dashes denote unavailable or less important results not reported. The best results of ours and other FSL methods are both highlighted in bold. Our method transforms the weaker models to become stronger than most of the other Res12 baselines; we can likely achieve even better performance using stronger base models.

Method (Res12/ConvNet)	Animals	CUB
ProtoNet	73.20 ± 0.63	46.38 ± 0.22
ProtoNet+1S-DAug-2 (Ours)	75.20 ± 0.65 (+2.00% ↑)	55.50 ± 0.24 (+9.12% ↑)
FEAT	79.37 ± 0.59	51.10 ± 0.24
FEAT+1S-DAug-2 (Ours)	<b>80.66</b> ± 0.62 (+1.23% ↑)	<b>61.55</b> ± 0.25 (+10.45% ↑)

Table 2: Inductive 5-way-1-shot accuracy (%) on Animals with Res12 backbones and CUB with ConvNet backbones.

Dataset	Method	ViTSmall	SwinTiny
MiniImagenet	ProtoNet Snell et al. (2017)	71.86	67.32
	ProtoNet+1S-DAug-1 (Ours)	80.42 (+8.56% ↑)	75.12 (+7.80% ↑)
	ProtoNet+1S-DAug-2 (Ours)	82.76 (+10.90% ↑)	77.82 (+10.50% ↑)
	ProtoNet+1S-DAug-3 (Ours)	83.66 (+11.80% ↑)	78.92 (+11.60% ↑)
CUB	ProtoNet Snell et al. (2017)	71.90	69.78
	ProtoNet+1S-DAug-1 (Ours)	75.72 (+3.82% ↑)	71.76 (+1.98% ↑)

Table 3: 5-way-1-shot accuracy (%) on miniImagenet/CUB with ViTSmall/SwinTiny backbones.

**Main Results.** Table 1/2/3 summarizes 5-way-1-shot and 5-way-5-shot results with Res12/ConvNet/ViTsmall/SwinTiny backbones, and our method with 1/2/3 additional augmentations are denoted as 1S-DAug-1/2/3 respectively. Note that we directly adopt 5-way-1-shot FSL models for the 5-way-5-shot evaluation, and Table 2 contains 5-way-1-shot results on CUB and Animals. As reported, test-time 1S-DAug consistently improves over the corresponding non-augmented baselines and over prior strong Res12 methods reported under the same backbone (e.g., on miniImagenet, FEAT improves from 63.31% to 69.25%, a maximum absolute gain of +5.94 percentage points, which achieves the highest among all the reported FSL works with Res12 backbones). The gains persist across both datasets (e.g., miniImagenet improves by at least +2.89%, tieredImagenet by +3.88%, CUB by +9.12%, and Animals by +1.23% on the 5-way-1-shot benchmarks), implying a high probability for: (i) the image-conditioned diffusion step preserves class-defining content sufficiently, and (ii) the shape tweak with noising creates diversity without compromising faithfulness.

## 6.2 ABLATION STUDIES

We conduct ablation studies for our method, and more analyses, including hyperparameter tuning (Section §B), efficiency studies (Section §A) and limitation (Section §C) are available in Appendix.

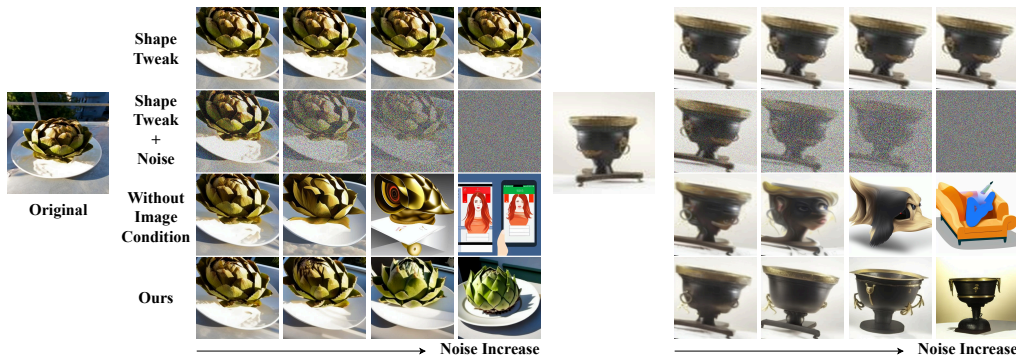


Figure 2: Effect of noise and conditioning. Qualitative ablation on a single input across increasing noise levels. Shape-only edits yield limited diversity; adding noise increases diversity but may reduce fidelity without conditioning. Attention-conditioned diffusion preserves class-defining content while enabling controlled pose/appearance changes; excessive noise without the image condition degrades faithfulness.

**Qualitative Analysis.** Figure 2 illustrates the effect of removing image conditioning, removing shape tweaks, and sweeping the noise level. With very small noise, changes are minimal; with very large noise and no image conditioning, generations may drift toward off-class content; removing shape tweaks reduces diversity and visible pose/layout variations. Additional visualizations are provided in the Appendix Section §K.

**Quantitative Analysis.** We further dissect the contribution of each component, including aggregation weight adjustment, image conditioning, noise magnitude, shape tweaking and diffusion generation using FEAT with a Res12 backbone on miniImagenet, in a controlled 5-way-1-shot setting with one augmented query and one augmented support per episode (Table 4). We first notice that reducing the emphasis on original samples via less aggregation weight would downgrade model accuracy slightly. This is expected, since the original samples are authentic images with the best quality. Besides, removing the image conditioning downgrades the performance severely, which mirrors our qualitative studies in Figure 2. When a small noise level ( $\eta=0.20$ ) is applied with shape tweaking and conditioning, diversity gain is limited, and distortion caused by shape tweaking may also stay and backfire, yielding 63% accuracy. Increasing the noise strength to a moderate level ( $\eta=0.70$ ) improves coverage while preserving class faithfulness, pushing performance to 67.1%, the best among diffusion-based rows. Pushing noise to the extreme ( $\eta=1.0$ ) still delivers reasonable performance (67.0%) when conditioning is enabled. Such a performance is enabled by the

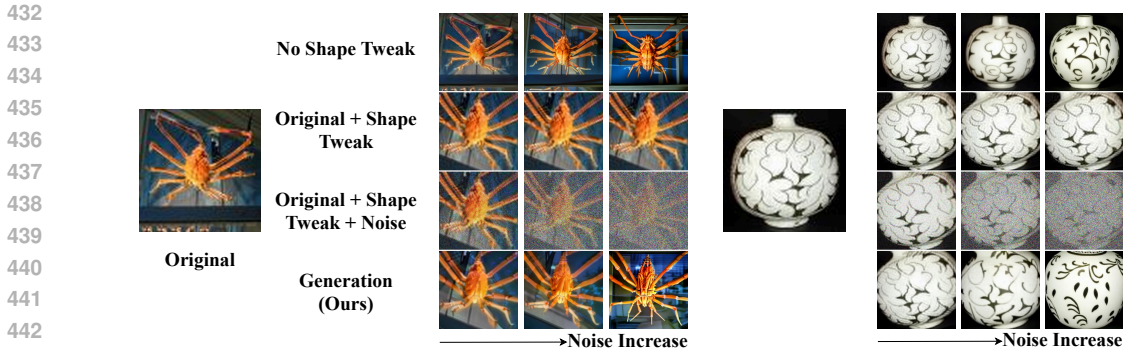


Figure 3: Effect of noise and shape tweak. Comparison across three settings: no shape tweak, shape tweak only, and shape + noise + attention-conditioned diffusion (ours). Increasing noise and including shape tweak expand diversity, and our full setting provides the best balance for both diversity and faithfulness.

greater diversity between the generation output and the original input, but the originality faithfulness is jeopardized, especially when confronted with rare object types. Therefore, generation from full noise should be discouraged. Besides, ablating shape tweaks reduces accuracy to 66.1%, confirming that geometric variation helps cover different pose/layout. Substituting true extra images of the same class (“Real/Oracle”) provides an upper bound of 78.0%, showing the headroom available with more independent samples. Meanwhile, traditional geometric test-time edits (rotations, affine warps, color jitter) only reach 57.89%, supporting the observation that such transforms add little new information and may distort original images.

Same Noise Level	Shape Tweak	Generation Techniques	Additional Adjustment	5-Way-1-Shot Accuracy
✓(0.70)	✓	✓	✓(0.3 original image weight)	66.79 ± 0.63
✓(0.70)	✓	✓	✓(remove image conditioning)	53.94 ± 0.62
×(0.20)	✓	✓	×	63.45 ± 0.64
✓(0.70)	✓	✓	×	67.08 ± 0.62
×(1.00)	-	✓	×	67.01 ± 0.68
✓(0.70)	×	✓	×	66.12 ± 0.66
-	✓(Real/Oracle)	×	×	77.99 ± 0.62
-	✓(Traditional)	×	×	57.89 ± 0.67

Table 4: Ablation of aggregation weight adjustment, shape tweak, noise level, and diffusion conditioning. ‘Traditional’ uses standard geometric edits; ‘Real’ substitutes actual additional images of the same object type. Our full setting (shape + controlled noise + attention-conditioned diffusion) outperforms classical test-time augmentation and approaches the oracle (true image) upper bound.

## 7 CONCLUSIONS AND FUTURE WORK

To conclude, we presented *IS-DAug*, a one-shot, test-time generative augmentation operator that synthesizes diverse yet faithful variants from a single image. By combining geometric perturbations with controlled noising and attention-conditioned denoising, the method maintains class-defining content while enhancing data diversity. As a plugin into standard FSL models, *IS-DAug* delivers consistent non-trivial accuracy gains under the 5-way-1/5-shot protocol across datasets, without FSL training, fine-tuning, or access to specific model parameters. This model-agnostic, data-side design makes the approach practical for modern deployments where models are large, fixed, or restricted. While our focus is test-time augmentation for FSL, the operator is also applicable to other downstream tasks like training-time augmentation and controllable image editing. Future work will explore inference speedups, more backbones and constraining diffusion-based image generator data to within the FSL training splits. Overall, we believe *IS-DAug* offers a useful building block for data-centric robustness with high potential in real-world low-label scenarios.

