

# FADRM: Fast and Accurate Data Residual Matching for Dataset Distillation

Jiacheng Cui<sup>\*1</sup>, Xinyue Bi<sup>\*2</sup>, Yaxin Luo<sup>1</sup>, Xiaohan Zhao<sup>1</sup>, Jiacheng Liu<sup>1</sup>, Zhiqiang Shen<sup>†1</sup>

<sup>1</sup>VILA Lab, MBZUAI    <sup>2</sup>University of Ottawa

\*Equal Contribution    †Corresponding Author

Code: [FADRM \(GitHub\)](#)

## Abstract

Residual connection has been extensively studied and widely applied at the model architecture level. However, its potential in the more challenging data-centric approaches remains unexplored. In this work, we introduce the concept of **Data Residual Matching** for the first time, leveraging data-level skip connections to facilitate data generation and mitigate data information vanishing. This approach maintains a balance between newly acquired knowledge through pixel space optimization and existing core local information identification within raw data modalities, specifically for the dataset distillation task. Furthermore, by incorporating optimization-level refinements, our method significantly improves computational efficiency, achieving superior performance while reducing training time and peak GPU memory usage by 50%. Consequently, the proposed method **Fast and Accurate Data Residual Matching for Dataset Distillation (FADRM)** establishes a new state-of-the-art, demonstrating substantial improvements over existing methods across multiple dataset benchmarks in both efficiency and effectiveness. For instance, with ResNet-18 as the student model and a 0.8% compression ratio on ImageNet-1K, the method achieves 47.7% test accuracy in single-model dataset distillation and 50.0% in multi-model dataset distillation, surpassing RDED by +5.7% and outperforming state-of-the-art multi-model approaches, EDC and CV-DD, by +1.4% and +4.0%.

## 1 Introduction

In recent years, the computer vision and natural language processing communities have predominantly focused on model-centric research, driving an unprecedented expansion in the scale of neural networks. Landmark developments such as LLMs and MLLMs in ChatGPT [25, 1], Gemini [35], DeepSeek [18] and other large-scale foundation models have shown the tremendous potential of deep learning architectures. However, as these models grow in complexity, the dependency on high-quality, richly informative datasets has become increasingly apparent, setting the stage for a paradigm shift towards data-centric approaches. Historically, the emphasis on building bigger and more complex models has often overshadowed the critical importance of the data. While model-centric strategies have

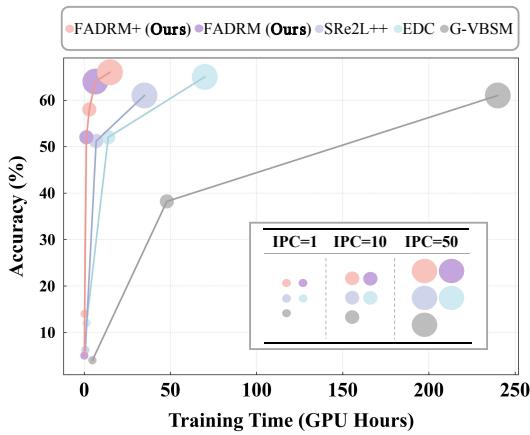


Figure 1: Total training hours on a single RTX-4090 vs. test set accuracy, comparing prior state-of-the-art methods with our proposed framework (+ denotes multi-model distillation).

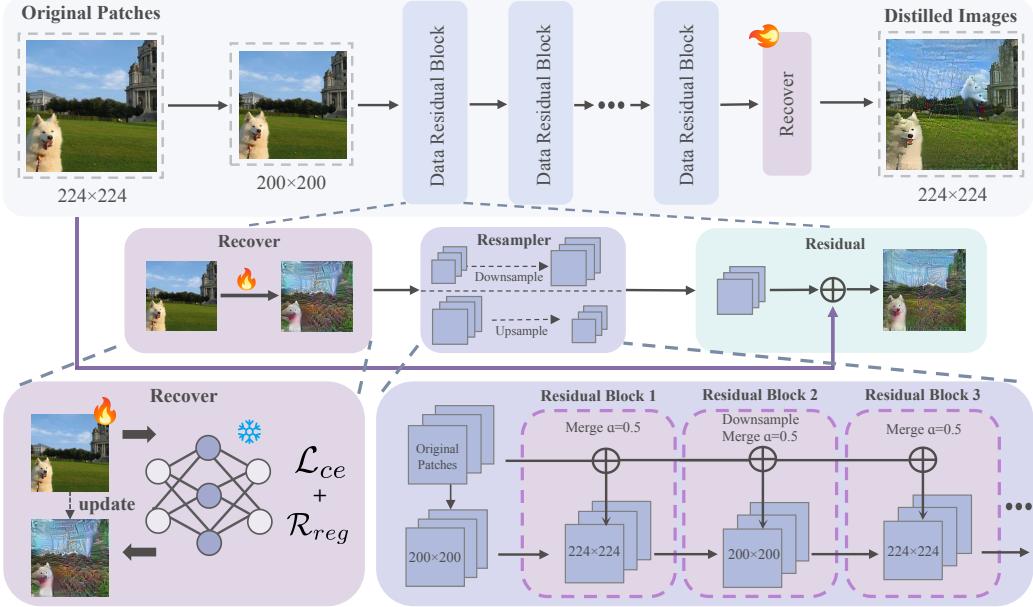


Figure 2: Overview of **FADRM**. It starts by downsampling the real data patches (both  $1 \times 1$  and  $2 \times 2$  [32] can be used as initialization and perform well in our experiments, meanwhile imposing downsampling to reduce cost). These downsampled images are subsequently processed through a series of proposed *Data Residual Blocks*. Each block utilizes a pretrained model to optimize the images within a predefined optimization budget, resamples them to a target resolution, and incorporates residual connections from the original patches via a mixing ratio  $\alpha$ . Finally, the images undergo an additional recovery stage, without residual connections, to produce the final distilled data.

delivered impressive results, they tend to overlook the benefits of optimizing data quality, which is essential for achieving higher performance with lower data demands. Recent advancements in data-centric research highlight the importance of improving information density, reducing the volume of required data, and expediting the training process of large-scale models, thus presenting a more holistic approach to performance enhancement.

Within this evolving landscape, dataset distillation [38], also called dataset condensation [14, 47, 43] has emerged as a pivotal area of research. The goal of dataset distillation is to compress large-scale datasets into smaller, highly informative subsets that retain the essential characteristics of the original data. This approach not only accelerates the training process of complex models but also mitigates the storage and computational challenges associated with massive datasets. Despite significant progress, many existing state-of-the-art methods in dataset distillation still struggle with issues related to scalability, generalization across diverse data resolutions, realism and robustness.

While residual connections have been well studied and widely implemented in the model architecture design field, primarily to prevent gradient vanishing and ensure effective feature propagation, their potential within data-centric paradigms remains largely unexplored. At the model level, residual connections help maintain the flow of gradients and enable deeper network architectures. In contrast, at the data level, similar connections can potentially prevent the loss of critical original dataset information and improve scalability and generalization across architectures during the data distillation process. This observation and design introduce a novel perspective on leveraging residual mechanisms beyond traditional model optimization, especially in the challenging domain of dataset optimization.

In this work, we introduce for the first time the concept of *Data Residual Matching* for dataset distillation. Our approach leverages data-level skip connections, a novel idea for data-centric task to prevent real data information vanishing in multi-block data synthesis architecture. We call our method **Fast and Accurate Data Residual Matching (FADRM)**, which, as shown in Fig. 2, employs a multi-resolution image recovery scheme that utilizes image resolution shrinkage and expansion in a residual manner, thereby capturing fine-grained details and facilitating the recovery of both global and local information. This balance between newly acquired knowledge through pixel space

optimization and the preservation of existing core local information within raw data modalities marks a significant advancement in dataset distillation. By integrating these data-level residual connections, our approach enhances the generalization and robustness of the distilled datasets.

Exhaustive empirical evaluations of our proposed **FADRM** on CIFAR-100 [15], Tiny-ImageNet [41], ImageNet-1k [8] and its subset demonstrate that it not only accelerates the dataset distillation process by 50% but also achieves superior accuracy that beats all previous state-of-the-art methods on both accuracy and generation speed. This approach effectively **bridges the gap between model-centric and data-centric paradigms**, providing a robust solution to the challenges inherent in high-quality data generation. Our contributions in this paper are as follows:

- We extend conventional residual connection from the model level to the data level area, and present for the first time a simple yet effective, theoretically grounded residual connection design for data generation to enhance data-centric task.
- We introduce a novel dataset distillation framework based on the proposed *data residual matching*, incorporating multi-scale residual connections in data synthesis to improve both efficiency and accuracy.
- Our approach achieves state-of-the-art results across multiple datasets, such as CIFAR-100, Tiny-ImageNet and ImageNet-1k, while being more efficient and requiring less computational cost than all previous methods.

## 2 Related Work

**Dataset Distillation** aims to synthesize a compact dataset that retains the critical information of a larger original dataset, enabling efficient training while maintaining performance comparable to the full dataset. Overall, the matching criteria include *Meta-Model Matching* [38, 24, 22, 49, 9, 12], *Gradient Matching* [47, 45, 17, 14, 48], *Trajectory Matching* [4, 7, 5, 10], Distribution Matching [46, 37, 20, 16, 26, 31, 40], and *Uni-level Global Statistics Matching* [43, 29, 30, 42, 6, 39]. Dataset distillation on large-scale datasets has recently attracted significant attention from the community. For a detailed overview, it can be referred to the newest survey works [28, 19] on this topic.

**Efficient Dataset Distillation.** Several methods improve the computational efficiency of dataset distillation. TESLA [7] accelerates MTT [4] via batched gradient computation, avoiding full graph storage and scaling to large datasets. DM [46] sidesteps bi-level optimization by directly matching feature distributions. SRe<sup>2</sup>L [43] adopts a Uni-Level Framework that aligns synthetic data with pretrained model statistics. G-VBMS [29] extends this by using lightweight model ensembles. EDC [30] further boosts efficiency through real data initialization, accelerating convergence.

**Residual Connection in Network Design.** Residual connections have played a pivotal role in advancing deep learning. Introduced in ResNet [11] to alleviate vanishing gradients, they enabled deeper networks by improving gradient flow. This idea was extended in Inception-ResNet [33] through multi-scale feature integration, and further generalized in DenseNet [13] via dense connectivity and feature reuse. Residual designs have also been central to Transformer architectures [36].

## 3 Approach

**Preliminaries.** Let the original dataset be denoted by  $\mathcal{O} = \{(x_i, y_i)\}_{i=1}^{|\mathcal{O}|}$ , and let the goal of *dataset distillation* be to construct a compact synthetic dataset  $\mathcal{C} = \{(\tilde{x}_j, \tilde{y}_j)\}_{j=1}^{|\mathcal{C}|}$ , with  $|\mathcal{C}| \ll |\mathcal{O}|$ , such that the model  $f_{\theta_C}$  trained on  $\mathcal{C}$  exhibits similar generalization behavior to the model  $f_{\theta_O}$  trained on  $\mathcal{O}$ . This objective can be formulated as minimizing the performance gap over the real data distribution:

$$\arg \min_{\mathcal{C}, |\mathcal{C}|} \sup_{(x, y) \sim \mathcal{O}} |\mathcal{L}(f_{\theta_O}(x), y) - \mathcal{L}(f_{\theta_C}(x), y)| \quad (1)$$

where the parameters  $\theta_O$  and  $\theta_C$  are obtained via empirical risk minimization:

$$\theta_O = \arg \min_{\theta} \mathbb{E}_{(x, y) \sim \mathcal{O}} [\mathcal{L}(f_{\theta}(x), y)], \quad \theta_C = \arg \min_{\theta} \mathbb{E}_{(\tilde{x}, \tilde{y}) \sim \mathcal{C}} [\mathcal{L}(f_{\theta}(\tilde{x}), \tilde{y})]. \quad (2)$$

The goal is to generate  $\mathcal{C}$  in order to maximize model performance with minimal data. Among existing methods, a notable class directly optimizes synthetic data without access to the original dataset,

referred to as *uni-level optimization*. While effective, this approach faces two key limitations: (1) progressive information loss during optimization, termed *information vanishing*, and (2) substantial computational and memory costs for large-scale synthesis, limiting real-world applicability.

**Information Vanishing.** In contrast to images distilled using bi-level frameworks, the information content in images generated by uni-level methods (e.g., EDC [30]) is fundamentally upper-bounded, as the original dataset is not utilized during synthesis (see Theorem 1). As optimization progresses, the information density initially increases but eventually deteriorates due to the accumulation of local feature loss. This degradation leads to information vanishing (see Fig. 3), which significantly reduces the fidelity of the distilled images and limits their effectiveness in downstream tasks.

**Theorem 1** (Proof in Appendix A.2). *Let  $f_\theta$  be a pretrained neural network on original dataset  $\mathcal{O}$  with fixed parameters and BatchNorm layers' mean and variance  $\mathbf{BN}^{RM}$  and  $\mathbf{BN}^{RV}$ . Let  $\tilde{x}$  denote an image optimized by minimizing the following loss:  $\mathcal{L}(\mathcal{C}) = \ell_{CE}(f_\theta(\tilde{x}), \tilde{y}) + \lambda(\sum_l \|\mu_l(\tilde{x}) - \mathbf{BN}_l^{RM}\|_2 + \sum_l \|\sigma_l^2(\tilde{x}) - \mathbf{BN}_l^{RV}\|_2)$ . Define:*

$$H(f_\theta) = \sup_{x \in \text{supp}(\mathcal{O})} H(f_\theta(x)), \quad (3)$$

as the maximum per-sample Shannon entropy of the network's output. Then, the mutual information between the optimized distilled dataset  $\mathcal{C} = \{(\tilde{x}_j, \tilde{y}_j)\}_{j=1}^{|\mathcal{C}|}$  and the original dataset  $\mathcal{O}$  is bounded by:

$$I(\mathcal{C}; \mathcal{O}) \leq |\mathcal{C}|H(f_\theta). \quad (4)$$

The insight of this theorem is that if the pretrained model  $f_\theta$  is overly confident on all inputs (low maximum entropy), then  $H(f_\theta)$  is small, and thus the distilled set, no matter how we optimize it, cannot encode a large amount of information about  $\mathcal{O}$ .

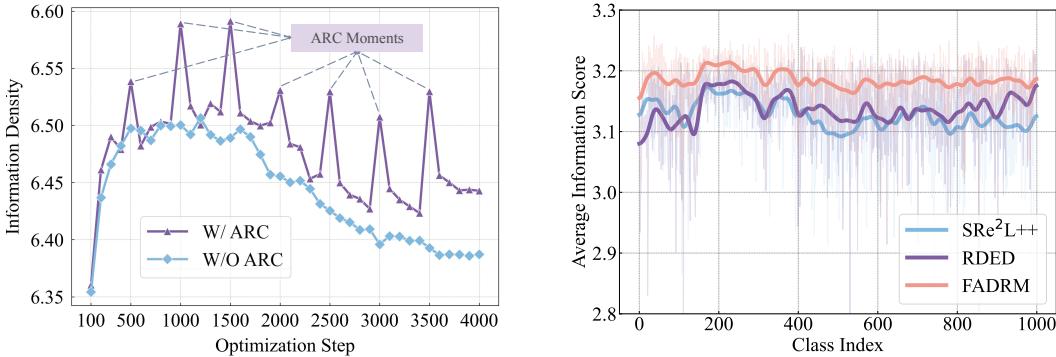


Figure 3: The above figures illustrate the phenomenon of *Information Vanishing*. The **Left** Figure shows the evolution of information density across optimization steps, quantified through feature-level entropy using a pretrained ResNet-18 [11], comparing uni-level optimization (W/O ARC) with our **FADRM** (W/ ARC). The gray lines highlight the information density enhancement achieved through residual connection. The **Right** Figure shows the comparison of information scores (higher is better) across different classes, measured by *pixel-level entropy*, among **FADRM**, SRe<sup>2</sup>L++, and RDED. All experiments are conducted on a distilled ImageNet-1k dataset with IPC=10.

**Computational Challenges.** Although uni-level frameworks exhibit scalability to large-scale datasets, the overall time required to generate a large distilled dataset remains prohibitively expensive. As illustrated in Fig. 1, EDC [30] requires nearly 70 hours to generate a 50 IPC distilled dataset, which limits its applicability in contexts involving repeated runs, large-scale data synthesis, or comprehensive empirical analysis. This motivates the need for more computationally efficient optimization strategies.

### 3.1 Overview of FADRM

The proposed **FADRM** framework, as illustrated in Fig. 2 and detailed in Algorithm 1, addresses the limitations of existing uni-level optimization frameworks by integrating three proposed components: (1) *MPT*: a mixed-precision training scheme that accelerates optimization and reduces computation

by casting model parameters to lower-precision formats, (2) *MRO*: a multiple resolution optimization that improves computational efficiency, and (3) *ARC*: an adjustable embedded residual mechanism designed to seamlessly integrate essential features from the original dataset. This framework ensures both efficiency and generation fidelity in the optimization process.

---

**Algorithm 1 FADRM: Residual Matching for Dataset Distillation**


---

**Require:** Recover model  $f_\theta$ , budget  $\mathcal{B}$ , real patches  $\mathbf{P}_s$ , merge ratio  $\alpha$ , downsampled resolutions  $D_{ds}$ , original resolutions  $D_{orig}$ , number of *ARCs*  $k$

**Ensure:** Distilled image  $\tilde{x}_{\mathcal{B}}$

- 1:  $b \leftarrow \lfloor \mathcal{B}/(k+1) \rfloor$ ,  $\tilde{x}_0 \leftarrow \text{RESAMPLE}(\mathbf{P}_s, D_{ds})$
- 2: **for**  $i = 1$  to  $k$  **do**
- 3:   **for**  $t = 1$  to  $b$  **do**
- 4:      $\tilde{x}_{(i-1)b+t} \leftarrow \text{GRADSTEP}(f_\theta, \tilde{x}_{(i-1)b+t-1})$    ▷ Optimize  $\tilde{x}$  to align the property of  $f_\theta$
- 5:   **end for**
- 6:      $\tilde{x}_{ib} \leftarrow \begin{cases} \text{RESAMPLE}(\tilde{x}_{ib}, D_{orig}), & \text{if Shape}(\tilde{x}_{ib}) = D_{ds} \\ \text{RESAMPLE}(\tilde{x}_{ib}, D_{ds}), & \text{otherwise} \end{cases}$
- 7:      $\tilde{x}_{ib} \leftarrow \alpha \tilde{x}_{ib} + (1-\alpha) \cdot \text{RESAMPLE}(\mathbf{P}_s, \text{Shape}(\tilde{x}_{ib}))$
- 8:   **end for**
- 9:   **for**  $t = 1$  to  $\mathcal{B} - kb$  **do**
- 10:      $\tilde{x}_{kb+t} \leftarrow \text{GRADSTEP}(f_\theta, \tilde{x}_{kb+t-1})$
- 11: **end for**
- 12: **return**  $\tilde{x}_{\mathcal{B}}$

---

### 3.2 Mixed Precision Training for Data Generation

Previous uni-level frameworks typically retain a fixed training pipeline, seeking efficiency through architectural or initialization-level changes. In contrast, we explicitly optimize the training process by incorporating Mixed Precision Training (MPT) [23]. Specifically, we convert the model parameters  $\theta$  from FP32 to FP16 and utilize FP16 for both logits computation and cross-entropy loss evaluation. To preserve numerical stability and ensure accurate distribution matching, we retain the computation of the divergence to the global statistics (Appendix B), as well as the gradients of the total loss with respect to  $\tilde{x}$  in FP32. By integrating *MPT*, our framework significantly reduces both GPU memory consumption and training time by approximately 50%, thereby significantly enhancing efficiency.

### 3.3 Multi-resolution Optimization

Multi-Resolution Optimization (*MRO*) enhances computational efficiency by optimizing images across multiple resolutions, unlike conventional methods that operate on a fixed input size. Naturally, low-resolution inputs can reduce computational cost for the model, they often come at the expense of performance. To mitigate this, our method periodically increases the data resolution back at specific stages, resulting in a mixed-resolution optimization process, as illustrated in Fig. 2 (bottom-right). This approach is particularly beneficial for large-scale datasets (e.g., ImageNet-1K), where direct high-resolution optimization is computationally inefficient. Notably, optimization time scales significantly with input size for large datasets but remains stable for smaller ones (input size  $\leq 64$ ). Thus, *MRO* is applied exclusively to large-scale datasets, as downscaling offers no efficiency gains for smaller ones. Specifically, given an initialized image  $\mathbf{P}_s \in \mathbb{R}^{D_{orig} \times D_{orig} \times C}$ , we first downsample it into a predefined resolution  $D_{ds}$  utilizing bilinear interpolation (detailed in Appendix C):

$$\tilde{x}_0 = \text{Resample}(\mathbf{P}_s, D_{ds}), \quad D_{ds} < D_{orig} \quad (5)$$

The downscaled images  $\tilde{x}_0$  are then optimized within a total budget  $b$ , yielding the refined version  $\tilde{x}_b$ . Subsequently,  $\tilde{x}_b$  are upscaled to their original dimensions:

$$\tilde{x}_b = \text{Resample}(\tilde{x}_b, D_{orig}) \quad (6)$$

The upscaled image  $\tilde{x}_b$  is further optimized within the same budget  $b$  to recover information lost during the downscaling and upscaling processes. This iterative procedure (downscaling optimization and upscaling optimization) is repeated until the total optimization budget  $\mathcal{B}$  is exhausted. To ensure

**MRO** achieves efficiency gain without compromising quality, selecting an appropriate  $D_{\text{ds}}$  is critical. Excessively small  $D_{\text{ds}}$  risks significant information loss, degrading distilled data's quality, while overly large  $D_{\text{ds}}$  offers negligible efficiency benefits. Thus,  $D_{\text{ds}}$  must be carefully calibrated to balance efficiency and effectiveness.

**Saved Computation by MRO.** Assume image-level convolution cost scales as  $\mathcal{O}(D^2C)$ . The baseline method performs all  $\mathcal{B}$  steps at full resolution  $D_{\text{orig}}$ , yielding:

$$\text{Cost}_{\text{baseline}} = \mathcal{B} \cdot \mathcal{O}(D_{\text{orig}}^2 C) \quad (7)$$

**FADRM** performs  $k$  alternating-resolution stages of  $b = \lfloor \mathcal{B}/(k+1) \rfloor$  steps, with approximately half at downsampled resolution  $D_{\text{ds}}$ . Let  $r = (D_{\text{ds}}/D_{\text{orig}})^2$ . The normalized cost is:

$$\frac{\text{Cost}_{\text{MRO}}}{\text{Cost}_{\text{baseline}}} = 1 - \frac{b}{\mathcal{B}} \cdot \left( \lceil \frac{k}{2} \rceil (1 - r) \right) \quad (8)$$

Under fixed  $\mathcal{B}, k$ , the cost ratio decreases linearly with  $1 - r$ . Smaller  $r$  (i.e., more aggressive downsampling) yields greater savings, but may compromise data fidelity. This trade-off highlights the role of  $r$  in balancing efficiency and representation quality during distillation.

### 3.4 Adjustable Residual Connection

In uni-level optimization, the absence of the original dataset leads to information vanishing which significantly degrades the feature representation of the distilled dataset. To mitigate this issue, we introduce Adjustable Residual Connection (*ARC*), a core mechanism that mitigates information vanishing (see Fig. 3) and improves the robustness of the distilled data (see Theroem 2). Essentially, *ARC* iteratively fuses the intermediate optimized image  $\tilde{x}_t \in \mathbb{R}^{D_t \times D_t \times C}$  at iteration  $t$  with the resized initialized data patches  $\tilde{\mathbf{P}}_t$ , which contain subtle details from the original dataset. Formally, the update rule is defined as:

$$\tilde{x}_t = \alpha \tilde{x}_t + (1 - \alpha) \text{Resample}(\mathbf{P}_s, D_t) \quad (9)$$

where  $\alpha \in [0, 1]$  is a tunable merge ratio governing the contribution of original dataset information. A smaller  $\alpha$  strengthens the integration of details from  $\mathbf{P}_s$ , whereas a larger  $\alpha$  prioritizes the preservation of the global features in the  $\tilde{x}_t$ . This trend is visualized in Fig. 4. *ARC* introduces a hyperparameter  $k$ , which determines the frequency of residual injections. Given a total optimization budget of  $\mathcal{B}$ , the training process is divided into  $k + 1$  segments, where residual connections occur after every  $b = \lfloor \mathcal{B}/(k + 1) \rfloor$  iterations. The update follows:

$$\tilde{\mathbf{P}}_{ib} = \text{Resample}(\mathbf{P}_s, D_{ib}), \quad \tilde{x}_{ib} = \alpha \tilde{x}_{ib} + (1 - \alpha) \tilde{\mathbf{P}}_{ib}. \quad (10)$$

where  $i \in \{1, 2, \dots, k\}$  denotes the index of the residual injection stage, and  $D_{ib}$  indicates the spatial resolution of the intermediate image at the corresponding iteration  $t = ib$ . The final phase consists of purely optimization without additional residual injections. Notably, *ARC* performs a per-element weighted fusion of two image tensors with negligible overhead. With a complexity of  $\mathcal{O}(H_t W_t C)$ , it scales linearly with the number of pixels and channels, making it well-suited for high-resolution data.