## REFERENCES

- 486  
487  
488 Yunwei Bai, Ying Kiat Tan, Shiming Chen, Yao Shu, and Tsuhan Chen. Fsl-rectifier: Rectify  
489 outliers in few-shot learning via test-time augmentation. In *Proceedings of the AAAI Conference*  
490 *on Artificial Intelligence*, volume 39, pp. 15462–15471, 2025.
- 491 Peter L. Bartlett and Shahar Mendelson. Rademacher and gaussian complexities: Risk bounds and  
492 structural results. *Journal of Machine Learning Research*, 3:463–482, 2002.
- 493 Yasser Benigmim, Subhankar Roy, Slim Essid, Vicky Kalogeiton, and Stéphane Lathuilière.  
494 One-shot unsupervised domain adaptation with personalized diffusion models. In *Proceed-*  
495 *ings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*  
496 *(CVPRW)*, pp. 698–708, June 2023. URL [https://openaccess.thecvf.com/  
497 content/CVPR2023W/GCV/papers/Benigmim\\_One-Shot\\_Unsupervised\\_  
498 Domain\\_Adaptation\\_With\\_Personalized\\_Diffusion\\_Models\\_CVPRW\\_2023\\_  
499 paper.pdf](https://openaccess.thecvf.com/content/CVPR2023W/GCV/papers/Benigmim_One-Shot_Unsupervised_Domain_Adaptation_With_Personalized_Diffusion_Models_CVPRW_2023_paper.pdf).
- 500 Wei-Yu Chen, Yen-Cheng Liu, Zsolt Kira, Yu-Chiang Frank Wang, and Jia-Bin Huang. A closer  
501 look at few-shot classification. *arXiv preprint arXiv:1904.04232*, 2019.
- 502 Yinbo Chen, Zhuang Liu, Huijuan Xu, Trevor Darrell, and Xiaolong Wang. Meta-baseline: Explor-  
503 ing simple meta-learning for few-shot learning. In *Proceedings of the IEEE/CVF International*  
504 *Conference on Computer Vision (ICCV)*, 2021.
- 506 Hao Cheng, Yufei Wang, Haoliang Li, Alex C Kot, and Bihan Wen. Disentangled feature repre-  
507 sentation for few-shot image classification. *IEEE transactions on neural networks and learning*  
508 *systems*, 35(8):10422–10435, 2023.
- 509 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas  
510 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszko-  
511 reit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at  
512 scale. In *International Conference on Learning Representations (ICLR)*, 2021.
- 514 Nikita Dvornik, Cordelia Schmid, and Julien Mairal. Diversity with cooperation: Ensemble methods  
515 for few-shot classification. In *Proceedings of the IEEE/CVF international conference on computer*  
516 *vision*, pp. 3723–3731, 2019.
- 517 Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation  
518 of deep networks. In *ICML*, 2017.
- 519 Bharath Hariharan and Ross Girshick. Low-shot visual recognition by shrinking and hallucinating  
520 features. In *Proceedings of the IEEE international conference on computer vision*, pp. 3018–3027,  
521 2017a.
- 522 Bharath Hariharan and Ross Girshick. Low-shot visual recognition by shrinking and hallucinating  
523 features. In *CVPR*, 2017b.
- 525 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-  
526 nition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp.  
527 770–778, 2016.
- 528 Ruifei He, Shuyang Sun, Xin Yu, Chuhui Xue, Wenqing Zhang, Philip Torr, Song Bai, and Xiaojuan  
529 Qi. Is synthetic data from generative models ready for image recognition? In *Proceedings*  
530 *of the International Conference on Learning Representations (ICLR)*, 2023. URL [https://  
531 openreview.net/forum?id=nUmCcZ5RKF](https://openreview.net/forum?id=nUmCcZ5RKF). Spotlight.
- 532 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *NeurIPS*,  
533 2020.
- 535 Ahmet Iscen, Giorgos Tolias, Yannis Avrithis, and Ondrej Chum. Transductive information maxi-  
536 mization for few-shot learning. In *CVPR*, 2020.
- 537 Bingyi Kang, Yu Xie, Marcus Rohrbach, Zhicheng Yan, Albert Gordo, Jiashi Feng, and Yannis  
538 Kalantidis. Decoupling representation and classifier for long-tailed recognition. In *International*  
539 *Conference on Learning Representations (ICLR)*, 2020.

- 540 Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative  
541 adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern*  
542 *recognition*, pp. 4401–4410, 2019.
- 543 Hyunwoo Lee, Hyo-Eun Kim, and Jin-Hwa Kim. Self-supervised label augmentation via input trans-  
544 formations. In *Proceedings of the 37th International Conference on Machine Learning (ICML)*,  
545 2020. SLA-AG variant (aggregation) results reported in the paper.
- 546 Kwonjoon Lee, Subhansu Maji, Avinash Ravichandran, and Stefano Soatto. Meta-learning with  
547 differentiable convex optimization. In *Proceedings of the IEEE/CVF Conference on Computer*  
548 *Vision and Pattern Recognition (CVPR)*, 2019.
- 549 Aoxue Li, Weiran Huang, Xu Lan, Jiashi Feng, Zhenguo Li, and Liwei Wang. Boosting few-shot  
550 learning with adaptive margin loss. In *Proceedings of the IEEE/CVF Conference on Computer*  
551 *Vision and Pattern Recognition (CVPR)*, pp. 12576–12584, 2020.
- 552 Zhenguo Li, Feng Zhou, Fei Chen, and Hang Li. Meta-sgd: Learning to learn quickly for few shot  
553 learning. In *NeurIPS*, 2017.
- 554 Jian Liang, Ran He, and Tieniu Tan. A comprehensive survey on test-time adaptation under dis-  
555 tribution shifts. *International Journal of Computer Vision*, 133(1):31–64, July 2024. ISSN  
556 1573-1405. doi: 10.1007/s11263-024-02181-w. URL [http://dx.doi.org/10.1007/](http://dx.doi.org/10.1007/s11263-024-02181-w)  
557 [s11263-024-02181-w](http://dx.doi.org/10.1007/s11263-024-02181-w).
- 558 Henry X Liu and Shuo Feng. Curse of rarity for autonomous vehicles. *nature communications*, 15  
559 (1):4808, 2024.
- 560 Ming-Yu Liu, Xun Huang, Arun Mallya, Tero Karras, Timo Aila, Jaakko Lehtinen, and Jan Kautz.  
561 Few-shot unsupervised image-to-image translation. In *arxiv*, 2019a.
- 562 Yanbin Liu, Juho Lee, Minseop Park, Saehoon Kim, Eunho Yang, Sung Ju Hwang, and Yi Yang.  
563 Learning to propagate labels: Transductive propagation network for few-shot learning. In *ICLR*,  
564 2019b.
- 565 Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo.  
566 Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of*  
567 *the IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 10012–10022, October  
568 2021.
- 569 Chenlin Meng, Yutong He, Jiaming Song, et al. SDEdit: Image synthesis and editing with stochastic  
570 differential equations. In *ICLR*, 2022.
- 571 Chong Mou, Yujun Wang, Yutong Bai, et al. T2i-adapter: Learning adapters to dig out more con-  
572 trollable ability for text-to-image diffusion models. In *CVPR*, 2023.
- 573 Renkun Ni, Micah Goldblum, Amr Sharaf, Kezhi Kong, and Tom Goldstein. Data augmentation for  
574 meta-learning. *arXiv preprint arXiv:2010.07092*, 2020.
- 575 Alex Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models.  
576 *arXiv:2102.09672*, 2021.
- 577 Alex Nichol, Joshua Achiam, and John Schulman. On first-order meta-learning algorithms.  
578 *arXiv:1803.02999*, 2018.
- 579 Maithra Raghu, Maruan Raghu, Samy Bengio, and Oriol Vinyals. Rapid learning or feature reuse?  
580 towards understanding the effectiveness of maml. In *ICLR*, 2020.
- 581 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. Stable  
582 diffusion v1.5. <https://huggingface.co/runwayml/stable-diffusion-v1-5>,  
583 2022a. Accessed: 2025-12-04.
- 584 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
585 resolution image synthesis with latent diffusion models, 2022b.

- 594 Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical  
595 image segmentation. In *International Conference on Medical Image Computing and*  
596 *Computer-Assisted Intervention (MICCAI)*, pp. 234–241. Springer, 2015.  
597
- 598 Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng  
599 Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei.  
600 ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision*  
601 (*IJCV*), 115(3):211–252, 2015. doi: 10.1007/s11263-015-0816-y.
- 602 Eli Schwartz, Leonid Karlinsky, Shai Avidan, and Alex M Bronstein. Delta-encoder: an effective  
603 sample synthesis method for few-shot object recognition. In *NeurIPS*, 2018.  
604
- 605 Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. *Ad-*  
606 *vances in neural information processing systems*, 30, 2017.
- 607 Jascha Sohl-Dickstein, Eric A Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsuper-  
608 vised learning using nonequilibrium thermodynamics. In *ICML*, 2015.  
609
- 610 Flood Sung, Yongxin Yang, Li Zhang, Tao Xiang, Philip HS Torr, and Timothy Hospedales. Learn-  
611 ing to compare: Relation network for few-shot learning. In *CVPR*, 2018.  
612
- 613 Brandon Trabucco, Kyle Doherty, Max A. Gurinas, and Ruslan Salakhutdinov. Effective data aug-  
614 mentation with diffusion models. In *Proceedings of the Twelfth International Conference on*  
615 *Learning Representations (ICLR)*, 2024. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=ZWzUA9zeAg)  
616 [ZWzUA9zeAg](https://openreview.net/forum?id=ZWzUA9zeAg). Poster.
- 617 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,  
618 Lukasz Kaiser, and Illia Polosukhin. Attention is all you need, 2017.  
619
- 620 Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Koray Kavukcuoglu, and Daan Wierstra. Match-  
621 ing networks for one shot learning. In *NeurIPS*, 2016.
- 622 C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie. cub. Technical Report CNS-TR-2011-  
623 001, California Institute of Technology, 2011.  
624
- 625 Xintao Wang, Liangbin Xie, Chao Dong, and Ying Shan. Real-esrgan: Training real-world blind  
626 super-resolution with pure synthetic data. In *International Conference on Computer Vision Work-*  
627 *shops (ICCVW)*, 2021.
- 628 Yaqing Wang, Quanming Yao, James T Kwok, and Lionel M Ni. Generalizing from a few examples:  
629 A survey on few-shot learning. *ACM computing surveys (csur)*, 53(3):1–34, 2020.  
630
- 631 Yu-Xiong Wang, Ross Girshick, Martial Hebert, and Bharath Hariharan. Low-shot learning from  
632 imaginary data. In *Proceedings of the IEEE conference on computer vision and pattern recogni-*  
633 *tion*, pp. 7278–7286, 2018.  
634
- 635 Weijian Xu, Yifan Xu, Huaijin Wang, and Zhuowen Tu. Attentional constellation nets for few-shot  
636 learning. In *International Conference on Learning Representations (ICLR)*, 2021.
- 637 Han-Jia Ye, Hexiang Hu, De-Chuan Zhan, and Fei Sha. Few-shot learning via embedding adaptation  
638 with set-to-set functions. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*  
639 (*CVPR*), pp. 8808–8817, 2020.  
640
- 641 Peng Ye, Huan Zhang, Xiaodan Liu, et al. Ip-adapter: Text compatible image prompt adapter for  
642 text-to-image diffusion models. *arXiv:2308.06721*, 2023.
- 643 Chi Zhang, Yujun Cai, Guosheng Lin, and Chunhua Shen. Deepemd: Few-shot image classifica-  
644 tion with differentiable earth mover’s distance and structured classifiers. In *Proceedings of the*  
645 *IEEE/CVF conference on computer vision and pattern recognition*, pp. 12203–12213, 2020.  
646
- 647 Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image  
diffusion models. *arXiv:2302.05543*, 2023.