**Theorem 2** (Proof in Appendix A.3). *Let  $\mathcal{H}$  be a class of functions  $h : \mathbb{R}^d \rightarrow \mathbb{R}$ , and let  $h$  be Lipschitz-continuous with constant  $L_h > 0$ , and the loss function  $\ell$  be Lipschitz-continuous with constant  $L_\ell > 0$  and bounded within a finite range  $[0, B]$ . Consider: 1. Optimized perturbation added to the original data:  $\tilde{\mathcal{C}}^{\text{res}} = \{\tilde{x}_i^{\text{res}}, \tilde{y}_i^{\text{res}}\}_{i=1}^n$ . 2. residual injected dataset (FADRM):  $\tilde{\mathcal{C}}_{\text{FADRM}} = \{\tilde{x}_i, \tilde{y}_i\}_{i=1}^n$ . 3. patches selected from the original dataset:  $\mathcal{O} = \{x_i, y_i\}_{i=1}^n$ . 4. discrepancy  $\Delta := \frac{1}{n} \sum_{i=1}^n \|\tilde{x}_i^{\text{res}} - x_i\|$ . Let  $h_{\text{res}} \in \mathcal{H}$  denote the hypothesis trained on  $\tilde{\mathcal{C}}^{\text{res}}$ , and  $h_{\text{FADRM}} \in \mathcal{H}$  be trained on  $\tilde{\mathcal{C}}_{\text{FADRM}}$ .*

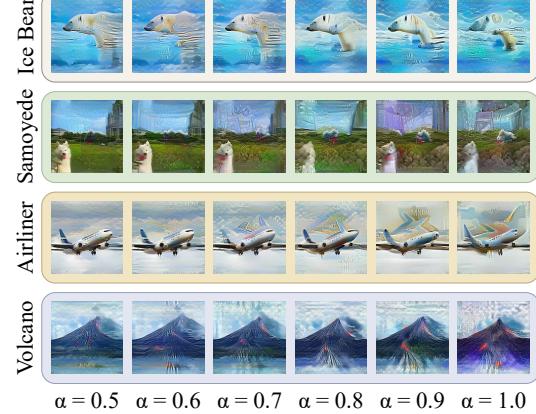


Figure 4: Visualization of the distilled images with varying merge ratios using **FADRM**.

Define the corresponding empirical risks:  $\widehat{\mathcal{L}}_{\text{res}} := \frac{1}{n} \sum_{i=1}^n \ell(h_{\text{res}}(\tilde{x}_i^{\text{res}}), \tilde{y}_i^{\text{res}})$ ,  $\widehat{\mathcal{L}}_{\text{FADRM}} := \frac{1}{n} \sum_{i=1}^n \ell(h_{\text{FADRM}}(\tilde{x}_i), \tilde{y}_i)$ . Suppose the following conditions hold:

$$\mathfrak{R}_n(\mathcal{H} \circ \mathcal{O}) - \mathfrak{R}_n(\mathcal{H} \circ \tilde{\mathcal{C}}^{\text{res}}) < -\frac{L_h \Delta (L_l + 2B\alpha)}{2B} \quad (11)$$

Where  $\alpha$  is the merge ratio. Under the condition specified in Equation (11), the generalization bound of  $h_{\text{FADRM}}$  is rigorously shown to be tighter than that of  $h_{\text{res}}$ , i.e.,

$$\widehat{\mathcal{L}}_{\text{FADRM}} + 2B \cdot \mathfrak{R}_n(\mathcal{H} \circ \tilde{\mathcal{C}}_{\text{FADRM}}) < \widehat{\mathcal{L}}_{\text{res}} + 2B \cdot \mathfrak{R}_n(\mathcal{H} \circ \tilde{\mathcal{C}}^{\text{res}}). \quad (12)$$

The insight of this theorem is that, when synthetic data is highly optimized and thus induces greater hypothesis complexity, combining it with the more structured and regular original data can lead to a tighter generalization bound, accounting for both empirical risk and Rademacher complexity.

## 4 Experiments

### 4.1 Datasets and Experimental Setup

**Datasets.** We conduct experiments across datasets with varying resolutions, including CIFAR-100 (32×32) [15], Tiny-ImageNet (64×64) [41], ImageNet-1K (224×224) [8], and their subsets.

**Baseline Methods.** To evaluate the effectiveness of our proposed framework, we conduct a comprehensive comparison against three state-of-the-art dataset distillation baselines. The first baseline is RDED [32], which selects cropped patches directly from the original dataset and is therefore categorized as involving full participation of the original data. The second method, EDC [30], retains a high degree of original data participation by optimizing selected patches with an extremely small learning rate, producing synthetic images that are close to the original samples. The third method, CV-DD [6], aligns global BatchNorm statistics with sufficient optimization by updating initialization, resulting in minimal original data involvement despite initialization from real patches. These baselines exhibit varying degrees of original data involvement, providing a solid basis for evaluating **FADRM**.

### 4.2 Main Results

**Results Analysis.** As shown in Table 1, our framework consistently achieves state-of-the-art performance across various settings. For instance, on ImageNet-1K with IPC=10 and ResNet-101 as the student model, the ensemble-enhanced variant **FADRM+** attains an accuracy of 58.1%, outperforming EDC and CV-DD by a substantial margin of +6.4%. Notably, RDED underperforms **FADRM**, underscoring the limitations of relying solely on the original dataset without further optimization. Furthermore, CV-DD is inferior to **FADRM+**, highlighting the drawbacks of largely excluding original data during synthesis. Lastly, the consistent outperformance of **FADRM+** over EDC validates the efficacy of our framework in harnessing original data via data-level residual connections.

**Efficiency Comparison.** Table 2 (Left) highlights the superior efficiency of our framework compared to existing Uni-level frameworks. Bi-level frameworks are excluded from this comparison due to their inherent limitations in scalability for large-scale datasets. Specifically, **FADRM+** achieves a reduction of 3.9 seconds per image in optimization time compared to EDC [30], culminating in a total computational saving of 54 hours when applied to the 50 IPC ImageNet-1K dataset. Similarly, **FADRM** demonstrates a 28.5 hours reduction in training time relative to SRe<sup>2</sup>L++ for the same task. Additionally, our framework significantly reduces peak memory usage compared to other frameworks, enabling efficient dataset distillation even in resource-constrained scenarios. These results underscore the scalability and computational efficiency of our approach, which not only accelerates large-scale dataset distillation but also substantially lowers associated computational costs.

### 4.3 Cross-Architecture Generalization

A fundamental criterion for evaluating the quality of distilled data is its ability to generalize across diverse network architectures, which significantly enhances its practical utility and adaptability in real-world applications. As illustrated in Table 2 (Right), **FADRM+** consistently outperforms all existing state-of-the-art methods across models of varying sizes and complexities, demonstrating superior generalization capabilities and robustness in diverse scenarios.

Dataset	IPC (Ratio)	ResNet18					ResNet50					ResNet101				
		RDED	EDC	CV-DD	FADRM	<b>FADRM+</b>	RDED	EDC	CV-DD	FADRM	<b>FADRM+</b>	RDED	EDC	CV-DD	FADRM	<b>FADRM+</b>
CIFAR-100	1 (0.2%)	17.1	39.7	28.3	31.8	<b>40.6</b>	10.9	36.1	28.7	27.3	<b>37.4</b>	11.2	32.3	29.0	29.2	<b>40.1</b>
	10 (2.0%)	56.9	63.7	62.7	67.4	<b>67.9</b>	41.6	62.1	61.5	66.5	<b>67.4</b>	54.1	61.7	63.8	68.3	<b>68.9</b>
	50 (10.0%)	66.8	68.6	67.1	71.0	<b>71.3</b>	64.0	69.4	68.2	71.5	<b>72.1</b>	67.9	68.5	67.6	71.9	<b>72.1</b>
	Whole Dataset	78.9				79.9				79.5				79.5		
Tiny-ImageNet	1 (0.2%)	11.8	39.2	30.6	28.6	<b>40.4</b>	8.2	35.9	25.1	28.4	<b>39.4</b>	9.6	40.6	28.0	27.9	<b>41.9</b>
	10 (2.0%)	41.9	51.2	47.8	48.9	<b>52.8</b>	38.4	50.2	43.8	47.3	<b>53.7</b>	22.9	51.6	47.4	47.8	<b>53.6</b>
	50 (10.0%)	58.2	57.2	54.1	56.4	<b>58.7</b>	45.6	58.8	54.7	57.0	<b>60.3</b>	41.2	58.6	54.1	57.2	<b>60.8</b>
	Whole Dataset	68.9				71.5				70.6				70.6		
ImageNette	1 (0.1%)	35.8	-	36.2	36.2	<b>39.2</b>	27.0	-	27.6	31.1	<b>31.9</b>	25.1	-	25.3	26.3	<b>29.3</b>
	10 (1.0%)	61.4	-	64.1	64.8	<b>69.0</b>	55.0	-	61.4	64.1	<b>68.1</b>	54.0	-	61.0	61.9	<b>63.7</b>
	50 (5.2%)	80.4	-	81.6	83.6	<b>84.6</b>	81.8	-	82.0	84.1	<b>85.4</b>	75.0	-	80.0	80.3	<b>82.3</b>
	Whole Dataset	93.8				89.8				89.3				89.3		
ImageWoof	1 (0.1%)	20.8	-	21.4	21.0	<b>22.8</b>	17.8	-	19.1	19.5	<b>19.9</b>	19.6	-	19.9	20.0	<b>21.8</b>
	10 (1.1%)	38.5	-	49.3	44.5	<b>57.3</b>	35.2	-	47.8	44.9	<b>54.1</b>	31.3	-	42.6	40.4	<b>51.4</b>
	50 (5.3%)	68.5	-	71.9	72.3	<b>72.6</b>	67.0	-	71.2	71.0	<b>71.7</b>	59.1	-	69.9	70.3	<b>70.6</b>
	Whole Dataset	88.2				77.8				82.7				82.7		
ImageNet-1k	1 (0.1%)	6.6	12.8	9.2	9.0	<b>14.7</b>	8.0	13.3	10.0	12.2	<b>16.2</b>	5.9	12.2	7.0	6.8	<b>14.1</b>
	10 (0.8%)	42.0	48.6	46.0	48.4	<b>50.9</b>	49.7	54.1	51.3	54.5	<b>57.5</b>	48.3	51.7	51.7	54.8	<b>58.1</b>
	50 (3.9%)	56.5	58.0	59.5	60.1	<b>61.2</b>	62.0	64.3	63.9	65.4	<b>66.9</b>	61.2	64.9	62.7	66.0	<b>67.0</b>
	Whole Dataset	72.3				78.6				79.8				79.8		

Table 1: **Post-evaluation performance comparison with SOTA baseline methods.** All experiments follow the training settings established in EDC [30]: 300 epochs for Tiny-ImageNet (IPC=10, 50), ImageNet-1K, and its subsets, and 1,000 epochs for CIFAR-100, Tiny-ImageNet (IPC=1). For fair comparison with single-model distillation (RDED) and ensemble-based methods (CV-DD, EDC), we evaluate both the single-model version (**FADRM** only utilized ResNet18 [11] for distillation) and the ensemble-enhanced version (**FADRM+**). This evaluation strategy ensures equitable benchmarking while maintaining methodological consistency across all experiments.

Method	Time Cost (s)	Peak Memory (GB)	Model	#Params	RDED	EDC	CV-DD	<b>FADRM+</b>
SRe <sup>2</sup> L++ [6]	2.52	5.3	ResNet18 [11]	11.7M	42.0	48.6	46.0	<b>50.0</b>
<b>FADRM</b>	<b>0.47</b>	<b>2.9</b>	ResNet50 [11]	25.6M	49.7	54.1	51.3	<b>57.5</b>
G-VBSM [29]	17.28	21.4	ResNet101 [11]	44.5M	48.3	51.7	51.7	<b>58.1</b>
CV-DD [6]	8.20	23.4	EfficientNet-B0 [34]	39.6M	42.8	51.1	43.2	<b>51.9</b>
EDC [30]	4.99	17.9	MobileNetV2 [27]	3.4M	34.4	45.0	39.0	<b>45.5</b>
<b>FADRM+</b>	<b>1.09</b>	<b>11.0</b>	ShuffleNetV2-0.5x [44]	1.4M	19.6	29.8	27.4	<b>30.2</b>
Swin-Tiny [21]			Swin-Tiny [21]	28.0M	29.2	38.3	-	<b>39.1</b>
Wide ResNet50-2 [11]			Wide ResNet50-2 [11]	68.9M	50.0	-	53.9	<b>59.1</b>
DenseNet121 [13]			DenseNet121 [13]	8.0M	49.4	-	50.9	<b>55.4</b>
DenseNet169 [13]			DenseNet169 [13]	14.2M	50.9	-	53.6	<b>58.5</b>
DenseNet201 [13]			DenseNet201 [13]	20.0M	49.0	-	54.8	<b>59.7</b>

Table 2: **Left:** Efficiency comparison between various optimization-based methods and our approach when distilling ImageNet-1k. The time cost is measured in seconds, representing the duration required to generate a single image on a single RTX 4090 GPU. **Right:** Top-1 accuracy (%) on ImageNet-1K for cross-architecture generalization with IPC=10.

#### 4.4 Ablation Study

**Impact of Patch Numbers for Initialization and Residuals.** To assess the effect of different patch configurations during both the initialization and residual injection stages, we conduct an ablation study, as shown in Table 3. The results suggest that both  $1 \times 1$  and  $2 \times 2$  patch settings are effective for generating distilled data. However, the  $1 \times 1$  configuration consistently delivers the better overall performance, making it the preferred choice in practice.

	IPC=10		IPC=50	
	<b>FADRM</b>	<b>FADRM+</b>	<b>FADRM</b>	<b>FADRM+</b>
$1 \times 1$	<b>48.4</b>	<b>50.9</b>	<b>60.1</b>	<b>61.2</b>
$2 \times 2$	47.7	50.0	59.8	60.1

Table 3: Comparison of student model (ResNet-18) generalization performance when trained on distilled datasets generated using  $1 \times 1$  and  $2 \times 2$  patch configurations during initialization and residual injection.

	FADRM		FADRM+		SRe <sup>2</sup> L++		G-VBSM	
	W/ MPT	W/O MPT	W/ MPT	W/O MPT	W/ MPT	W/O MPT	W/ MPT	W/O MPT
ResNet-18 (Student)	47.7 %	47.8 %	50.0 %	49.6 %	43.1 %	43.1 %	30.5 %	30.7 %
Efficiency	0.26 ms	0.63 ms	0.58 ms	0.96 ms	0.26 ms	0.63 ms	2.65 ms	4.32 ms
Peak GPU Memory	2.9 GB	5.3 GB	12 GB	23 GB	2.9 GB	5.3 GB	11.8 GB	21.4 GB

Table 4: Comparison of model generalization performance, optimization efficiency (milliseconds per image per iteration, measured under 100 batch size and 224 as input size), and peak GPU memory usage with and without mixed precision training under ImageNet-1K IPC=10.

Configuration	Accuracy (%)	Time Cost (s)	Configuration	Accuracy (%)	Time Cost (s)
<b>FADRM (W/O ARC + W/O MRO)</b>	46.4	0.52	<b>FADRM+ (W/O ARC + W/O MRO)</b>	48.7	1.16
<b>FADRM (W/O ARC + W/ MRO)</b>	46.2	0.47	<b>FADRM+ (W/O ARC + W/ MRO)</b>	48.2	1.09
<b>FADRM (W/ ARC (<math>\alpha = 0.9</math>) + W/ MRO)</b>	45.7	0.47	<b>FADRM+ (W/ ARC (<math>\alpha = 0.9</math>) + W/ MRO)</b>	48.5	1.09
<b>FADRM (W/ ARC (<math>\alpha = 0.8</math>) + W/ MRO)</b>	46.4	0.47	<b>FADRM+ (W/ ARC (<math>\alpha = 0.8</math>) + W/ MRO)</b>	48.0	1.09
<b>FADRM (W/ ARC (<math>\alpha = 0.7</math>) + W/ MRO)</b>	47.6	0.47	<b>FADRM+ (W/ ARC (<math>\alpha = 0.7</math>) + W/ MRO)</b>	48.9	1.09
<b>FADRM (W/ ARC (<math>\alpha = 0.6</math>) + W/ MRO)</b>	47.3	0.47	<b>FADRM+ (W/ ARC (<math>\alpha = 0.6</math>) + W/ MRO)</b>	49.3	1.09
<b>FADRM (W/ ARC (<math>\alpha = 0.5</math>) + W/ MRO)</b>	<b>47.7</b>	0.47	<b>FADRM+ (W/ ARC (<math>\alpha = 0.5</math>) + W/ MRO)</b>	<b>50.0</b>	1.09
<b>FADRM (W/ ARC (<math>\alpha = 0.4</math>) + W/ MRO)</b>	47.4	0.47	<b>FADRM+ (W/ ARC (<math>\alpha = 0.4</math>) + W/ MRO)</b>	49.5	1.09

Table 5: Performance comparison of ResNet-18 as the student model trained on distilled ImageNet-1K (IPC=10) datasets generated with different merge ratios ( $\alpha$ ), fixed  $D_{ds}=200$  and  $k=3$ . The efficiency is measured in seconds per image generation. **Left** presents the ablation results for single-model distillation, while **Right** shows the corresponding results for multi-model distillation.

$k$	1	2	3	4	5	6	$D_{ds}$	160	180	200	224
ImageNet-1k	47.1	47.6	<b>47.8</b>	47.6	47.4	47.3	Post Eval (%)	47.2	47.5	47.7	47.7
CIFAR-100	59.2	60.9	<b>61.5</b>	60.5	59.4	57.9	Time Cost (s)	0.42	0.44	0.47	0.52

Table 6: **Left** presents the ablation results for  $k$  (frequency of residual connections) using **FADRM** with  $D_{ds}=200$ ,  $\alpha=0.5$ , while **Right** shows the ablation results for  $D_{ds}$  on ImageNet-1k IPC=10. Efficiency is measured as the total time required to optimize a single image under a fixed budget of 2,000 optimization iterations using **FADRM** with  $\alpha=0.5$ ,  $k=3$ .

**Impact of Mixed Precision Training (MPT).** Our ablation study in Table 4 shows that MPT preserves distilled dataset quality while significantly reducing peak memory usage and improving optimization efficiency, making it an effective strategy for accelerating distillation.

**Impact of Components in FADRM.** To evaluate the impact of individual components (*MRO* and *ARC*) in our framework, we conduct a comprehensive ablation study. As demonstrated in Table 5, the results reveal that the initial integration of *MRO* results in a performance decline compared to the baseline (W/O *ARC* and W/O *MRO*). This decline is primarily attributed to the loss of crucial details during the resampling process. However, by incorporating *ARC* and reducing the merge ratio  $\alpha$  (thereby assigning higher priority to the original patches during merging), the performance is significantly enhanced compared to the baseline, which struggles with the issue of information vanishing. The optimal performance is achieved with a merge ratio of 0.5 for both **FADRM** and **FADRM+**, indicating that an equal combination of the original patch and the intermediate optimized input produces the most favorable outcomes. Crucially, the results confirm that *ARC* effectively mitigates information vanishing and compensates for missing details during the resampling process, thereby facilitating the deployment of a fast yet highly effective framework.

**Impact of Downsampled Input Size in MRO.** To determine the optimal downsampled input size ( $D_{ds}$ ) for *MRO*, we conduct an ablation study, as presented in Table 6 (Right). Our results demonstrate that  $D_{ds}=200$  achieves the most optimal performance. Notably, using other sizes leads to a degradation in the quality of the distilled dataset compared to optimizing with the original input size of 224.

**Impact of varying  $k$ .** To investigate the impact of  $k$  on the quality of the distilled dataset, we conduct an ablation study as presented in Table 6 (Left). The results demonstrate that  $k = 3$  yields the optimal configuration. Notably, we observe a positive correlation between the  $k$  and post-evaluation performance when increasing  $k$  from one to three. However, beyond  $k = 3$ , performance decreases as  $k$  increases, indicating that excessive residual connections can introduce redundant local details over global structures, which can lead to suboptimal results.

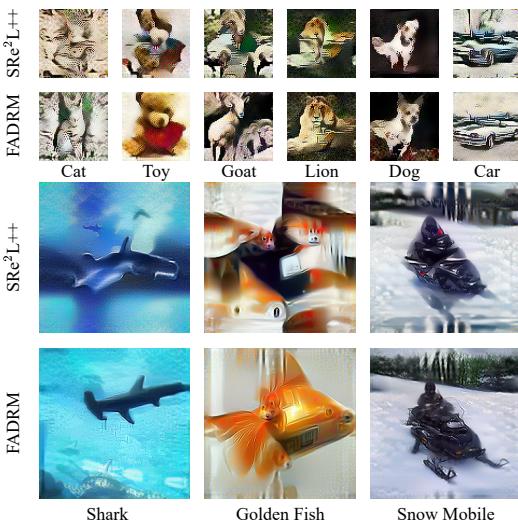


Figure 5: Visualization of dataset distilled by **FADRM** and **SRe<sup>2</sup>L++** on Tiny-ImageNet (top two rows) and ImageNet-1k (bottom two rows).

## 4.5 Distilled Image Visualization

Fig. 5 compares distilled data from **FADRM** and SRe<sup>2</sup>L++ [6], both using ResNet-18 with identical initial patch images, differing only in **FADRM**'s incorporation of residual connections. As shown, **FADRM** effectively preserves the critical features of the original patches and retains significantly more details than SRe<sup>2</sup>L++. This highlights the advantage of residual connections in enhancing information density and improving the quality of distilled data.

## 4.6 Application: Continual Learning

Leveraging continual learning to verify the effectiveness of distilled dataset generalization has been widely used in prior work [43, 29, 46]. Following these protocols and utilizing the class-incremental learning framework as in DM [46], we conduct an evaluation on Tiny-ImageNet IPC=50 using a 5-step and 10-step incremental, as shown in Fig. 6. The results indicate that **FADRM** consistently surpasses RDED, demonstrating its effectiveness.

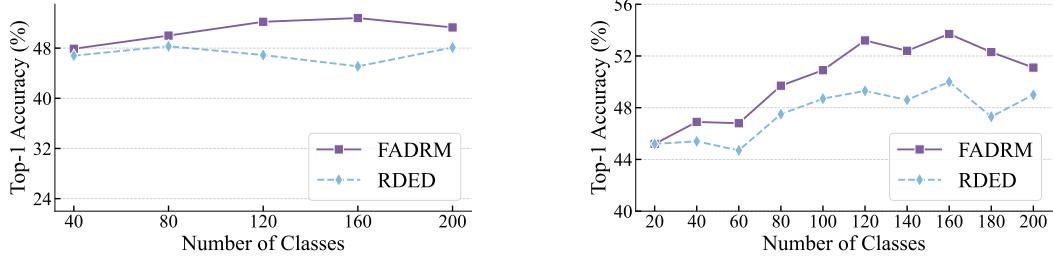


Figure 6: Five-step and Ten-step class-incremental learning on Tiny-ImageNet with IPC=50.

## 5 Conclusion

We proposed **FADRM**, a novel framework for dataset distillation designed to generate high-quality distilled datasets with significantly reduced computational overhead. Our work identifies and addresses the critical challenge of vanishing information, a fundamental limitation in *Uni-Level Framework* that heavily undermines the information density of distilled datasets. To address this, we introduce *data-level residual connections*, a novel mechanism that balances the operations of preserving critical original features and integrating new information, enriching the distilled dataset with both original and new condensed features and increasing its overall information density. Furthermore, by integrating parameter mixed precision training and input multi-resolution optimization, our framework achieves significant reductions in both Peak GPU memory consumption and training time. Extensive experiments demonstrate that **FADRM** outperforms existing state-of-the-art methods in both efficiency and accuracy across multiple benchmark datasets. For future work, we aim to extend the idea of *data-level residual connections* to broader modalities and applications of dataset distillation tasks.

## References

- [1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- [2] Peter L Bartlett and Shahar Mendelson. Rademacher and gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3(Nov):463–482, 2002.
- [3] Normand J Beaudry and Renato Renner. An intuitive proof of the data processing inequality. *arXiv preprint arXiv:1107.0740*, 2011.
- [4] George Cazenavette, Tongzhou Wang, Antonio Torralba, Alexei A Efros, and Jun-Yan Zhu. Dataset distillation by matching training trajectories. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4750–4759, 2022.

- [5] Mingyang Chen, Bo Huang, Junda Lu, Bing Li, Yi Wang, Minhao Cheng, and Wei Wang. Dataset distillation via adversarial prediction matching. *arXiv preprint arXiv:2312.08912*, 2023.
- [6] Jiacheng Cui, Zhaoyi Li, Xiaochen Ma, Xinyue Bi, Yixin Luo, and Zhiqiang Shen. Dataset distillation via committee voting. *arXiv preprint arXiv:2501.07575*, 2025.
- [7] Justin Cui, Ruochen Wang, Si Si, and Cho-Jui Hsieh. Scaling up dataset distillation to imagenet-1k with constant memory. In *International Conference on Machine Learning*, pages 6565–6590. PMLR, 2023.
- [8] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [9] Zhiwei Deng and Olga Russakovsky. Remember the past: Distilling datasets into addressable memories for neural networks. *arXiv preprint arXiv:2206.02916*, 2022.
- [10] Ziyao Guo, Kai Wang, George Cazenavette, Hui Li, Kaipeng Zhang, and Yang You. Towards lossless dataset distillation via difficulty-aligned trajectory matching. In *The Twelfth International Conference on Learning Representations*, 2024.
- [11] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [12] Yang He, Lingao Xiao, Joey Tianyi Zhou, and Ivor Tsang. Multisize dataset condensation. *ICLR*, 2024.
- [13] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4700–4708, 2017.
- [14] Jang-Hyun Kim, Jinuk Kim, Seong Joon Oh, Sangdoo Yun, Hwanjun Song, Joonhyun Jeong, Jung-Woo Ha, and Hyun Oh Song. Dataset condensation via efficient synthetic-data parameterization. In *International Conference on Machine Learning*, pages 11102–11118. PMLR, 2022.
- [15] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. Technical report, University of Toronto, Toronto, ON, Canada, 2009.
- [16] Hae Beom Lee, Dong Bok Lee, and Sung Ju Hwang. Dataset condensation with latent space knowledge factorization and sharing. *arXiv preprint arXiv:2208.10494*, 2022.
- [17] Saehyung Lee, Sanghyuk Chun, Sangwon Jung, Sangdoo Yun, and Sungroh Yoon. Dataset condensation with contrastive signals. In *International Conference on Machine Learning*, pages 12352–12364. PMLR, 2022.
- [18] Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024.
- [19] Ping Liu and Jiawei Du. The evolution of dataset distillation: Toward scalable and generalizable solutions. *arXiv preprint arXiv:2502.05673*, 2025.
- [20] Songhua Liu, Kai Wang, Xingyi Yang, Jingwen Ye, and Xinchao Wang. Dataset distillation via factorization. *Advances in Neural Information Processing Systems*, 35:1100–1113, 2022.
- [21] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10012–10022, 2021.
- [22] Noel Loo, Ramin Hasani, Alexander Amini, and Daniela Rus. Efficient dataset distillation using random feature approximation. *arXiv preprint arXiv:2210.12067*, 2022.