648 Zangwei Zheng, Pengtai Xu, Xuan Zou, Da Tang, Zhen Li, Chenguang Xi, Peng Wu, Leqi Zou, Yijie  
649 Zhu, Ming Chen, et al. Cowclip: reducing ctr prediction model training time from 12 hours to 10  
650 minutes on 1 gpu. In *Proceedings of the AAAI conference on artificial intelligence*, volume 37,  
651 pp. 11390–11398, 2023.

652 Zhi-Hua Zhou. *Ensemble methods: foundations and algorithms*. CRC press, 2012.

654 Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation  
655 using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference  
656 on computer vision*, pp. 2223–2232, 2017.

657 Imtiaz Masud Ziko, Jose Dolz, Eric Granger, and Ismail Ben Ayed. Laplacian regularized few-shot  
658 learning. In *ICML*, 2020.

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

**Ethics Statement.** We affirm adherence to the ICLR Code of Ethics. Our study uses only standard, publicly available benchmarks (miniImagenet, tieredImagenet, CUB, and an animal-face dataset referenced in prior work) under their respective licenses. No new data were collected, no personally identifiable or sensitive information is involved, and no human-subject interventions were conducted; institutional review board approval was therefore not required. The proposed method is a test-time, data-side augmentation wrapper intended to improve recognition robustness in few-shot settings; it does not require access to underlying model parameters. Potential risks include misuse of generative models to synthesize misleading content and amplification of dataset biases. To mitigate these risks, we (i) confine generation to class-faithful, small perturbations of inputs, (ii) evaluate on public benchmarks with well-documented splits, and (iii) plan to release code with conservative defaults and documentation describing appropriate use and limitations. We disclose no conflicts of interest or external sponsorship affecting the work.

**Reproducibility Statement.** We take reproducibility seriously. The paper specifies the method mathematically (Section § 4), the experimental set-up (encoders, classifier, datasets, and episodic protocol), and all evaluation details (Section § 6); ablations and qualitative analyses are provided to validate design choices. The Appendix Section §I details data preprocessing (including upscaling and shape-tweak parameters), hyperparameters (noise level, conditioning strength, denoising steps), and the exact episodic sampling procedure (5-way-1-shot, 15,000 queries, 95% confidence intervals). Upon paper acceptance, we will release a repository containing: training scripts for ProtoNet/FEAT under the stated backbones, inference scripts for our augmentation, configuration files for table, deterministic seeds, and instructions to download datasets and reproduce numbers end-to-end on a single GPU (the hardware we report). Where we use pretrained weights or models, we provide pointers or scripts to obtain them. Together, these materials enable exact regeneration of all reported tables and figures.

## A EFFICIENCY

	0.25 Noise	0.50 Noise	0.75 Noise	1.00 Noise
<b>Generation Time (s) ↓</b>	0.41	0.68	0.92	1.42

Table 5: Per-image generation time across noise levels. Higher noise entails more denoising compute. Measured on a single GPU; see experimental setup for hardware details.

We record the wall-clock running time for our inference script, and the results are reported in Table 5. All experiments are run on a single NVIDIA RTX A5000 GPU. As expected, runtime scales with the noising level. Higher noising (and more denoising steps) produces larger edits and requires longer inference, whereas lower noising is faster. As we start from the noisy image halfway, the inference cost is generally lower than that of standard diffusion process starting from pure noise.

## B HOW MUCH AND WHERE TO AUGMENT

We ablate support-only, query-only, and joint support+query augmentation under ProtoNet (Res12) on miniImagenet. See our ablation table (Table 4) for the full grid. Accuracy improves monotonically as we add a number of augmented copies to *both* support and queries (e.g., from 60.01% with no augmentation to 64.94% with +3/3, yielding a +4.93 absolute gain). Adding only support copies while leaving queries un-augmented can underperform due to distribution mismatch between prototype construction and query embeddings (e.g., with +3/0 accuracy is 61.82%, well below the 64.55% achieved when queries are matched with +3/1). This suggests that matched augmentation on both sides yields the largest benefit.

## C LIMITATIONS

Our approach has two main limitations. First, some 1S-DAug evaluation involves pretrained components, including ViTSmall Dosovitskiy et al. (2021) and SwinTiny Liu et al. (2021) encoders pretrained on ImageNet-1k Russakovsky et al. (2015) and a Stable-Diffusion-v1.5 generator Rombach

756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809

		Query			
		+0	+1	+2	+3
Support	+0	60.01 ± 0.65	60.24 ± 0.69	60.09 ± 0.70	60.31 ± 0.70
	+1	60.20 ± 0.66	62.90 ± 0.66	63.05 ± 0.67	63.16 ± 0.67
	+2	61.80 ± 0.64	64.53 ± 0.66	64.61 ± 0.66	64.87 ± 0.66
	+3	61.82 ± 0.64	64.55 ± 0.65	64.72 ± 0.66	64.94 ± 0.66

Table 6: Inductive 5-way-1-shot accuracy (mean ± 95% CI) as a function of the number of augmented copies for supports (rows) and queries (columns).

et al. (2022a). Therefore, we cannot fully rule out potential data leakage; future work will explore strict constraint within the few-shot training splits. Second, the diffusion-based augmentation introduces inference overhead compared to running the backbone alone or using classical geometric augmentations. Reducing this computational cost (e.g., via lighter generative backbones or faster denoising schedules) is an important direction for future work.

## D MORE RELATED WORK

**Few-shot Learning.** FSL methods commonly fall into metric-, model-, and augmentation-based families. *Metric-based* methods learn an embedding where queries are classified by proximity to supports or class prototypes, including Matching Networks (Vinyals et al., 2016), Prototypical Networks (Snell et al., 2017), Relation Networks (Sung et al., 2018), and episodic feature adaptation such as FEAT (Ye et al., 2020). Strong baselines refine this recipe with improved training protocols and heads, e.g., Baseline++ (Chen et al., 2019), Meta-Baseline (Chen et al., 2021), MetaOpt-Net (Lee et al., 2019), and transductive inference methods such as TPN (Liu et al., 2019b), LaplacianShot (Ziko et al., 2020), and TIM (Isken et al., 2020). *Model-based* approaches emphasize rapid parameter adaptation from few examples, e.g., gradient-based meta-learning with MAML (Finn et al., 2017), Meta-SGD (Li et al., 2017), Reptile (Nichol et al., 2018), and ANIL (Raghu et al., 2020). *Augmentation-based* approaches increase training diversity via feature or image synthesis—e.g., feature hallucination (Hariharan & Girshick, 2017b) and delta-based example synthesizers (Schwartz et al., 2018). These are primarily *training-time* techniques that rely on base-class supervision; by contrast, few-shot *test-time* augmentation must produce high-quality, class-faithful variants for unseen classes without retraining or labels.

Test-time generative augmentation for FSL remains limited. Bai et al. (2025) uses an adversarial image-to-image translator to combine the geometric “shape” of one image with the class-defining “style” of another (i.e., FUNIT (Liu et al., 2019a)) for inference-time augmentation. While a useful proof of concept, the dataset scope is narrow and failure arises on more complex, diverse categories, reflecting the difficulty of preserving content under large structural gaps.

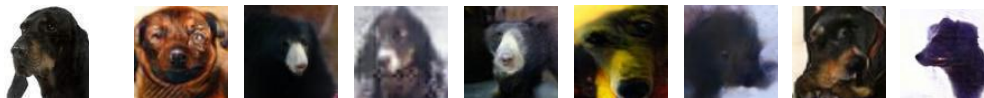
**Diffusion Models.** Early diffusion models established iterative denoising as a competitive generative paradigm (Sohl-Dickstein et al., 2015; Ho et al., 2020), with subsequent improvements to training and sampling (Nichol & Dhariwal, 2021). Latent-space diffusion amortizes computation via a learned autoencoder, enabling high-resolution synthesis (Rombach et al., 2022b). Attention-based conditioning adapters inject external signals into cross-attention without retraining the denoiser, supporting controllable editing and image-conditioned generation, including image-prompt adapters (Ye et al., 2023), general adapters (Mou et al., 2023), and control modules such as ControlNet (Zhang et al., 2023). These advances in stability, controllability, and fidelity make diffusion well-suited for few-shot test-time augmentation. Editing-by-denoising constructs variants by adding controlled noise to a source image and running the reverse process with conditioning (Meng et al., 2022). However, in this setup, too little noise yields small changes, and too much sacrifices faithfulness (e.g., changes object type), which is not suitable for data augmentation.