- [23] Paulius Micikevicius, Sharan Narang, Jonah Alben, Gregory Diamos, Erich Elsen, David Garcia, Boris Ginsburg, Michael Houston, Oleksii Kuchaiev, Ganesh Venkatesh, et al. Mixed precision training. *arXiv preprint arXiv:1710.03740*, 2017.
- [24] Timothy Nguyen, Roman Novak, Lechao Xiao, and Jaehoon Lee. Dataset distillation with infinitely wide convolutional networks. *Advances in Neural Information Processing Systems*, 34:5186–5198, 2021.
- [25] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language understanding by generative pre-training. *OpenAI Blog*, 2018.
- [26] Ahmad Sajedi, Samir Khaki, Ehsan Amjadian, Lucy Z. Liu, Yuri A. Lawryshyn, and Konstantinos N. Plataniotis. Datadam: Efficient dataset distillation with attention matching. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 17097–17107, October 2023.
- [27] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4510–4520, 2018.
- [28] Xinyi Shang, Peng Sun, Zhiqiang Shen, Tao Lin, and Jing-Hao Xue. Dataset distillation in the era of large-scale data: Methods, analysis, and future directions. 2025.
- [29] Shitong Shao, Zeyuan Yin, Muxin Zhou, Xindong Zhang, and Zhiqiang Shen. Generalized large-scale data condensation via various backbone and statistical matching. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16709–16718, 2024.
- [30] Shitong Shao, Zikai Zhou, Huanran Chen, and Zhiqiang Shen. Elucidating the design space of dataset condensation. *arXiv preprint arXiv:2404.13733*, 2024.
- [31] Donghyeok Shin, Seungjae Shin, and Il-Chul Moon. Frequency domain-based dataset distillation. *Advances in Neural Information Processing Systems*, 36, 2024.
- [32] Peng Sun, Bei Shi, Daiwei Yu, and Tao Lin. On the diversity and realism of distilled dataset: An efficient dataset distillation paradigm. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9390–9399, 2024.
- [33] Christian Szegedy, Sergey Ioffe, Vincent Vanhoucke, and Alexander Alemi. Inception-v4, inception-resnet and the impact of residual connections on learning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 31, 2017.
- [34] Mingxing Tan and Quoc Le. Efficientnetv2: Smaller models and faster training. In *International conference on machine learning*, pages 10096–10106. PMLR, 2021.
- [35] Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- [36] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [37] Kai Wang, Bo Zhao, Xiangyu Peng, Zheng Zhu, Shuo Yang, Shuo Wang, Guan Huang, Hakan Bilen, Xinchao Wang, and Yang You. Cafe: Learning to condense dataset by aligning features. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.
- [38] Tongzhou Wang, Jun-Yan Zhu, Antonio Torralba, and Alexei A Efros. Dataset distillation. *arXiv preprint arXiv:1811.10959*, 2018.
- [39] Lingao Xiao and Yang He. Are large-scale soft labels necessary for large-scale dataset distillation? *arXiv preprint arXiv:2410.15919*, 2024.

- [40] Eric Xue, Yijiang Li, Haoyang Liu, Yifan Shen, and Haohan Wang. Towards adversarially robust dataset distillation by curvature regularization. *arXiv preprint arXiv:2403.10045*, 2024.
- [41] Lian Yao, Yin Li, and Li Fei-Fei. Image classification using deep convolutional neural networks. [https://cs231n.stanford.edu/reports/2015/pdfs/yle\\_project.pdf](https://cs231n.stanford.edu/reports/2015/pdfs/yle_project.pdf), 2015. CS231n: Convolutional Neural Networks for Visual Recognition, Stanford University, Course Project Report.
- [42] Zeyuan Yin and Zhiqiang Shen. Dataset distillation via curriculum data synthesis in large data era. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856. URL <https://openreview.net/forum?id=PlaZD2nGCl>.
- [43] Zeyuan Yin, Eric Xing, and Zhiqiang Shen. Squeeze, recover and relabel: Dataset condensation at imagenet scale from a new perspective. *Advances in Neural Information Processing Systems*, 36, 2024.
- [44] Xiangyu Zhang, Xinyu Zhou, Mengxiao Lin, and Jian Sun. Shufflenet: An extremely efficient convolutional neural network for mobile devices. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6848–6856, 2018.
- [45] Bo Zhao and Hakan Bilen. Dataset condensation with differentiable siamese augmentation. In *International Conference on Machine Learning*, pages 12674–12685. PMLR, 2021.
- [46] Bo Zhao and Hakan Bilen. Dataset condensation with distribution matching. In *IEEE/CVF Winter Conference on Applications of Computer Vision, WACV 2023, Waikoloa, HI, USA, January 2-7, 2023*, 2023.
- [47] Bo Zhao, Konda Reddy Mopuri, and Hakan Bilen. Dataset condensation with gradient matching. In *International Conference on Learning Representations*, 2021.
- [48] Binglin Zhou, Linhao Zhong, and Wentao Chen. Improve cross-architecture generalization on dataset distillation. *arXiv preprint arXiv:2402.13007*, 2024.
- [49] Yongchao Zhou, Ehsan Nezhadarya, and Jimmy Ba. Dataset distillation using neural feature regression. *Advances in Neural Information Processing Systems*, 35:9813–9827, 2022.

# Appendix of FADRM

## Contents

<b>A Theoretical Derivation</b>	<b>15</b>
<b>B Optimization Details</b>	<b>18</b>
<b>C Resampling via Bilinear Interpolation</b>	<b>19</b>
<b>D Limitations</b>	<b>19</b>
<b>E Experimental Setup</b>	<b>19</b>
<b>F Hyper-Parameters Setting</b>	<b>20</b>
<b>G Additional Distilled Data Visualization</b>	<b>22</b>

## A Theoretical Derivation

### A.1 Preliminary

**Lemma 1** (Data Processing Inequality [3]). *Let  $X \rightarrow Y \rightarrow Z$  form a Markov chain. Then the mutual information between  $X$  and  $Z$  is upper bounded by that between  $X$  and  $Y$ :*

$$I(X; Z) \leq I(X; Y). \quad (13)$$

In particular, no post-processing of  $Y$  can increase the information that  $Y$  contains about  $X$ .

**Theorem 3** (Temperature-scaled KL divergence is bounded and Lipschitz-continuous). *Fix integers  $k \geq 2$  and a constant  $C > 0$ . For any temperature  $T > 0$  let*

$$z = (z_1, \dots, z_k) \in [-C, C]^k, \quad q_i^{(T)} = \frac{\exp(z_i/T)}{\sum_{j=1}^k \exp(z_j/T)} \quad (i = 1, \dots, k). \quad (14)$$

Let  $p = (p_1, \dots, p_k) \in \Delta_k$  be an arbitrary target probability vector (e.g. it may come from another soft-max with its own temperature). Define the loss

$$\ell(p, q^{(T)}(z)) := \text{KL}(p \parallel q^{(T)}(z)) = \sum_{i=1}^k p_i \log \frac{p_i}{q_i^{(T)}}. \quad (15)$$

Then the following hold:

1. (Bounded range) For every admissible pair  $(p, z)$ ,

$$0 \leq \ell(p, q^{(T)}(z)) \leq B, \quad B := \log k + \frac{2C}{T}. \quad (16)$$

2. ( $\ell_\infty$ -Lipschitz continuity in logits) The map  $z \mapsto \ell(p, q^{(T)}(z))$  is  $L$ -Lipschitz w.r.t. the  $\ell_\infty$  norm with  $L = \frac{1}{T}$ . Consequently, it is  $\sqrt{k}/T$ -Lipschitz w.r.t. the Euclidean norm.

*proof of Theorem 3.* **(i) Boundedness.** Write

$$\text{KL}(p \parallel q^{(T)}) = \sum_{i=1}^k p_i \log p_i - \sum_{i=1}^k p_i \log q_i^{(T)}. \quad (17)$$

Since  $x \mapsto x \log x$  is non-positive on  $[0, 1]$ , the first term is at most 0, so

$$\text{KL}(p \parallel q^{(T)}) \leq - \sum_{i=1}^k p_i \log q_i^{(T)}. \quad (18)$$

For the soft-max,  $\log q_i^{(T)} = z_i/T - \log Z$ , where  $Z := \sum_{j=1}^k \exp(z_j/T)$ . Hence

$$-\sum_{i=1}^k p_i \log q_i^{(T)} = -\frac{1}{T} \sum_{i=1}^k p_i z_i + \log Z. \quad (19)$$

Because each  $z_i \in [-C, C]$  and  $\sum_i p_i = 1$ ,

$$-\frac{1}{T} \sum_i p_i z_i \leq \frac{C}{T}. \quad (20)$$

Moreover,  $z_i \leq C$  implies  $Z \leq k \exp(C/T)$  and thus  $\log Z \leq \log k + \frac{C}{T}$ . Combining the two parts yields the desired upper bound  $\log k + 2C/T$ . Non-negativity of KL divergence gives the lower bound 0.

**(ii) Lipschitz continuity.** Differentiate  $\ell$  w.r.t.  $z_i$ :

$$\partial_{z_i} \ell(p, q^{(T)}(z)) = -\frac{p_i - q_i^{(T)}}{T}. \quad (21)$$

Because  $|p_i - q_i^{(T)}| \leq 1$ , we have  $|\partial_{z_i} \ell| \leq 1/T$  for every coordinate. Thus  $\|\nabla_z \ell\|_\infty \leq 1/T$ , and by the mean-value theorem,

$$|\ell(p, q^{(T)}(z)) - \ell(p, q^{(T)}(z'))| \leq \frac{1}{T} \|z - z'\|_\infty, \quad \forall z, z' \in [-C, C]^k, \quad (22)$$

so  $L = 1/T$  in the  $\ell_\infty$  norm. Since  $\|v\|_2 \leq \sqrt{k}\|v\|_\infty$ , the Euclidean Lipschitz constant is at most  $\sqrt{k}/T$ .  $\square$

**Lemma 2** (Generalization Bound via Rademacher Complexity [2]). *Let  $\mathcal{H}$  be a class of functions mapping  $\mathcal{X} \rightarrow [0, B]$ , and let  $S = \{x_1, \dots, x_n\}$  be an i.i.d. sample from distribution  $\mathcal{D}$ . Then, for any  $\delta > 0$ , with probability at least  $1 - \delta$ , the following inequality holds for all  $h \in \mathcal{H}$ :*

$$\mathbb{E}_{x \sim \mathcal{D}}[h(x)] \leq \frac{1}{n} \sum_{i=1}^n h(x_i) + 2\mathfrak{R}_n(\mathcal{H}) + B \sqrt{\frac{\log(1/\delta)}{2n}} \quad (23)$$

**Lemma 3** (Empirical Risk Proximity). *Let  $\tilde{x}_i := \alpha \tilde{x}_i^{\text{res}} + (1 - \alpha)x_i$  with  $\alpha \in (0, 1)$ , and let the corresponding datasets be  $\tilde{\mathcal{C}}^{\text{res}} := \{(\tilde{x}_i^{\text{res}}, y_i)\}_{i=1}^n$ ,  $\tilde{\mathcal{C}}_{\text{FADRM}} := \{(\tilde{x}_i, y_i)\}_{i=1}^n$ . Then for any model  $h \in \mathcal{H}$ , the empirical risk difference is bounded by a negligible value:*

$$|\hat{\mathcal{L}}_{\text{res}}(h) - \hat{\mathcal{L}}_{\text{FADRM}}(h)| \leq L_\ell L_h (1 - \alpha) \cdot \Delta_1, \quad \text{where } \Delta_1 := \frac{1}{n} \sum_{i=1}^n \|\tilde{x}_i^{\text{res}} - x_i\|. \quad (24)$$

*Proof of Lemma 3.* We begin by computing the pointwise difference in the loss:

$$|\ell(h(\tilde{x}_i^{\text{res}}), y_i) - \ell(h(\tilde{x}_i), y_i)|. \quad (25)$$

Since  $\ell$  is  $L_\ell$ -Lipschitz in the model output, and  $h$  is  $L_h$ -Lipschitz in the input, we have:

$$|\ell(h(\tilde{x}_i^{\text{res}}), y_i) - \ell(h(\tilde{x}_i), y_i)| \leq L_\ell \cdot |h(\tilde{x}_i^{\text{res}}) - h(\tilde{x}_i)| \leq L_\ell L_h \cdot \|\tilde{x}_i^{\text{res}} - \tilde{x}_i\|. \quad (26)$$

Note that:

$$\tilde{x}_i = \alpha \tilde{x}_i^{\text{res}} + (1 - \alpha)x_i \Rightarrow \tilde{x}_i^{\text{res}} - \tilde{x}_i = (1 - \alpha)(\tilde{x}_i^{\text{res}} - x_i), \quad (27)$$

so:

$$\|\tilde{x}_i^{\text{res}} - \tilde{x}_i\| = (1 - \alpha)\|\tilde{x}_i^{\text{res}} - x_i\|. \quad (28)$$

Therefore,

$$|\ell(h(\tilde{x}_i^{\text{res}}), y_i) - \ell(h(\tilde{x}_i), y_i)| \leq L_\ell L_h (1 - \alpha) \|\tilde{x}_i^{\text{res}} - x_i\|. \quad (29)$$

Averaging over  $n$  samples:

$$|\hat{\mathcal{L}}_{\text{res}}(h) - \hat{\mathcal{L}}_{\text{FADRM}}(h)| \leq \frac{1}{n} \sum_{i=1}^n L_\ell L_h (1 - \alpha) \|\tilde{x}_i^{\text{res}} - x_i\| = L_\ell L_h (1 - \alpha) \cdot \Delta_1. \quad (30)$$

$\square$

**Corollary 1** (Lipschitz Convex Combination Bound). *Let  $h : \mathbb{R}^d \rightarrow \mathbb{R}$  be an  $L$ -Lipschitz function. For any  $x, y \in \mathbb{R}^d$  and  $\alpha \in (0, 1)$ , define  $z = \alpha x + (1 - \alpha)y$ . Then:*

$$|h(z) - (\alpha h(x) + (1 - \alpha)h(y))| \leq L\alpha(1 - \alpha)\|x - y\| \quad (31)$$

In particular, this implies:

$$h(z) \leq \alpha h(x) + (1 - \alpha)h(y) + L\alpha(1 - \alpha)\|x - y\| \quad (32)$$

$$h(z) \geq \alpha h(x) + (1 - \alpha)h(y) - L\alpha(1 - \alpha)\|x - y\| \quad (33)$$

## A.2 Bounded Information in BN-Aligned Synthetic Data

*Proof of Theorem 1.* Let  $\mathcal{O}$  denote the original dataset. From it, a pretrained model  $f_\theta$  is derived, which includes BatchNorm statistics  $\{\mu_l, \sigma_l^2\}$ . Each synthetic image  $\tilde{x}_j$  in the distilled dataset  $\mathcal{C}$  is generated by minimizing an objective function depending only on  $f_\theta$  and a fixed label  $\tilde{y}_j$ .

We assume that each  $\tilde{x}_j$  is generated independently given  $f_\theta$ , and that  $f_\theta$  is a deterministic function of  $\mathcal{O}$ . Then, for each sample  $(\tilde{x}_j, \tilde{y}_j)$ , we have the Markov chain:

$$\mathcal{O} \rightarrow f_\theta \rightarrow \tilde{x}_j, \quad (34)$$

By applying Lemma 1, we get:

$$I(\tilde{x}_j; \mathcal{O}) \leq I(f_\theta; \mathcal{O}) = H(f_\theta), \quad (35)$$

Now, by the chain rule of mutual information:

$$I(\mathcal{C}; \mathcal{O}) = I(\{\tilde{x}_j, \tilde{y}_j\}_{j=1}^{|\mathcal{C}|}; \mathcal{O}) \leq \sum_{j=1}^{|\mathcal{C}|} I(\tilde{x}_j; \mathcal{O}) \leq |\mathcal{C}| \cdot H(f_\theta), \quad (36)$$

where we used the fact that  $\tilde{y}_j$  is fixed and independent of  $\mathcal{O}$  and the independence assumption across samples. Thus, the total information that the synthetic dataset  $\mathcal{C}$  can retain about the original dataset  $\mathcal{O}$  is bounded by the product of its size and the entropy of the model  $f_\theta$ .  $\square$

## A.3 ARC improves the robustness of the distilled images

*Proof of Theorem 2.* Let  $\tilde{x}_i^{res}$  be a perturbation generated via distribution (running statistics) matching and prediction (cross entropy) matching, and let  $x_i$  be a real image from the original dataset.

Define the residual-injected sample  $\tilde{x}_i$  as:

$$\tilde{x}_i := \alpha \tilde{x}_i^{res} + (1 - \alpha)x_i, \quad \alpha \in (0, 1) \quad (37)$$

Define the datasets:

- $\tilde{\mathcal{C}}^{res} = \{\tilde{x}_i^{res}, \tilde{y}_i^{res}\}_{i=1}^n$ : perturbation generated via distribution (running statistics) matching and prediction (cross entropy) matching,
- $\mathcal{O} = \{x_i, y_i\}_{i=1}^n$ : selected patches from the original dataset,
- $\tilde{\mathcal{C}}_{FADRM} = \{\tilde{x}_i, \tilde{y}_i\}_{i=1}^n$ : residual-injected dataset.

We begin by bounding the Rademacher complexity of the residual-injected dataset  $\tilde{\mathcal{C}}_{FADRM} = \{\tilde{x}_i\}_{i=1}^n$ , where  $\tilde{x}_i = \alpha \tilde{x}_i^{res} + (1 - \alpha)x_i$ , using Lemma 2 and Corollary 1.

From the definition:

$$\mathfrak{R}_n(\mathcal{H} \circ \tilde{\mathcal{C}}_{FADRM}) = \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i h(\tilde{x}_i) \right] \quad (38)$$

By Corollary 1, we have for each term:

$$h(\tilde{x}_i) \leq \alpha h(\tilde{x}_i^{res}) + (1 - \alpha)h(x_i) + \varepsilon_i, \quad \text{where } |\varepsilon_i| \leq L_h \cdot \alpha(1 - \alpha) \|\tilde{x}_i^{res} - x_i\| \quad (39)$$

Therefore:

$$\sum_{i=1}^n \sigma_i h(\tilde{x}_i) \leq \sum_{i=1}^n \sigma_i (\alpha h(\tilde{x}_i^{res}) + (1 - \alpha)h(x_i)) + \sum_{i=1}^n |\sigma_i \varepsilon_i| \quad (40)$$

Using  $|\sigma_i| = 1$ , we get:

$$\sum_{i=1}^n |\sigma_i \varepsilon_i| \leq L_h \alpha(1 - \alpha) \sum_{i=1}^n \|\tilde{x}_i^{res} - x_i\| = n \cdot L_h \alpha(1 - \alpha) \cdot \Delta \quad (41)$$

Divide by  $n$ , take supremum and expectation:

$$\mathfrak{R}_n(\mathcal{H} \circ \tilde{\mathcal{C}}_{\text{FADRM}}) \leq \alpha \cdot \mathfrak{R}_n(\mathcal{H} \circ \tilde{\mathcal{C}}^{\text{res}}) + (1 - \alpha) \cdot \mathfrak{R}_n(\mathcal{H} \circ \mathcal{O}) + L_h \alpha (1 - \alpha) \cdot \Delta \quad (42)$$

Rearrange the Inequality:

$$\mathfrak{R}_n(\mathcal{H} \circ \tilde{\mathcal{C}}_{\text{FADRM}}) - \mathfrak{R}_n(\mathcal{H} \circ \tilde{\mathcal{C}}^{\text{res}}) \leq (1 - \alpha) [\mathfrak{R}_n(\mathcal{H} \circ \mathcal{O}) - \mathfrak{R}_n(\mathcal{H} \circ \tilde{\mathcal{C}}^{\text{res}})] + L_h \alpha (1 - \alpha) \cdot \Delta \quad (43)$$

Multiply  $2B$  on both sides and add a negligible positive value  $\epsilon$  to the LHS:

$$2B \cdot [\mathfrak{R}_n(\mathcal{H} \circ \tilde{\mathcal{C}}_{\text{FADRM}}) - \mathfrak{R}_n(\mathcal{H} \circ \tilde{\mathcal{C}}^{\text{res}})] < 2B(1 - \alpha) \cdot [\mathfrak{R}_n(\mathcal{H} \circ \mathcal{O}) - \mathfrak{R}_n(\mathcal{H} \circ \tilde{\mathcal{C}}^{\text{res}})] + 2BL_h \alpha (1 - \alpha) \cdot \Delta + \epsilon \quad (44)$$

As validated in Theorem 3, when  $T > 0$ , KL-divergence becomes a bounded  $B$ -range loss, which we then apply Lemma 2 to formulate generalization error:

$$\mathcal{L}_{\text{gen}}(h) \leq \hat{\mathcal{L}}(h) + 2B \cdot \mathfrak{R}_n(\mathcal{H} \circ S) \quad (45)$$

Apply to both models:

$$\mathcal{L}_{\text{gen}}(h_{\text{res}}) \leq \hat{\mathcal{L}}_{\text{res}} + 2B \cdot \mathfrak{R}_n(\mathcal{H} \circ \tilde{\mathcal{C}}^{\text{res}}) \quad (46)$$

$$\mathcal{L}_{\text{gen}}(h_{\text{FADRM}}) \leq \hat{\mathcal{L}}_{\text{FADRM}} + 2B \cdot \mathfrak{R}_n(\mathcal{H} \circ \tilde{\mathcal{C}}_{\text{FADRM}}) \quad (47)$$

Recall the lower bound for the difference of two ERMs established in Lemma 3, we then have:

$$\hat{\mathcal{L}}_{\text{res}}(h) - \hat{\mathcal{L}}_{\text{FADRM}}(h) \geq -L_\ell L_h (1 - \alpha) \cdot \Delta, \quad \text{where } \Delta := \frac{1}{n} \sum_{i=1}^n \|\tilde{x}_i^{\text{res}} - x_i\|. \quad (48)$$

Given the assumption (11), we can then derive:

$$-L_\ell L_h (1 - \alpha) \cdot \Delta > 2B \left\{ (1 - \alpha) [\mathfrak{R}_n(\mathcal{H} \circ \mathcal{O}) - \mathfrak{R}_n(\mathcal{H} \circ \tilde{\mathcal{C}}_{\text{res}})] + L + h\alpha (1 - \alpha) \cdot \Delta \right\} + \epsilon \quad (49)$$

where the RHS in Equation (49) is the upper bound for the difference in Rademacher Complexity, we then derive the following inequality:

$$\hat{\mathcal{L}}_{\text{res}} - \hat{\mathcal{L}}_{\text{FADRM}} > 2B \cdot [\mathfrak{R}_n(\mathcal{H} \circ \tilde{\mathcal{C}}_{\text{FADRM}}) - \mathfrak{R}_n(\mathcal{H} \circ \tilde{\mathcal{C}}^{\text{res}})] \quad (50)$$

which shows:

$$\hat{\mathcal{L}}_{\text{res}} + 2B \cdot \mathfrak{R}_n(\mathcal{H} \circ \tilde{\mathcal{C}}^{\text{res}}) > \hat{\mathcal{L}}_{\text{FADRM}} + 2B \cdot \mathfrak{R}_n(\mathcal{H} \circ \tilde{\mathcal{C}}_{\text{FADRM}}) \quad (51)$$

□

## B Optimization Details

Formally, the optimization process adheres to the principle of aligning the synthesized data with both the predictive behavior and the statistical distribution captured by a pretrained model  $f_\theta$ . Specifically, given a synthesized image  $\tilde{x}_t$  at iteration  $t$ , the optimization objective is defined as:

$$\arg \min_{\tilde{x}_t} \mathcal{L}(f_\theta(\tilde{x}_t), \tilde{y}) + \mathcal{D}_{\text{global}}(\tilde{x}_t), \quad (52)$$

where  $\mathcal{L}(f_\theta(\tilde{x}_t), \tilde{y})$  enforces consistency with the target predictions, while  $\mathcal{D}_{\text{global}}(\tilde{x}_t)$  ensures alignment with the statistical distribution. Importantly, the parameters of  $f_\theta$  remain fixed throughout the optimization, and only  $\tilde{x}_t$  is updated.

The prediction alignment term is formulated as the cross-entropy loss computed over the synthesized batch:

$$\mathcal{L}(f_\theta(\tilde{x}_t), \tilde{y}) = -\frac{1}{N} \sum_{n=1}^N \sum_{i=1}^C \tilde{y}_{n,i} \log f_\theta(\tilde{x}_t)_{n,i}, \quad (53)$$

where  $N$  denotes the batch size, and  $C$  represents the total number of classes. The alignment to the distribution in pretrained model is calculated as follows:

$$\begin{aligned}\mathcal{D}_{\text{global}}(\tilde{x}_t) &= \sum_l \|\mu_l(\tilde{x}_t) - \mathbb{E}[\mu_l | \mathcal{O}]\|_2 \\ &\quad + \sum_l \|\sigma_l^2(\tilde{x}_t) - \mathbb{E}[\sigma_l^2 | \mathcal{O}]\|_2 \\ &= \sum_l \|\mu_l(\tilde{x}_t) - \mathbf{BN}_l^{\text{RM}}\|_2 \\ &\quad + \sum_l \|\sigma_l^2(\tilde{x}_t) - \mathbf{BN}_l^{\text{RV}}\|_2,\end{aligned}$$

where  $\mathcal{O}$  denotes the original dataset, and  $l$  indexes the layers of the model. The terms  $\mathbf{BN}_l^{\text{RM}}$  and  $\mathbf{BN}_l^{\text{RV}}$  correspond to the running mean and running variance of the Batch Normalization (BN) statistics at layer  $l$ . By minimizing  $\mathcal{D}_{\text{global}}(\tilde{x}_t)$ , the synthesized data is encouraged to exhibit statistical characteristics consistent with the original dataset, thereby preserving global information.