**GAN/Diffusion-based Data Augmentation.** Adversarial generators have long been used for data expansion and translation. Few-shot translation frameworks (e.g., FUNIT (Liu et al., 2019a) and derivatives), unpaired mappers (CycleGAN (Zhu et al., 2017)), and class-conditional generators (StyleGAN families (Karras et al., 2019)) can expand training sets but face limitations for FSL test-

---

810 **Algorithm 1** 1S-DAug (Single image  $x$ )  
811 **Require:** image  $x$ ; steps  $T$ ; schedule  $(\beta_t)$ ; user noise  $\eta \in [0, 1]$ ; conditioning weight  $\lambda_{\text{img}}$ ; number  
812 of variants  $K$ ; optional text prior  $p$   
813 1: **Geometric seed:** sample a shape tweak  $T_\psi$  and set  $x_{\text{geom}} \leftarrow T_\psi(x)$   
814 2: **Noising entry:** compute  $t_0$  from  $\eta$  via equation equation 2; draw  $x_{t_0} \sim q(\cdot | x_{\text{geom}})$  using  
815 equation equation 1  
816 3: **Working state:** set  $z_{t_0} \leftarrow x_{t_0}$  (pixel space) or  $z_{t_0} \leftarrow \text{Enc}(x_{t_0})$  (latent variant)  
817 4: **for**  $k = 1$  to  $K$  **do**  
818 5:     **for**  $t = t_0, t_0 - 1, \dots, 1$  **do**  
819 6:         Form  $K_t, V_t$  by equation equation 7 using  $x$  (and  $p$  if used); set  $Q_t = W_Q z_t$  and compute  
820  $c_t \leftarrow A_t(Q_t, K_t, V_t)$   
821 7:         **Reverse step:** update  $z_{t-1} \leftarrow \mu_\varphi(z_t, t, c_t) + \sigma_t \epsilon$  via equation equation 3, with  $\epsilon \sim$   
822  $\mathcal{N}(0, I)$   
823 8:     **end for**  
824 9:     **Decode:**  $\tilde{x}^{(k)} \leftarrow \text{Dec}(z_0)$  ▷ identity if denoising in pixel space  
825 10: **end for**  
826 11: **return**  $\{\tilde{x}^{(k)}\}_{k=1}^K$

---



827  
828  
829  
830  
831 **Intended Class** **Failure Generation (GAN)**  
832  
833  
834  
835  
836 Figure 4: Failure modes of GAN-based image-to-image translation. Examples where image-to-  
837 image GAN translation fails to preserve the intended class. Rows contain target class and failed  
838 GAN outputs with typical artifacts.  
839  
840

841 time augmentation, including training stability, mode coverage, and faithfulness for *unseen* classes  
842 without supervision. For diffusion-based generators, there have been a few recent works focusing on  
843 using the diffusion-based synthetic images for downstream tasks other than few-shot learning. How-  
844 ever, these works rely on fine-tuning with text prompts or a handful of extra target-class samples,  
845 not suitable for test-time augmentation (He et al., 2023; Benigmim et al., 2023). In contrast, we do  
846 not rely on any label or additional target-class samples, and the strict set-up fulfills the requirement  
847 for challenging downstream tasks like FSL test-time augmentation.  
848

849 **E ALGORITHM SUMMARY**

850 Algorithmic summary of 1SDAug is available as Algorithm 1.  
851  
852

853 **F GAN FAILURE CASE ILLUSTRATION**

854  
855 Figure 4 illustrates the failure cases of GAN-based image-to-image translation models (Liu et al.,  
856 2019a).  
857

858 **G PROOF OF PROPOSITION RISK DECOMPOSITION**

859  
860 Recall that  $R(g) = \frac{1}{4} \mathbb{E}[(g(x) - y)^2]$  and  $f, f_A : \mathcal{X} \rightarrow \{-1, 1\}$ ,  $\tilde{f} = \frac{1}{2}(f + f_A)$ . Since  $f^2(x) =$   
861  $f_A^2(x) = y^2 = 1$ ,  
862

$$863 R(f) = \frac{1}{4} \mathbb{E}[(f - y)^2] = \frac{1}{4} \mathbb{E}[f^2 - 2fy + y^2] = \frac{1}{2} - \frac{1}{2} \mathbb{E}[f(x)y]. \tag{15}$$

For  $\tilde{f}$ ,

$$R(\tilde{f}) = \frac{1}{4} \mathbb{E} \left[ \left( \frac{f + f_A}{2} - y \right)^2 \right] = \frac{1}{16} \mathbb{E}[(f + f_A)^2] - \frac{1}{4} \mathbb{E}[(f + f_A)y] + \frac{1}{4} \mathbb{E}[y^2]. \quad (16)$$

Expanding  $(f + f_A)^2 = f^2 + 2ff_A + f_A^2$  and using  $f^2 = f_A^2 = y^2 = 1$ ,

$$R(\tilde{f}) = \frac{1}{16} \mathbb{E}[2 + 2ff_A] - \frac{1}{4} \mathbb{E}[fy + f_A y] + \frac{1}{4} = \frac{1}{8} (1 + \mathbb{E}[f(x)f_A(x)]) - \frac{1}{4} \mathbb{E}[f(x)y] - \frac{1}{4} \mathbb{E}[f_A(x)y] + \frac{1}{4}. \quad (17)$$

Subtracting  $R(f)$  gives

$$R(\tilde{f}) - R(f) = \frac{1}{4} (\mathbb{E}[f(x)y] - \mathbb{E}[f_A(x)y]) + \frac{1}{8} (\mathbb{E}[f(x)f_A(x)] - 1). \quad (18)$$

## H ADDITIONAL THEORETICAL RESULTS

Throughout this appendix we use the pairwise reduction of Section 5.1. Each input is a query-prototype pair  $x = (q, p)$  with label  $y \in \{-1, 1\}$ . For a fixed encoder  $\Phi_\theta$  and prototype  $p$  we define the difference feature  $\Omega_\theta(x) := \Phi_\theta(q) - p$  and the Euclidean score  $g_\theta(x) := -\|\Omega_\theta(x)\|_2^2$ . The 0-1 pairwise risk of  $\theta$  is

$$R_{\text{cls}}(\theta) := \mathbb{P}(y g_\theta(x) \leq 0),$$

which coincides with the classification risk used in Section 5.3.

Where the arguments apply to generic real-valued predictors, we write  $g$  for a function in a class  $\mathcal{G}$  and specialize to the encoder score class  $\mathcal{G} = \{x \mapsto g_\theta(x) : \theta \in \Theta\}$  at the end.

### H.1 EPISODIC EUCLIDEAN MODEL AND TEST-TIME AUGMENTATION (DETAILS)

An  $N$ -way  $K$ -shot episode has support set  $S = \{(s_{c,k}, c) : c = 1, \dots, N, k = 1, \dots, K\}$  and query set  $Q = \{(q_j, y_j)\}_{j=1}^{n_q}$  with  $y_j \in \{1, \dots, N\}$ . A single encoder  $\Phi_\theta : \mathcal{X} \rightarrow \mathbb{R}^d$  produces features  $z_{c,k} := \Phi_\theta(s_{c,k})$  and  $z_q := \Phi_\theta(q)$ . Class prototypes and the Euclidean classifier are

$$p_c := \frac{1}{K} \sum_{k=1}^K z_{c,k}, \quad \hat{y}(q) := \arg \min_c \|\Phi_\theta(q) - p_c\|_2^2,$$

where  $p_c$  corresponds to  $p_c$  in Section 4. For analysis we reduce episodes to binary query-prototype pairs:  $x = (q, p)$  with label  $y \in \{-1, 1\}$ , where  $y = +1$  if  $p$  is the prototype of the true class of  $q$  and  $y = -1$  otherwise. We work with the difference feature  $\Omega_\theta(x) := \Phi_\theta(q) - p$  and Euclidean score  $g_\theta(x) := -\|\Omega_\theta(x)\|_2^2$ , and assume a uniform pre-augmentation radius bound

$$\|\Omega_\theta(x)\|_2 \leq r_0 \quad \text{for all } x \text{ and } \theta,$$

as in Section 5.1.

At test time, 1S-DAug uses support-side feature averaging and query-side logit averaging as described in Section 4, yielding (possibly aggregated) prototypes  $p_c$ . For the theory we only need query-side notation. Let  $A_k(q)$  be the  $k$ -th augmented view of a query  $q$  ( $k = 0, \dots, K_a$ ),  $z^{(k)}(q) := \Phi_\theta(A_k(q))$ , and  $(\alpha_k)_{k=0}^{K_a}$  convex weights with  $\sum_k \alpha_k = 1$ . We define the averaged query embedding

$$\bar{z}_{\text{qry}}(q) := \sum_{k=0}^{K_a} \alpha_k z^{(k)}(q),$$

the aggregated difference feature  $\bar{\Omega}_\theta(x) := \bar{z}_{\text{qry}}(q) - p$ , and the aggregated score  $\bar{g}_\theta(x) := -\|\bar{\Omega}_\theta(x)\|_2^2$ , where  $p$  denotes a (possibly aggregated) prototype. For squared Euclidean scores, logit averaging over per-view scores  $g_\theta^{(k)}(q, p) = -\|z^{(k)}(q) - p\|_2^2$  is equivalent (up to a class-independent constant) to using  $\bar{z}_{\text{qry}}(q)$ ; see Appendix H.3. We therefore express all bounds in terms of  $\bar{\Omega}_\theta(x)$  and  $\bar{g}_\theta(x)$ .

## 918 H.2 MARGIN-BASED BOUND AND RADEMACHER COMPLEXITY

919  
920 For completeness we recall the margin loss and Rademacher complexity used in Section 5. For a  
921 score function  $g_\theta$  and a pair  $(x, y) \sim D$  with  $y \in \{-1, 1\}$ , define the signed margin  $u_\theta(x) :=$   
922  $y g_\theta(x)$ . Fix a margin parameter  $\rho > 0$  and the piecewise-linear margin loss

$$923 \tau_\rho(t) := \begin{cases} 1, & t \leq 0, \\ 1 - t/\rho, & 0 < t < \rho, \\ 0, & t \geq \rho. \end{cases}$$

924  
925  
926  
927 The (population) margin risk and empirical margin risk on a sample  $S = \{(x_i, y_i)\}_{i=1}^m$  are

$$928 R_\rho(\theta) := \mathbb{E}[\tau_\rho(u_\theta(x))], \quad \widehat{R}_{S,\rho}(\theta) := \frac{1}{m} \sum_{i=1}^m \tau_\rho(y_i g_\theta(x_i)).$$

929  
930 Let the empirical 0–1 pairwise risk on  $S$  be

$$931 \widehat{R}_{\text{cls},S}(\theta) := \frac{1}{m} \sum_{i=1}^m \mathbf{1}\{y_i g_\theta(x_i) \leq 0\}.$$

932 Since  $\mathbf{1}\{t \leq 0\} \leq \tau_\rho(t)$ , we have

$$933 R_{\text{cls}}(\theta) \leq R_\rho(\theta) \quad \text{and} \quad \widehat{R}_{\text{cls},S}(\theta) \leq \widehat{R}_{S,\rho}(\theta).$$

934  
935 Let  $\mathcal{G} := \{g_\theta : \theta \in \Theta\}$  be the encoder score class. Its empirical Rademacher complexity is

$$936 \widehat{\mathfrak{R}}_S(\mathcal{G}) := \mathbb{E}_\sigma \left[ \sup_{g \in \mathcal{G}} \frac{1}{m} \sum_{i=1}^m \sigma_i g(x_i) \right],$$

937 where  $\sigma_i$  are independent Rademacher variables. A standard contraction argument with  $\tau_\rho$  yields  
938 Theorem 1 in the main text:

939 **Theorem 2** (Restatement of Theorem 1). *For any  $\rho > 0$  and  $\delta > 0$ , with probability at least  $1 - \delta$   
940 over  $S \sim D^m$ , every  $\theta$  satisfies*

$$941 R_{\text{cls}}(\theta) \leq \widehat{R}_{S,\rho}(\theta) + \frac{2}{\rho} \widehat{\mathfrak{R}}_S(\mathcal{G}) + \sqrt{\frac{\log(1/\delta)}{2m}}.$$

942  
943 **Lemma 2** (Rademacher complexity of the encoder score class). *Suppose each encoder  $\Phi_\theta$  is re-  
944 alised by a feedforward network with 1-Lipschitz nonlinearities and layer spectral norms  $\|W_\ell\|_2 \leq$   
945  $s_\ell$  such that  $\prod_{\ell=1}^L s_\ell \leq L_{\text{enc}}$ , and that  $\|\Omega_\theta(x)\|_2 \leq r$  for all  $x$  and  $\theta$ . Then there exists a constant  
946  $C_{\text{enc}} > 0$  (depending only on the architecture) such that, for any sample  $S$ ,*

$$947 \widehat{\mathfrak{R}}_S(\mathcal{G}) \leq C_{\text{enc}} L_{\text{enc}} \frac{r}{\sqrt{m}}.$$

948  
949 *In particular, in the setting of Section 5.3, one may take  $r = r_0$  before augmentation and  $r = \hat{r}$  after  
950 augmentation.*

951  
952 *Proof sketch.* By assumption, each encoder  $\Phi_\theta$  is realised by a feedforward network with 1-  
953 Lipschitz nonlinearities and layer spectral norms  $\|W_\ell\|_2 \leq s_\ell$  satisfying  $\prod_{\ell=1}^L s_\ell \leq L_{\text{enc}}$ . Standard  
954 results on spectral-norm control of deep networks (e.g. via composing linear maps and 1-Lipschitz  
955 activations) imply that  $\Phi_\theta$  is  $L_{\text{enc}}$ -Lipschitz with respect to the input  $\ell_2$ -norm, i.e.

$$956 \|\Phi_\theta(x) - \Phi_\theta(x')\|_2 \leq L_{\text{enc}} \|x - x'\|_2 \quad \text{for all } x, x'.$$

957  
958 Since  $p$  does not depend on  $q$  for a fixed pair  $x = (q, p)$ , the difference feature  $\Omega_\theta(x) = \Phi_\theta(q) - p$   
959 is also  $L_{\text{enc}}$ -Lipschitz in  $q$ . On the domain where  $\|\Omega_\theta(x)\|_2 \leq r$  for all  $x$  and  $\theta$ , the score  $g_\theta(x) =$   
960  $-\|\Omega_\theta(x)\|_2^2$  is  $2rL_{\text{enc}}$ -Lipschitz in  $x$ : the gradient of  $g_\theta$  with respect to  $\Omega_\theta$  has norm  $2\|\Omega_\theta(x)\|_2 \leq$   
961  $2r$ , and  $\Omega_\theta$  itself is  $L_{\text{enc}}$ -Lipschitz.

Let  $\mathcal{G} = \{g_\theta : \theta \in \Theta\}$  be the corresponding score class. It is therefore contained in a class of Lipschitz real-valued functions whose Lipschitz constant is bounded by  $2rL_{\text{enc}}$ . Standard covering-number/Rademacher arguments for such Lipschitz function classes (see, e.g., generic chaining bounds for  $L$ -Lipschitz functions on an  $\ell_2$ -ball) yield that there exists a constant  $C > 0$ , depending only on the input dimension and hence only on the architecture, such that for any sample  $S$  of size  $m$ ,

$$\widehat{\mathfrak{R}}_S(\mathcal{G}) \leq C (2rL_{\text{enc}}) \frac{1}{\sqrt{m}}.$$

Setting  $C_{\text{enc}} := 2C$  gives

$$\widehat{\mathfrak{R}}_S(\mathcal{G}) \leq C_{\text{enc}} L_{\text{enc}} \frac{r}{\sqrt{m}},$$

which is the claimed bound.  $\square$

### H.3 EQUIVALENCE OF LOGIT AVERAGING AND FEATURE AVERAGING

For completeness we record the standard equivalence between logit averaging and feature averaging for squared Euclidean scores.

Fix a query  $q$  and a prototype  $p$ , and let

$$z^{(k)}(q) := \Phi_\theta(A_k(q))$$

denote the encoded query under augmentation  $k$ , with convex weights  $\alpha_k \geq 0$ ,  $\sum_k \alpha_k = 1$ . The per-view scores are

$$g_\theta^{(k)}(q, p) := -\|z^{(k)}(q) - p\|_2^2, \quad \tilde{g}_\theta(q, p) := \sum_k \alpha_k g_\theta^{(k)}(q, p).$$

Define the averaged query feature

$$\bar{z}(q) := \sum_k \alpha_k z^{(k)}(q), \quad g_\theta^{\text{avg}}(q, p) := -\|\bar{z}(q) - p\|_2^2.$$

Then

$$\begin{aligned} \tilde{g}_\theta(q, p) &= -\sum_k \alpha_k \|z^{(k)}(q) - p\|_2^2 \\ &= -\sum_k \alpha_k \left( \|z^{(k)}(q)\|_2^2 - 2\langle z^{(k)}(q), p \rangle + \|p\|_2^2 \right) \\ &= -\sum_k \alpha_k \|z^{(k)}(q)\|_2^2 + 2\left\langle \sum_k \alpha_k z^{(k)}(q), p \right\rangle - \|p\|_2^2 \\ &= -\sum_k \alpha_k \|z^{(k)}(q)\|_2^2 - \|\bar{z}(q)\|_2^2 + \|\bar{z}(q)\|_2^2 + 2\langle \bar{z}(q), p \rangle - \|p\|_2^2 \\ &= -\|\bar{z}(q) - p\|_2^2 - \sum_k \alpha_k \|z^{(k)}(q)\|_2^2 + \|\bar{z}(q)\|_2^2. \end{aligned}$$

The last two terms depend on  $q$  and the set of views  $\{z^{(k)}(q)\}$  but not on  $p$ . Hence, for fixed  $q$ , comparing classes by  $\tilde{g}_\theta(q, p_c)$  is equivalent to comparing them by

$$g_\theta^{\text{avg}}(q, p_c) = -\|\bar{z}(q) - p_c\|_2^2.$$

In other words, query-side logit averaging with squared Euclidean scores induces exactly the same class ranking as nearest-prototype classification in the feature space of the averaged query embedding  $\bar{z}(q)$ .

This justifies working, for analysis, with the aggregated difference feature

$$\bar{\Omega}_\theta(x) := \bar{z}(q) - p$$

and its radius, even though the implementation performs logit averaging rather than explicit feature averaging on the query side.

#### 1026 H.4 STABILITY OF THE EMPIRICAL MARGIN RISK

1027 Let  $S = \{(x_i, y_i)\}_{i=1}^m$  be the original sample and let  $\widehat{S}$  denote the same pairs  $(x_i, y_i)$  but with the  
1029 empirical margin risk evaluated using the aggregated difference features  $\widehat{\Omega}_\theta(x_i)$ . Write

$$1030 \Omega_i := \Omega_\theta(x_i), \quad \widehat{\Omega}_i := \widehat{\Omega}_\theta(x_i),$$

1032 and recall that  $g_\theta(x) = -\|\Omega_\theta(x)\|_2^2$ . Abusing notation slightly, we write  $g_\theta(\Omega)$  for the score evalu-  
1033 ated at a difference feature  $\Omega$ , i.e.  $g_\theta(\Omega) := -\|\Omega\|_2^2$ . Assume that  $g_\theta(\cdot)$  is  $L_g$ -Lipschitz with respect  
1034 to its feature argument on the radius- $r_0$  ball in  $\mathbb{R}^d$ .

1035 **Lemma 3** (Empirical margin stability). *For any  $\rho > 0$ ,*

$$1036 |\widehat{R}_{\widehat{S},\rho}(\theta) - \widehat{R}_{S,\rho}(\theta)| \leq \frac{L_g}{\rho m} \sum_{i=1}^m \|\widehat{\Omega}_i - \Omega_i\|_2. \quad (19)$$

1039 *In particular,*

$$1040 |\widehat{R}_{\widehat{S},\rho}(\theta) - \widehat{R}_{S,\rho}(\theta)| \leq \frac{L_g}{\rho} \cdot \left( \frac{1}{m} \sum_{i=1}^m \|\widehat{\Omega}_i - \Omega_i\|_2 \right). \quad (20)$$

1044 *Proof.* For each  $i$ ,

$$1045 |\tau_\rho(y_i g_\theta(\widehat{\Omega}_i)) - \tau_\rho(y_i g_\theta(\Omega_i))| \leq \frac{1}{\rho} |y_i g_\theta(\widehat{\Omega}_i) - y_i g_\theta(\Omega_i)| = \frac{1}{\rho} |g_\theta(\widehat{\Omega}_i) - g_\theta(\Omega_i)|.$$

1048 Since  $g_\theta$  is  $L_g$ -Lipschitz in its feature argument,

$$1049 |g_\theta(\widehat{\Omega}_i) - g_\theta(\Omega_i)| \leq L_g \|\widehat{\Omega}_i - \Omega_i\|_2.$$

1051 Combining the two displays gives

$$1052 |\tau_\rho(y_i g_\theta(\widehat{\Omega}_i)) - \tau_\rho(y_i g_\theta(\Omega_i))| \leq \frac{L_g}{\rho} \|\widehat{\Omega}_i - \Omega_i\|_2.$$

1055 Averaging over  $i$  and using the triangle inequality yields equation 19.  $\square$

1057 For the Euclidean score class considered here we can take  $L_g \leq 2r_0$ , since  $z \mapsto -\|z\|_2^2$  has gradient  
1058 of norm  $2\|z\|_2$  and is therefore  $2r_0$ -Lipschitz on the radius- $r_0$  ball.