## C Resampling via Bilinear Interpolation

Given an original image  $I : \mathbb{Z}^2 \rightarrow \mathbb{R}^C$  defined on discrete pixel coordinates, the continuous extension  $\tilde{I} : \mathbb{R}^2 \rightarrow \mathbb{R}^C$  at non-integer location  $(i', j') \in \mathbb{R}^2$  is computed via bilinear interpolation as follows:

$$\tilde{I}(i', j') = \sum_{m=0}^1 \sum_{n=0}^1 w_{m,n} \cdot I(i+m, j+n), \quad (54)$$

where  $i = \lfloor i' \rfloor$ ,  $j = \lfloor j' \rfloor$ ,  $\alpha = i' - i \in [0, 1)$ ,  $\beta = j' - j \in [0, 1)$ , and the interpolation weights are defined by:

$$w_{m,n} = (1 - m + (-1)^m \alpha)(1 - n + (-1)^n \beta). \quad (55)$$

Explicitly, Equation (54) expands to:

$$\begin{aligned}\tilde{I}(i', j') &= (1 - \alpha)(1 - \beta) \cdot I(i, j) + \alpha(1 - \beta) \cdot I(i+1, j) \\ &\quad + (1 - \alpha)\beta \cdot I(i, j+1) + \alpha\beta \cdot I(i+1, j+1),\end{aligned} \quad (56)$$

This interpolation scheme can be viewed as a separable approximation to the continuous image function, with weights derived from tensor-product linear basis functions over the unit square. It preserves differentiability with respect to the fractional coordinates  $(i', j')$ , making it particularly amenable to gradient-based optimization frameworks.

## D Limitations

While FADRM offers substantial improvements in computational efficiency and performance for dataset distillation, it also introduces several limitations. First, the method relies on the assumption that residual signals between synthetic and real data capture critical learning dynamics, which may not generalize across domains with highly abstract or non-visual modalities such as natural language or time-series data. Second, the use of distilled datasets can inadvertently reinforce biases present in the original data if not carefully audited, potentially leading to fairness concerns in downstream applications. From a broader societal perspective, while FADRM reduces the computation and resource demands of training large models, thereby contributing positively to sustainability, it may also facilitate the deployment of powerful models in low-resource or surveillance scenarios without adequate ethical oversight. Thus, responsible deployment and continued research into bias mitigation and cross-domain generalization are essential to ensure the safe and equitable application of FADRM.

## E Experimental Setup

Our method strictly follows the training configuration established in EDC to ensure a fair and consistent comparison across all evaluated approaches. Additionally, we re-run RDED and CV-DD

under the same configuration and report the highest performance obtained between their original setup and the EDC configuration. This methodology guarantees a rigorous and equitable evaluation by accounting for potential variations in training dynamics across different settings.

To establish an upper bound on performance across different backbone architectures (representing the results achieved when training models on the full original dataset) we adopt the hyperparameters specified in Table 7. These hyperparameters are carefully chosen to ensure full model convergence while effectively mitigating the risk of overfitting, thereby providing a reliable reference for evaluating the performance of distilled datasets.

Hyperparameters for Training the Original Dataset	
Optimizer	SGD
Learning Rate	0.1
Weight Decay	1e-4
Momentum	0.9
Batch Size	128
Loss Function	Cross-Entropy
Epochs	300
Augmentation	RandomResizedCrop, Horizontal Flip, CutMix

Table 7: Hyperparameters for Training the Original Dataset.

## F Hyper-Parameters Setting

In summary, the synthesis of distilled data follows consistent hyperparameter configurations, as outlined in Table 8. Variations in hyperparameters are introduced exclusively during two phases: (1) the model Pre-training phase. and (2) the post-evaluation phase. These adjustments are carefully tailored based on the scale of the models and the specific characteristics of the datasets used. During the post-evaluation phase, we evaluate a total of four hyperparameter combinations, as detailed in Table 9. Among these, the parameter  $\eta$  plays a critical role in controlling the decay rate of the learning rate, as defined by the cosine learning rate schedule in Equation 57. Specifically, a larger value of  $\eta$  results in a slower decay rate, thereby preserving a higher learning rate for a longer duration during training.

$$\text{Learning Rate} = 0.5 \times \left( 1 + \cos \left( \pi \frac{\text{step}}{\text{epochs} \times \eta} \right) \right) \quad (57)$$

### F.1 CIFAR-100

This subsection outlines the hyperparameter configurations employed in the CIFAR-100 experiments, providing the necessary details to ensure reproducibility in future research.

**Pre-training phase.** Table 10 provides a comprehensive summary of the hyperparameters employed for training the models on the original CIFAR-100 dataset for generating the distilled dataset.

**Evaluation Phase.** Table 11 outlines the hyperparameter configurations employed for the post-evaluation phase on the Distilled CIFAR-100 dataset.

### F.2 Tiny-ImageNet

This part describes the hyperparameter settings used in the Tiny-ImageNet experiments, offering comprehensive details to facilitate reproducibility for future studies.

**Pre-training phase.** Table 12 presents a detailed overview of the hyperparameters used for model training on the original Tiny-ImageNet dataset.

**Evaluation Phase.** Table 13 details the hyperparameter settings used during the post-evaluation phase on the Distilled Tiny-ImageNet dataset.

Hyperparameter	Value
Optimizer	Adam
Learning rate	0.25
Beta	(0.5, 0.9)
Epsilon	$1 \times 10^{-8}$
Batch Size	100 or 10 (if $C < 100$ )
Iterations Budgets ( $\mathcal{B}$ )	2,000
Merge Ratio ( $\alpha$ )	0.5
Number of ARC ( $k$ )	3
Downsampled Size ( $D_{ds}$ )	200 (ImageNet-1k and Its subsets), Original Input Size (CIFAR-100, Tiny-ImageNet)
<b>FADRM Model (<math>R</math>)</b>	ResNet18
<b>FADRM+ Model (<math>R</math>)</b>	ResNet18 DenseNet121 Shuf- fleNetV2 MobileNetV2
Scheduler	Cosine Annealing
Augmentation	RandomResizedCrop, Horizontal Flip

Table 8: Hyperparameters for generating the distilled datasets.

Setting	Learning Rate	$\eta$
S1	0.001	1
S2	0.001	2
S3	0.0005	1
S4	0.0005	2

Table 9: Hyperparameter settings with learning rate and  $\eta$ .

### E.3 ImageNette

This subsection describes the hyperparameter settings utilized in the ImageNette experiments, offering detailed information to facilitate reproducibility for subsequent studies.

**Pre-training phase.** Table 14 summarizes the hyperparameters used for training models on the original ImageNette dataset to generate the distilled dataset, ensuring clarity and reproducibility.

**Evaluation Phase.** Table 15 details the hyperparameter settings applied during the post-evaluation phase on the Distilled ImageNette dataset.

### E.4 ImageWoof

This section describes the hyperparameter settings used in the ImageWoof experiments, offering detailed information to facilitate reproducibility for future studies.

**Pre-training phase.** Table 16 presents a detailed overview of the hyperparameters utilized for training models on the original ImageWoof dataset to produce the distilled dataset.

**Evaluation Phase.** Table 17 presents the hyperparameter settings utilized during the post-evaluation stage on the Distilled Imagewoof dataset, detailing the configurations applied for performance assessment.

### E.5 ImageNet-1k

This subsection outlines the hyperparameter configurations employed in the ImageNet1k experiments, providing the necessary details to ensure reproducibility in future research.

**Pre-training phase.** For ImageNet-1K, we employed the official PyTorch pretrained models, which have been extensively trained on the full ImageNet-1K dataset.

Hyperparameters for Model Pre-training	
Optimizer	SGD
Learning Rate	0.1
Weight Decay	1e-4
Momentum	0.9
Batch Size	128
Epoch	50
Scheduler	Cosine Annealing
Augmentation	RandomCrop, Horizontal Flip
Loss Function	Cross-Entropy

Table 10: Hyperparameters for CIFAR-100 Pre-trained Models.

Hyperparameters for Post-Eval on R18, R50 and R101	
Optimizer	Adamw
S1	IPC1 (R50), IPC50 (R18,R50)
S2	IPC10 (R18, R50)
S3	IPC1 (R101), IPC10 (R101), IPC50 (R101)
S4	IPC1 (R18)
Soft Label Generation	BSSL
Loss Function	KL-Divergence
Batch Size	16
Epochs	1000
Augmentation	RandomResizedCrop, Horizontal Flip, CutMix

Table 11: Hyperparameters for post-evaluation task on ResNet18, ResNet50 and ResNet101 for CIFAR-100.

Hyperparameters for Model Pre-training	
Optimizer	SGD
Learning Rate	0.1
Weight Decay	1e-4
Momentum	0.9
Batch Size	64
Epoch	150
Scheduler	Cosine Annealing
Augmentation	RandomCrop, Horizontal Flip
Loss Function	Cross-Entropy

Table 12: Hyperparameters for Tiny-ImageNet Pre-trained Models.

**Evaluation Phase.** Table 18 provides a detailed overview of the hyperparameter settings used during the post-evaluation phase on the Distilled ImageNet-1k dataset.

## G Additional Distilled Data Visualization

Additional visualizations of the distilled data generated by **FADRM** are provided in Fig. 7 (CIFAR-100), Fig. 9 (Tiny-ImageNet), Fig. 11 (ImageNette), Fig. 13 (ImageWoof), and Fig. 15 (ImageNet-1K). Furthermore, enhanced versions - **FADRM+** are presented in Fig. 8 (CIFAR-100), Fig. 10 (Tiny-ImageNet), Fig. 12 (ImageNette), Fig. 14 (ImageWoof), and Fig. 16 (ImageNet-1K).

<b>Hyperparameters for Post-Eval on R18, R50 and R101</b>	
Optimizer	Adamw
S1	IPC50 (R18)
S2	IPC1 (R18) IPC10 (R18)
S3	IPC50 (R50, R101)
S4	IPC1 (R50, R101) IPC10 (R50, R101)
Soft Label Generation	BSSL
Loss Function	KL-Divergence
Batch Size	16
Epochs	300 (IPC10, IPC50), 1000 (IPC1)
Augmentation	RandomResizedCrop, Horizontal Flip, CutMix

Table 13: Hyperparameters for post-evaluation task on ResNet18, ResNet50 and ResNet101 for Tiny-ImageNet.

<b>Hyperparameters for Model Pre-training</b>	
Optimizer	SGD
Learning Rate	0.01
Weight Decay	1e-4
Momentum	0.9
Batch Size	128
Epoch	300
Scheduler	Cosine Annealing
Augmentation	RandomReizeCrop, Horizontal Flip
Loss Function	Cross-Entropy

Table 14: Hyperparameters for ImageNette Pre-trained Models.

<b>Hyperparameters for Post-Eval on R18, R50 and R101</b>	
Optimizer	Adamw
S2	IPC50 (R101)
S3	IPC10 (R18, R50) IPC50(R50)
S4	IPC1(R18, R50, R101) IPC10 (R101) IPC50 (R18)
Soft Label Generation	BSSL
Loss Function	KL-Divergence
Batch Size	16
Epochs	300
Augmentation	RandomResizedCrop, Horizontal Flip, CutMix

Table 15: Hyperparameters for post-evaluation task on ResNet18, ResNet50 and ResNet101 for ImageNette.

Hyperparameters for Model Pre-training	
Optimizer	SGD
Learning Rate	0.1
Weight Decay	1e-4
Momentum	0.9
Batch Size	128
Epoch	50
Scheduler	Cosine Annealing
Augmentation	RandomResizeCrop, Horizontal Flip
Loss Function	Cross-Entropy

Table 16: Hyperparameters for ImageWoof Pre-trained Models.

Hyperparameters for Post-Eval on R18, R50 and R101	
Optimizer	Adamw
S1	IPC1 (R101)
S2	IPC50 (R18)
S3	IPC10 (R18, R50) IPC50 (R50, R101)
S4	IPC1 (R18, R50) IPC10 (R101)
Soft Label Generation	BSSL
Loss Function	KL-Divergence
Batch Size	16
Epochs	300
Augmentation	RandomResizedCrop, Horizontal Flip, CutMix

Table 17: Hyperparameters for post-evaluation task on ResNet18, ResNet50 and ResNet101 for ImageWoof.

Hyperparameters for Post-Eval on R18, R50 and R101	
Optimizer	Adamw
S1	IPC50 (R18, R50)
S2	IPC1 (R18) IPC10 (R18, R50, R101)
S3	IPC50 (R101)
S4	IPC1 (R50, R101)
Soft Label Generation	BSSL
Loss Function	KL-Divergence
Batch Size	16
Epochs	300
Augmentation	RandomResizedCrop, Horizontal Flip, CutMix

Table 18: Hyperparameters for post-evaluation task on ResNet18, ResNet50 and ResNet101 for ImageNet-1k.