#### 1060 H.5 PROTOTYPE-TYPICALITY AND EMPIRICAL MARGIN REDUCTION

1061 To make the effect of augmentation on margins more explicit, we consider a simplified linear sur-  
1062 surrogate acting on encoder-induced features for the binary pairwise problem. Recall that each pair  
1063  $x = (q, p)$  is labelled  $y \in \{-1, +1\}$ , where  $y = +1$  denotes a ‘‘same-class’’ (correct-prototype) pair  
1064 and  $y = -1$  a ‘‘different-class’’ pair. Let  $\Omega(x) \in \mathbb{R}^d$  denote any fixed representation of pairs (for  
1065 example,  $\Omega_\theta(x)$  for a given encoder  $\Phi_\theta$ ), and let  $h(z) = w^\top z + b$  be a linear head on this feature  
1066 space. For each pairwise label  $y \in \{-1, +1\}$ , let  $\mu_y \in \mathbb{R}^d$  be a prototype of the corresponding  
1067 pairwise class in  $\Omega$ -space (e.g., the conditional mean  $\mu_y := \mathbb{E}[\Omega(x) | y]$ ), and write  $\mu_+ := \mu_{+1}$  and  
1068  $\mu_- := \mu_{-1}$ .

1069 **Lemma 4** (Prototype-aligned linear head). *Assume there exists  $\alpha > 0$  with  $w = \alpha(\mu_+ - \mu_-)$ . Then*

$$1070 h(\mu_+) - h(\mu_-) = \alpha \|\mu_+ - \mu_-\|_2^2 > 0. \quad (21)$$

1072 **Lemma 5** (Aggregation toward prototypes). *For each  $(x_i, y_i)$ , suppose the aggregated feature sat-  
1073 isfies*

$$1074 \widehat{\Omega}_i = (1 - \lambda_i) \Omega_i + \lambda_i \mu_{y_i} \quad \text{for some } \lambda_i \in [0, 1], \quad (22)$$

1075 where  $\Omega_i := \Omega(x_i)$ . Then

$$1076 \|\widehat{\Omega}_i - \mu_{y_i}\|_2 \leq \|\Omega_i - \mu_{y_i}\|_2. \quad (23)$$

1078 *Proof.* We have  $\widehat{\Omega}_i - \mu_{y_i} = (1 - \lambda_i)(\Omega_i - \mu_{y_i})$ , so  $\|\widehat{\Omega}_i - \mu_{y_i}\|_2 = (1 - \lambda_i) \|\Omega_i - \mu_{y_i}\|_2 \leq$   
1079  $\|\Omega_i - \mu_{y_i}\|_2$ .  $\square$

Define the original, prototype and aggregated margins

$$m_i^0 := y_i h(\Omega_i), \quad m_i^\mu := y_i h(\mu_{y_i}), \quad m_i^{\text{agg}} := y_i h(\widehat{\Omega}_i). \quad (24)$$

**Lemma 6** (Margin interpolation). *Under the assumptions of Lemma 5, for all  $i$ ,*

$$m_i^{\text{agg}} = (1 - \lambda_i)m_i^0 + \lambda_i m_i^\mu. \quad (25)$$

*In particular, if  $m_i^\mu \geq m_i^0$  for all  $i$ , then  $m_i^{\text{agg}} \geq m_i^0$  for all  $i$ , with strict inequality whenever  $\lambda_i > 0$  and  $m_i^\mu > m_i^0$ .*

*Proof.* By linearity,

$$h(\widehat{\Omega}_i) = h((1 - \lambda_i)\Omega_i + \lambda_i \mu_{y_i}) = (1 - \lambda_i)h(\Omega_i) + \lambda_i h(\mu_{y_i}),$$

and multiplying by  $y_i$  yields  $m_i^{\text{agg}} = (1 - \lambda_i)m_i^0 + \lambda_i m_i^\mu$ . The monotonicity statement follows immediately.  $\square$

**Proposition 2** (Empirical margin risk reduction). *Assume:*

- (i)  $h$  is linear and satisfies Lemma 4;
- (ii) For each  $i$ ,  $\widehat{\Omega}_i$  satisfies Lemma 5 for some  $\lambda_i \in [0, 1]$ ;
- (iii) For each  $i$ ,  $m_i^\mu \geq m_i^0$ .

Then for any  $\rho > 0$ ,

$$\widehat{R}_{\widehat{S}, \rho}(h) \leq \widehat{R}_{S, \rho}(h). \quad (26)$$

*If in addition there exists  $i$  with  $\lambda_i > 0$ ,  $m_i^\mu > m_i^0$ , and  $m_i^0 < \rho$ ,  $m_i^{\text{agg}} < \rho$ , then the inequality is strict.*

*Proof.* By Lemma 6,  $m_i^{\text{agg}} \geq m_i^0$  for all  $i$ . Since  $\tau_\rho$  is non-increasing,  $\tau_\rho(m_i^{\text{agg}}) \leq \tau_\rho(m_i^0)$  for all  $i$ , hence

$$\widehat{R}_{\widehat{S}, \rho}(h) = \frac{1}{m} \sum_{i=1}^m \tau_\rho(m_i^{\text{agg}}) \leq \frac{1}{m} \sum_{i=1}^m \tau_\rho(m_i^0) = \widehat{R}_{S, \rho}(h).$$

For strict inequality, if for some  $i$  we have  $\lambda_i > 0$  and  $m_i^\mu > m_i^0$ , then  $m_i^{\text{agg}} > m_i^0$ . If also  $m_i^0 < \rho$  and  $m_i^{\text{agg}} < \rho$ , then  $\tau_\rho$  is strictly decreasing on  $(0, \rho)$ , so  $\tau_\rho(m_i^{\text{agg}}) < \tau_\rho(m_i^0)$  for that  $i$  and the average strictly decreases.  $\square$

This linear surrogate analysis explains how augmentation that moves pairwise examples toward their class-typical prototypes in  $\Omega$ -space tends to increase margins and reduce empirical margin risk, which is also corroborated in prior work’s analysis Bai et al. (2025). In the Euclidean prototype model of the main text, a similar effect arises when query and support features are averaged within class and remain well aligned.

## H.6 PROBABILISTIC RADIUS REDUCTION WITH MULTIPLE AUGMENTATIONS

For each  $i = 1, \dots, m$  and  $k = 0, \dots, M$ , let  $\{\Omega_i^{(k)} \in \mathbb{R}^d\}$  be i.i.d. feature vectors with radii  $R_i^{(k)} := \|\Omega_i^{(k)}\|_2$ . Define

$$R_{\text{orig}}^{\max} := \max_{1 \leq i \leq m} R_i^{(0)}, \quad R_{\text{aug}}^{\max} := \max_{1 \leq i \leq m} \max_{1 \leq k \leq M} R_i^{(k)}. \quad (27)$$

Define the aggregated feature and corresponding radii

$$\widehat{\Omega}_i := \frac{1}{M+1} \sum_{k=0}^M \Omega_i^{(k)}, \quad \widehat{R}_i := \|\widehat{\Omega}_i\|_2, \quad \widehat{R}^{\max} := \max_{1 \leq i \leq m} \widehat{R}_i. \quad (28)$$

**Proposition 3** (Radius reduction with  $M$  augmentations). *Assume  $\{\Omega_i^{(k)}\}_{i,k}$  are i.i.d. from a continuous distribution on  $\mathbb{R}^d$ . Then*

$$\mathbb{P}(\widehat{R}^{\max} < R_{\text{orig}}^{\max}) \geq \mathbb{P}(R_{\text{aug}}^{\max} < R_{\text{orig}}^{\max}) = \frac{1}{M+1}. \quad (29)$$

1134 *Proof.* Consider all  $m(M + 1)$  radii  $\{R_i^{(k)}\}_{i,k}$ , which are i.i.d. on  $\mathbb{R}_+$ . Let  $R_{\text{all}}^{\max} := \max_{i,k} R_i^{(k)}$ .  
 1135 By continuity, this maximum is almost surely unique and, by symmetry, each of the  $m(M + 1)$  radii  
 1136 is equally likely to be the maximum. The event  $\{R_{\text{aug}}^{\max} < R_{\text{orig}}^{\max}\}$  occurs exactly when the unique  
 1137 maximum lies among the  $m$  original radii  $\{R_i^{(0)}\}$ , hence  
 1138

$$1139 \mathbb{P}(R_{\text{aug}}^{\max} < R_{\text{orig}}^{\max}) = \frac{m}{m(M + 1)} = \frac{1}{M + 1}.$$

1142 Now fix an outcome in this event. Let  $I$  be the (unique) index such that  $R_I^{(0)} = R_{\text{orig}}^{\max}$ . Then for  
 1143 every  $i$  and  $k \geq 1$ ,  $R_i^{(k)} < R_{\text{orig}}^{\max}$ , and for all  $j \neq I$ ,  $R_j^{(0)} < R_{\text{orig}}^{\max}$ . For  $i = I$ , the average of  
 1144 the  $M + 1$  vectors  $\{\Omega_I^{(k)}\}$  has strictly smaller norm than the largest of them: not all  $\Omega_I^{(k)}$  are equal  
 1145 (almost surely, by continuity), and the Euclidean norm is strictly convex, so  
 1146

$$1147 \widehat{R}_I = \left\| \frac{1}{M + 1} \sum_{k=0}^M \Omega_I^{(k)} \right\|_2 < R_{\text{orig}}^{\max}.$$

1148 For any  $j \neq I$ , each  $\|\Omega_j^{(k)}\|_2$  is strictly less than  $R_{\text{orig}}^{\max}$ , hence by the triangle inequality  
 1149

$$1150 \widehat{R}_j = \left\| \frac{1}{M + 1} \sum_{k=0}^M \Omega_j^{(k)} \right\|_2 \leq \frac{1}{M + 1} \sum_{k=0}^M \|\Omega_j^{(k)}\|_2 < R_{\text{orig}}^{\max}.$$

1151 Thus  $\widehat{R}^{\max} < R_{\text{orig}}^{\max}$  on this event, so  
 1152

$$1153 \mathbb{P}(\widehat{R}^{\max} < R_{\text{orig}}^{\max}) \geq \mathbb{P}(R_{\text{aug}}^{\max} < R_{\text{orig}}^{\max}) = \frac{1}{M + 1}.$$

□

## 1162 H.7 TRAINING-TIME VERSUS TEST-TIME AUGMENTATION

1163 Let  $P_{\text{base}}$  and  $P_{\text{novel}}$  denote the base and novel pairwise distributions on  $(x, y)$ , and let the corre-  
 1164 sponding 0–1 pairwise risks for the encoder score  $g_\theta$  be  
 1165

$$1166 R_{\text{base}}(\theta) := \mathbb{P}_{(x,y) \sim P_{\text{base}}}(y g_\theta(x) \leq 0), \quad R_{\text{novel}}(\theta) := \mathbb{P}_{(x,y) \sim P_{\text{novel}}}(y g_\theta(x) \leq 0).$$

1167 **Training-time augmentation on base classes.** Let  
 1168

$$1169 S_{\text{base}} = \{(x_i, y_i)\}_{i=1}^{m_{\text{base}}} \sim P_{\text{base}}^{m_{\text{base}}}. \quad (30)$$

1170 At training time, generate  $M_{\text{tr}}$  augmentations per example:  
 1171

$$1172 x_i^{(0)} = x_i, \quad x_i^{(k)} = A_k^{\text{tr}}(x_i), \quad k = 1, \dots, M_{\text{tr}}, \quad (31)$$

1173 with difference features  $\Omega_\theta(x_i^{(k)})$ . Assume a radius bound  
 1174

$$1175 \|\Omega_\theta(x_i^{(k)})\|_2 \leq r_{0,\text{base}}^{\text{aug}} \quad \text{for all } i, k.$$

1176 The augmented training sample is  
 1177

$$1178 S_{\text{base}}^{\text{tr-aug}} = \{(x_i^{(k)}, y_i) : i = 1, \dots, m_{\text{base}}, k = 0, \dots, M_{\text{tr}}\}, \quad (32)$$

1179 of size  $m_{\text{base}}(M_{\text{tr}} + 1)$ .  
 1180

1181 **Proposition 4** (Training-time augmentation bound on base distribution). *Under the assumptions  
 1182 above, there exists  $C_{\text{enc}} > 0$  (as in Lemma 2) such that for any  $\rho > 0$  and  $\delta > 0$ , with probability at  
 1183 least  $1 - \delta$  over  $S_{\text{base}} \sim P_{\text{base}}^{m_{\text{base}}}$ , every encoder  $\theta$  satisfies*  
 1184

$$1185 R_{\text{base}}(\theta) \leq \widehat{R}_{S_{\text{base}}^{\text{tr-aug}}, \rho}(\theta) + \frac{2C_{\text{enc}}L_{\text{enc}}}{\rho} \frac{r_{0,\text{base}}^{\text{aug}}}{\sqrt{m_{\text{base}}(M_{\text{tr}} + 1)}} + \sqrt{\frac{\log(1/\delta)}{2m_{\text{base}}(M_{\text{tr}} + 1)}}. \quad (33)$$

1188 *Proof.* Apply Theorem 1 to the score class  $\mathcal{G}$  on the augmented sample  $S_{\text{base}}^{\text{tr-aug}}$ , with sample size  
 1189  $m_{\text{base}}(M_{\text{tr}} + 1)$  and radius parameter  $r_{0,\text{base}}^{\text{aug}}$ . Lemma 2 bounds the empirical Rademacher com-  
 1190 plexity by  $C_{\text{enc}}L_{\text{enc}}r_{0,\text{base}}^{\text{aug}}/\sqrt{m_{\text{base}}(M_{\text{tr}} + 1)}$ , yielding equation 33.  $\square$   
 1191

1192 To relate base and novel risks, assume there exists a discrepancy functional  $\text{disc}_{\mathcal{G}}$  such that for all  
 1193 encoders  $\theta$ ,

$$1194 R_{\text{novel}}(\theta) \leq R_{\text{base}}(\theta) + \text{disc}_{\mathcal{G}}(P_{\text{base}}, P_{\text{novel}}). \quad (34)$$

1196 Combining equation 34 with Proposition 4 upper-bounds  $R_{\text{novel}}(\theta)$  via a term controlled by training-  
 1197 time augmentation on  $P_{\text{base}}$  plus the discrepancy.

1198 **Test-time augmentation on novel classes.** For the novel distribution, consider a labeled sample

$$1200 S_{\text{novel}} = \{(x_i, y_i)\}_{i=1}^{m_{\text{novel}}} \sim P_{\text{novel}}^{m_{\text{novel}}}. \quad (35)$$

1202 At test time, generate  $M$  augmentations per example, form aggregated difference features  $\widehat{\Omega}_{\theta}(x_i)$ ,  
 1203 and let  $\widehat{S}_{\text{novel}}$  denote the sample with these aggregated features. Let  $\widehat{r}_{\text{novel}}$  be a radius bound for the  
 1204 aggregated novel features, i.e.,  
 1205

$$1206 \|\widehat{\Omega}_{\theta}(x_i)\|_2 \leq \widehat{r}_{\text{novel}} \quad \text{for all } i.$$

1208 Applying Theorem 1 and Lemma 2 directly to  $\widehat{S}_{\text{novel}}$  yields, with probability at least  $1 - \delta$ ,

$$1210 R_{\text{novel}}(\theta) \leq \widehat{R}_{\widehat{S}_{\text{novel}},\rho}(\theta) + \frac{2C_{\text{enc}}L_{\text{enc}}}{\rho} \frac{\widehat{r}_{\text{novel}}}{\sqrt{m_{\text{novel}}}} + \sqrt{\frac{\log(1/\delta)}{2m_{\text{novel}}}}. \quad (36)$$

1213 Proposition 3 shows that, in an idealized i.i.d. model with  $M$  augmentations per example, the max-  
 1214 imum radius of novel features decreases after augmentation with probability at least  $1/(M + 1)$ ,  
 1215 which in turn tends to reduce  $\widehat{r}_{\text{novel}}$  relative to the original radius bound. Together with Proposi-  
 1216 tion 2 (for a linear surrogate) and Lemma 3, this implies that test-time augmentation tends to both  
 1217 shrink the complexity term and reduce the empirical margin risk on the target distribution  $P_{\text{novel}}$   
 1218 itself.

1219 In contrast, training-time augmentation mainly tightens the bound on  $R_{\text{base}}(\theta)$ ; any bound on  
 1220  $R_{\text{novel}}(\theta)$  obtained via equation 34 still contains the discrepancy  $\text{disc}_{\mathcal{G}}(P_{\text{base}}, P_{\text{novel}})$ , which can  
 1221 be large when base and novel classes differ substantially. This highlights the comparative advantage  
 1222 of Test-time augmentation for few-shot generalisation under distribution shift.  
 1223

## 1224 I MORE IMPLEMENTATION DETAILS

### 1226 I.1 FSL MODEL TRAINING (PROTO NET AND FEAT)

1228 We follow the public FEAT repository for episodic training and evaluation, including 5-way,  
 1229 1-shot meta-training; 15 queries/class for both train and evaluation; Euclidean distance for  
 1230 classification; Res12 as the default backbone. Key arguments (defaults shown where applic-  
 1231 able) are exposed by `train_fsl.py`: `task_setup` {dataset={miniImagenet, tieredImagenet, CUB,  
 1232 Animals}, way=5, shot=1, query=15; eval\_way=5, eval\_shot=1, eval\_query=15}; `optimization`  
 1233 {max\_epoch=400, episodes\_per\_epoch=100, num\_eval\_episodes=200, lr =  $10^{-4}$  (with pre-trained  
 1234 weights), lr\_scheduler=step, step\_size=20, gamma=0.2, momentum=0.9, weight\_decay =  $5 \cdot 10^{-4}$ };  
 1235 `model` {model\_class  $\in$  {ProtoNet, FEAT}, backbone\_class  $\in$  {ConvNet, Res12}, use\_euclidean (Eu-  
 1236 clidean distances), temperature=1 (ProtoNet)/64 (FEAT), lr\_mul=10 for the set-to-set head}. Exam-  
 1237 ple FEAT commands for Res12 on tieredImagenet use `lr =  $2 \cdot 10^{-4}$ , step_size  $\in$  {20, 40},  $\gamma = 0.5$ ,`  
 1238 and temperatures `temperature = 64, temperature2  $\in$  {64, 32}`; we mirror this recipe for our FEAT  
 1239 runs and use the same episodic protocol for ProtoNet.

1240 Concretely, in our re-trains we use:

- 1241 • **Backbones:** Res12 (all models except CUB), ConvNet (CUB).

- **ProtoNet:** `model_class = ProtoNet`, Euclidean distances, `max_epoch = 400`, `episodes_per_epoch = 100`, `lr = 1e-4` (pretrained), step scheduler with `step_size = 20`,  $\gamma = 0.2$ ; `temperature = 1` unless otherwise tuned on validation; `momentum = 0.9`, `weight_decay = 5 \cdot 10^{-4}`. (All other task/eval counts as above.)
- **FEAT:** `model_class = FEAT`, `lr = 2e-4`, `lr_mul = 10` for the Transformer head, step scheduler with `step_size` in  $\{20, 40\}$  and  $\gamma = 0.5$ ; `temperature = 64`, `temperature2`  $\in \{64, 32\}$ ; Euclidean distances enabled; same episodic counts as ProtoNet.

At evaluation, we sample 15,000 queries and report mean accuracy with 95% confidence intervals, matching the repository’s evaluation practice.

## I.2 1S-DAUG (ONE-SHOT TEST-TIME AUGMENTATION) CONFIGURATION

We implement 1S-DAug as an image-conditioned, SDEdit-style denoising pipeline with optional image-prompt adapters and a light, class-preserving geometric pre-edit (“shape tweak”). The script exposes the following arguments (defaults in `.`), which we fix across all main tables unless the ablation states otherwise:

**Core diffusion/editing.** We use `stable-diffusion-v1.5` from the `diffusers` library as the base image generator. `-noise-level`  $\in [0, 1] \cdot 0.7$ : entry point on the diffusion trajectory (larger = more rewrite, smaller = higher faithfulness); `-steps`  $\cdot 20$ : denoising steps; `-cfg`  $\cdot 9.0$ : guidance scale; `-seed` (optional) for deterministic replication; attention slicing is enabled; VAE tiling can be toggled for large images.

**Backbone generator.** `-model` `runwayml/stable-diffusion-v1-5`; images are fed at  $512 \times 512$  resolution. Benchmarks originally at  $84 \times 84$  are upsampled before augmentation using Real-ESRGAN (Wang et al., 2021).

**Image-prompt adapter (optional).** `-ip_adapter` (on/off), `-ip_repo` `h94/IP-Adapter`, `-ip_scale`  $\cdot 0.8$  controls conditioning strength.

**Shape tweaking (geometric seed).** Enabled via `-shape_aug`; parameters: `-shape_aug_rotate`  $\cdot 20^\circ$  (uniform in  $[-R, +R]$ ); `-shape_aug_stretch`  $\cdot 0.20$  (anisotropic scales  $s_x, s_y \in [1 - S, 1 + S]$ ); `-shape_aug_translate`  $\cdot 0.025$  (fraction of width/height); `-shape_aug_persp`  $\cdot 0.12$  (corner jitter fraction for a single-view projective warp). Intermediate augmented images can be saved for inspection with `-save_aug`.

**I/O and batching.** Single-image mode or directory batch mode; recursive directory traversal and extension override are supported; per-image runtime and peak memory are logged (used in our efficiency table).

**Recommended ranges (used in ablations).** Noise levels  $\in \{0.25, 0.5, 0.75, 1.0\}$ ; shape tweaks at the defaults above or slightly weaker for fine-grained datasets. When noise is very small, diversity is limited; when very large, fidelity drops unless image conditioning is active (consistent with our qualitative/quantitative ablations).

## J USE OF LARGE LANGUAGE MODELS

**Writing assistance.** Yes—large language models (LLMs) were used to aid and polish writing (e.g., improving clarity, tightening tone, harmonizing notation, and converting prose to  $\LaTeX$ ). Substantive technical content, mathematical formulations, and experimental design were authored by the authors; LLM outputs were treated as drafts and were edited for accuracy and consistency with our contributions. No text was accepted without human verification.

**Retrieval and discovery.** Yes—LLMs were used for literature discovery and organization (e.g., surfacing related work candidates, clustering themes, and drafting citation lists). All citations included in the paper were validated by the authors against the original sources; bibliographic meta-

1296 data and claims were cross-checked manually. LLMs were not used to generate experimental results  
1297 or to fabricate evidence.  
1298

1299 **Scope and safeguards.** LLMs were not used to generate, alter, or select experimental data; to  
1300 tune hyperparameters automatically; or to produce figures or tables beyond cosmetic wording. All  
1301 code and analyses were implemented and executed by the authors, and all numbers reported in the  
1302 paper come from our runs. Prompts contained only non-sensitive project information and public  
1303 references, and no proprietary or personally identifying data were included. Where LLM-assisted  
1304 text appears (e.g., phrasing of method and related-work passages), it was reviewed for factual faith-  
1305 fulness and edited for technical precision.  
1306

## 1307 K 1S-DAUG VISUALIZATION

1308

1309 Figure 5 and Figure 6 illustrate more image generation results of our proposed method.  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349

1350  
1351  
1352  
1353  
1354  
1355  
1356  
1357  
1358  
1359  
1360  
1361  
1362  
1363  
1364  
1365  
1366  
1367  
1368  
1369  
1370  
1371  
1372  
1373  
1374  
1375  
1376  
1377  
1378  
1379  
1380  
1381  
1382  
1383  
1384  
1385  
1386  
1387  
1388  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1400  
1401  
1402  
1403

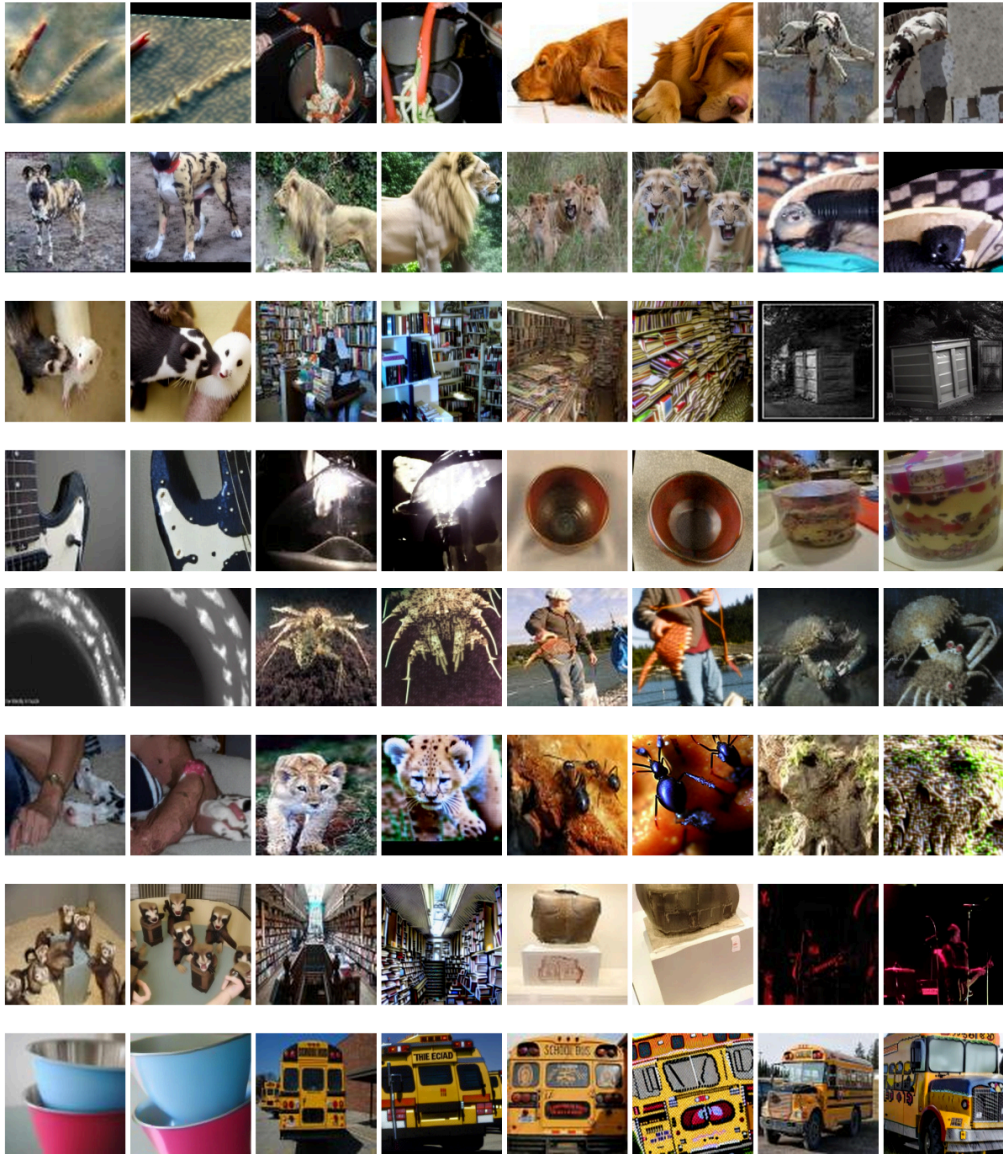


Figure 5: More qualitative results from 1S-DAug. Each pair contains the original image followed by our synthesis. All visualization pairs are random without cherry-picking.

1404  
1405  
1406  
1407  
1408  
1409  
1410  
1411  
1412  
1413  
1414  
1415  
1416  
1417  
1418  
1419  
1420  
1421  
1422  
1423  
1424  
1425  
1426  
1427  
1428  
1429  
1430  
1431  
1432  
1433  
1434  
1435  
1436  
1437  
1438  
1439  
1440  
1441  
1442  
1443  
1444  
1445  
1446  
1447  
1448  
1449  
1450  
1451  
1452  
1453  
1454  
1455  
1456  
1457

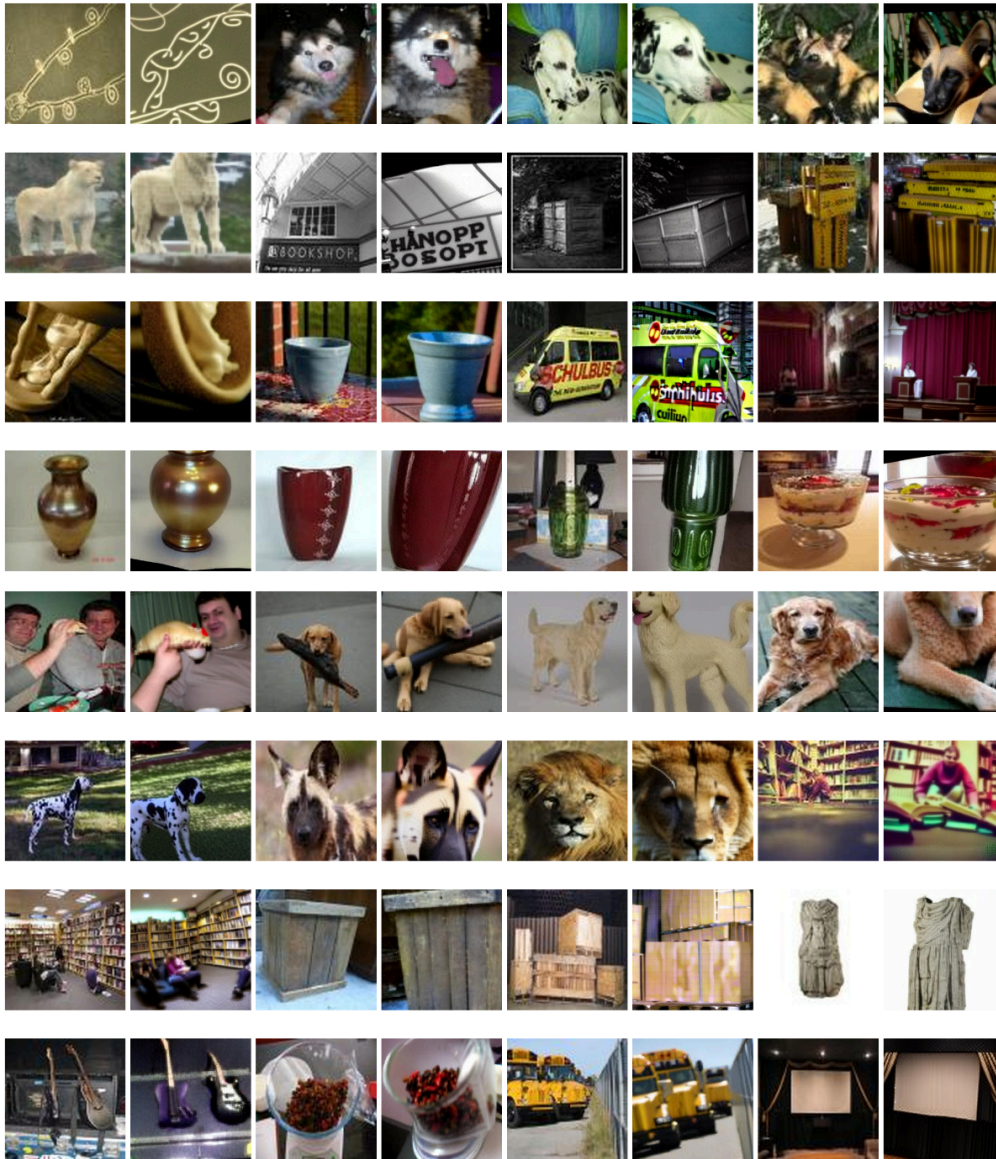


Figure 6: More qualitative results from 1S-DAug. Each pair contains the original image followed by our synthesis. All visualization pairs are random without cherry-picking.