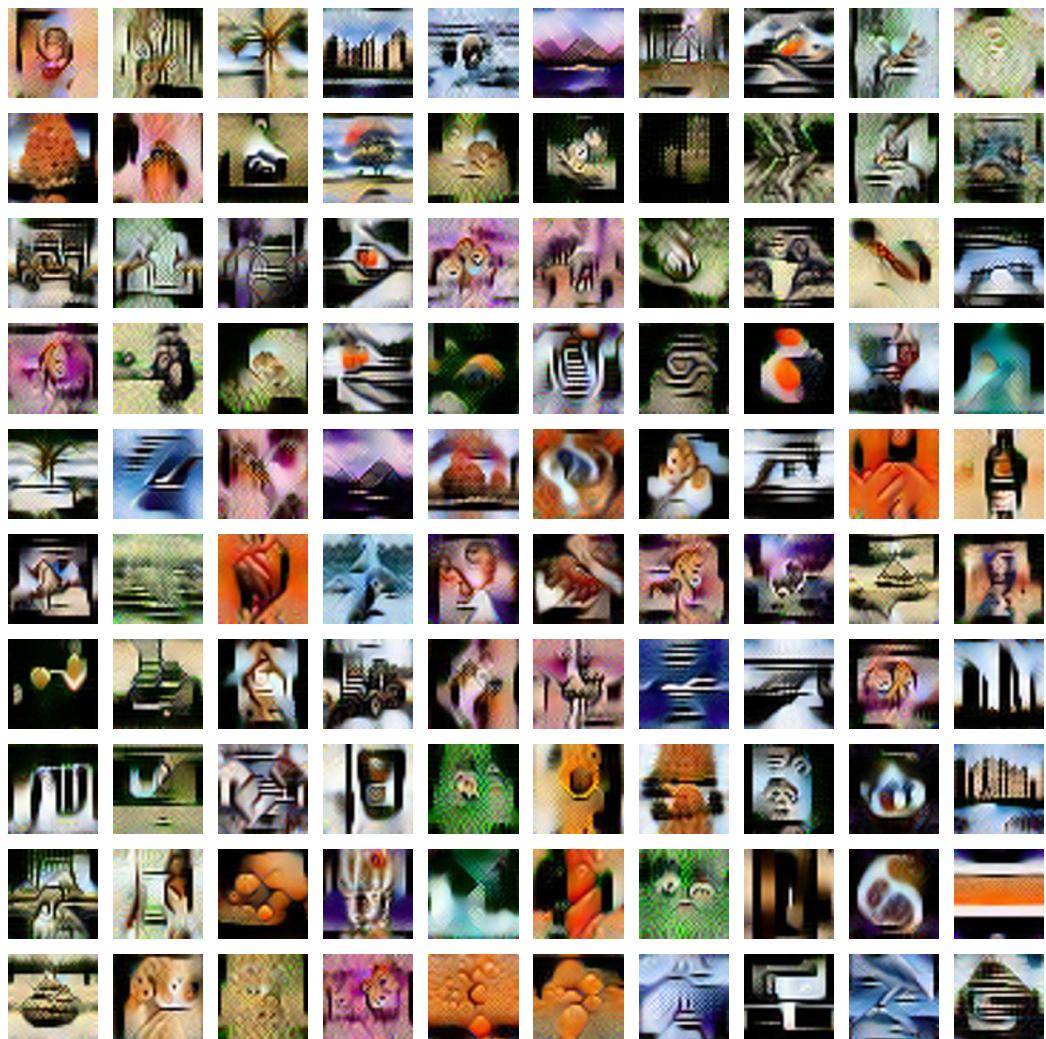


Figure 7: Visualization of synthetic data on CIFAR-100 generated by **FADRM**.

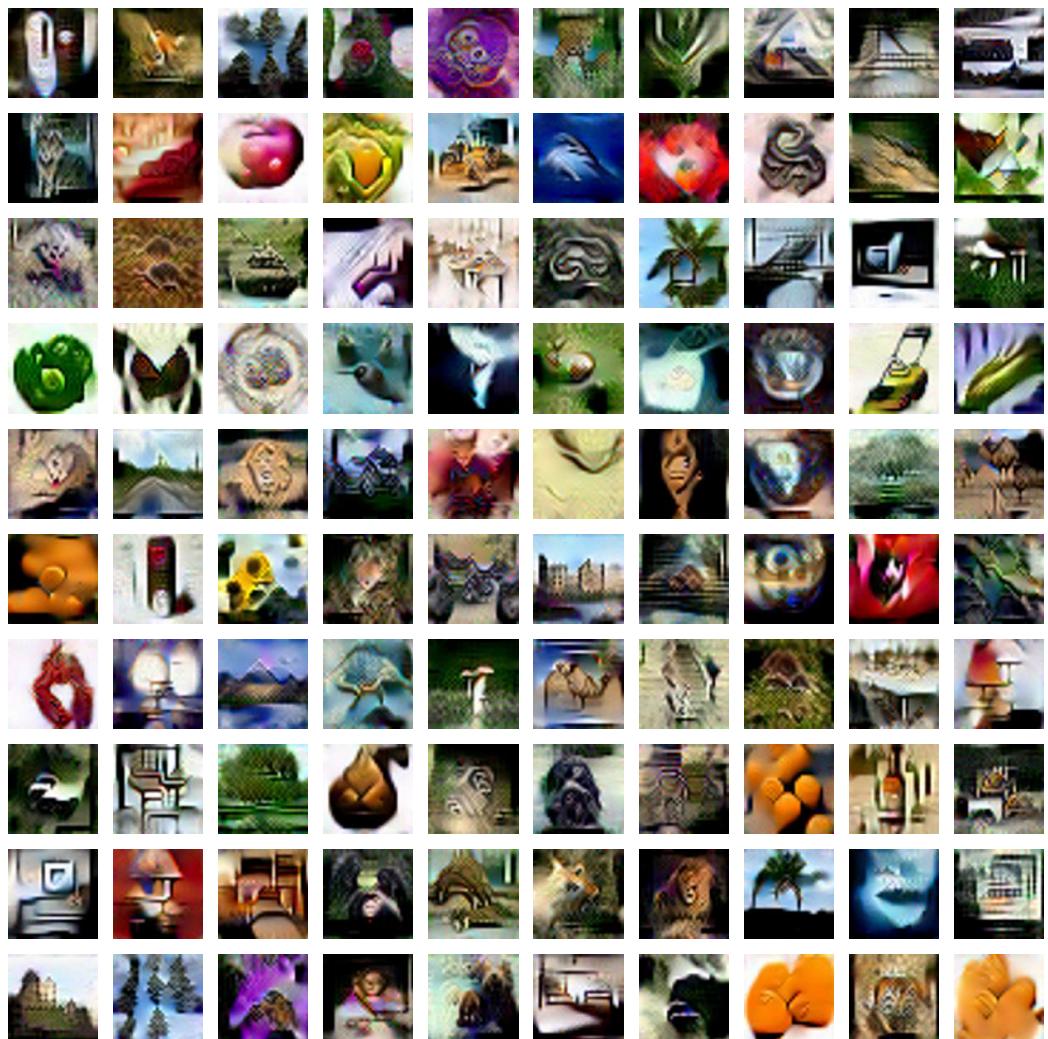


Figure 8: Visualization of synthetic data on CIFAR-100 generated by **FADRM+**.

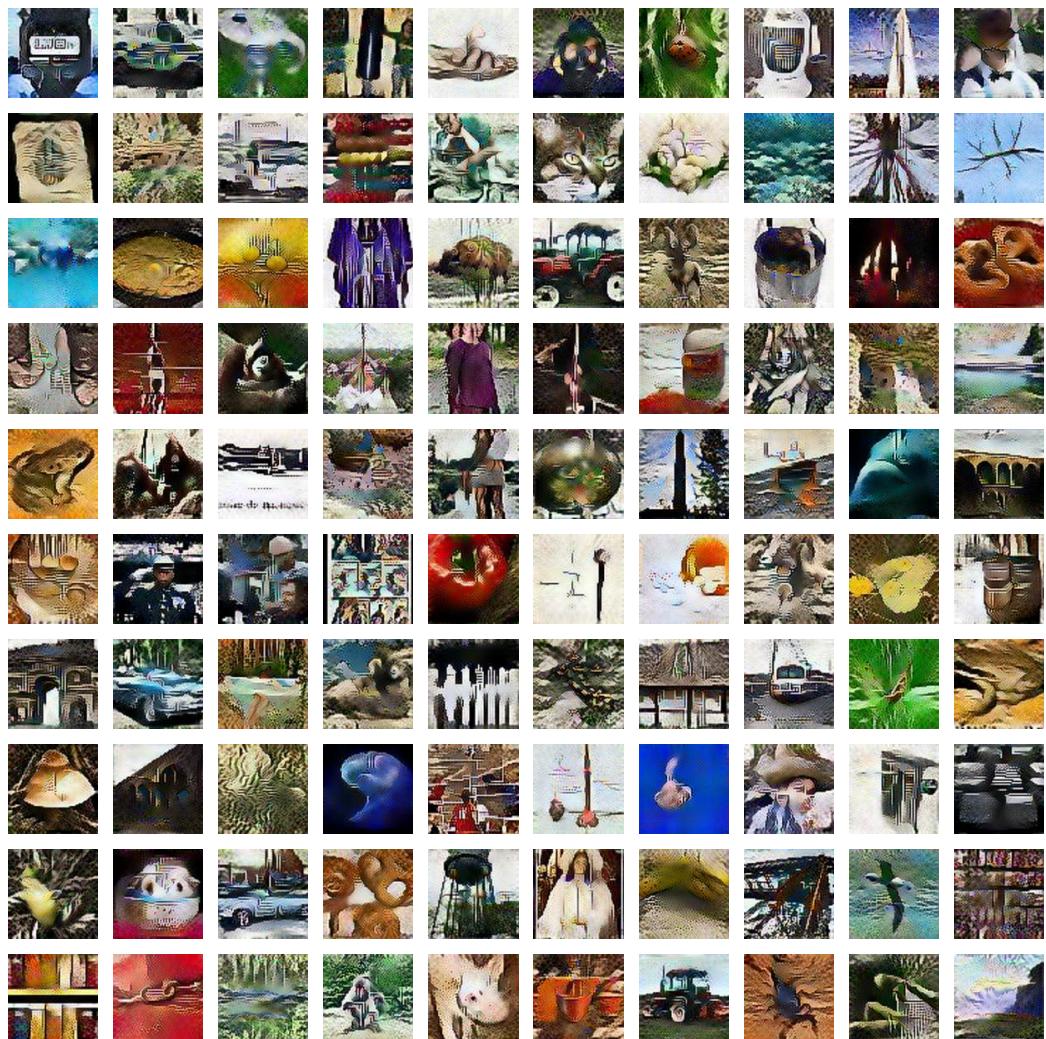


Figure 9: Visualization of synthetic data on Tiny-ImageNet generated by **FADRM**.

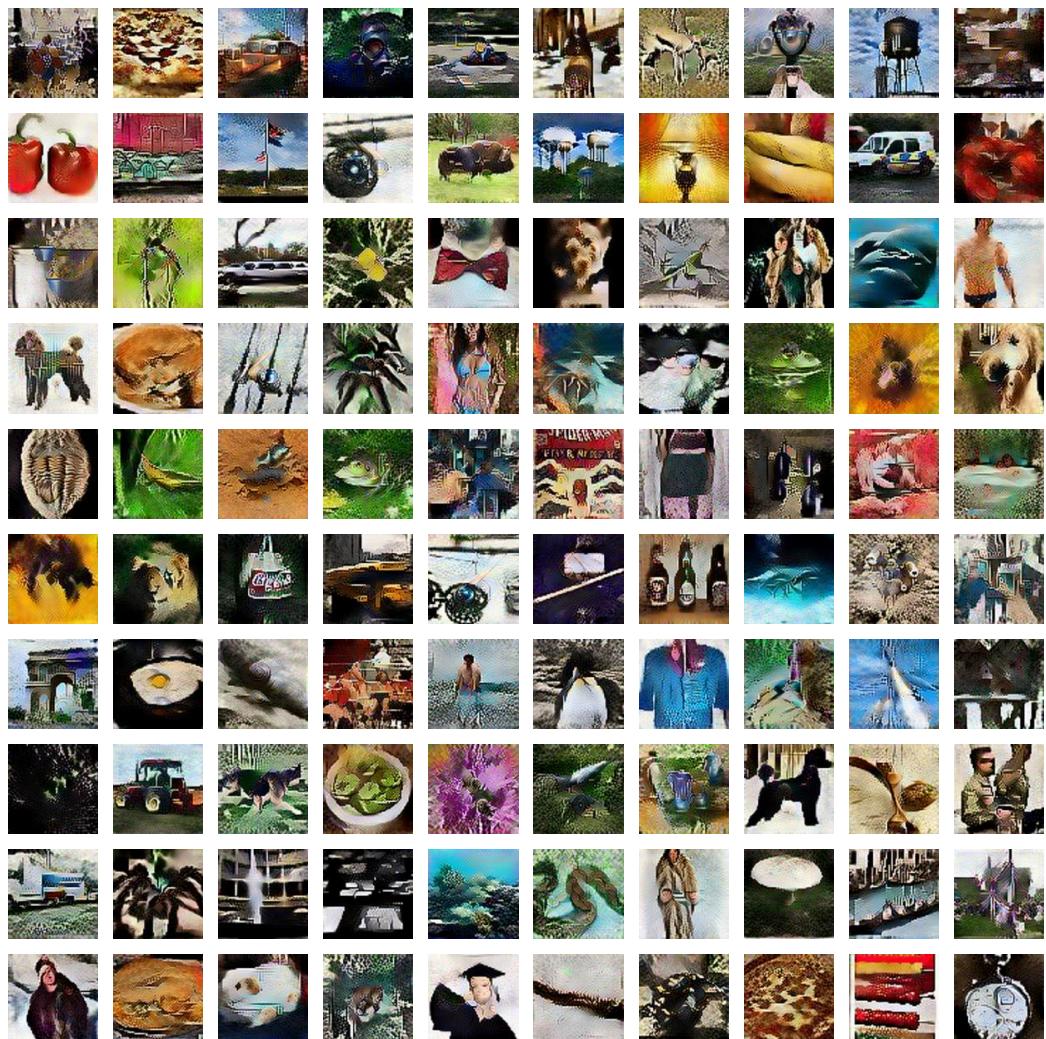


Figure 10: Visualization of synthetic data on Tiny-ImageNet generated by **FADRM+**.

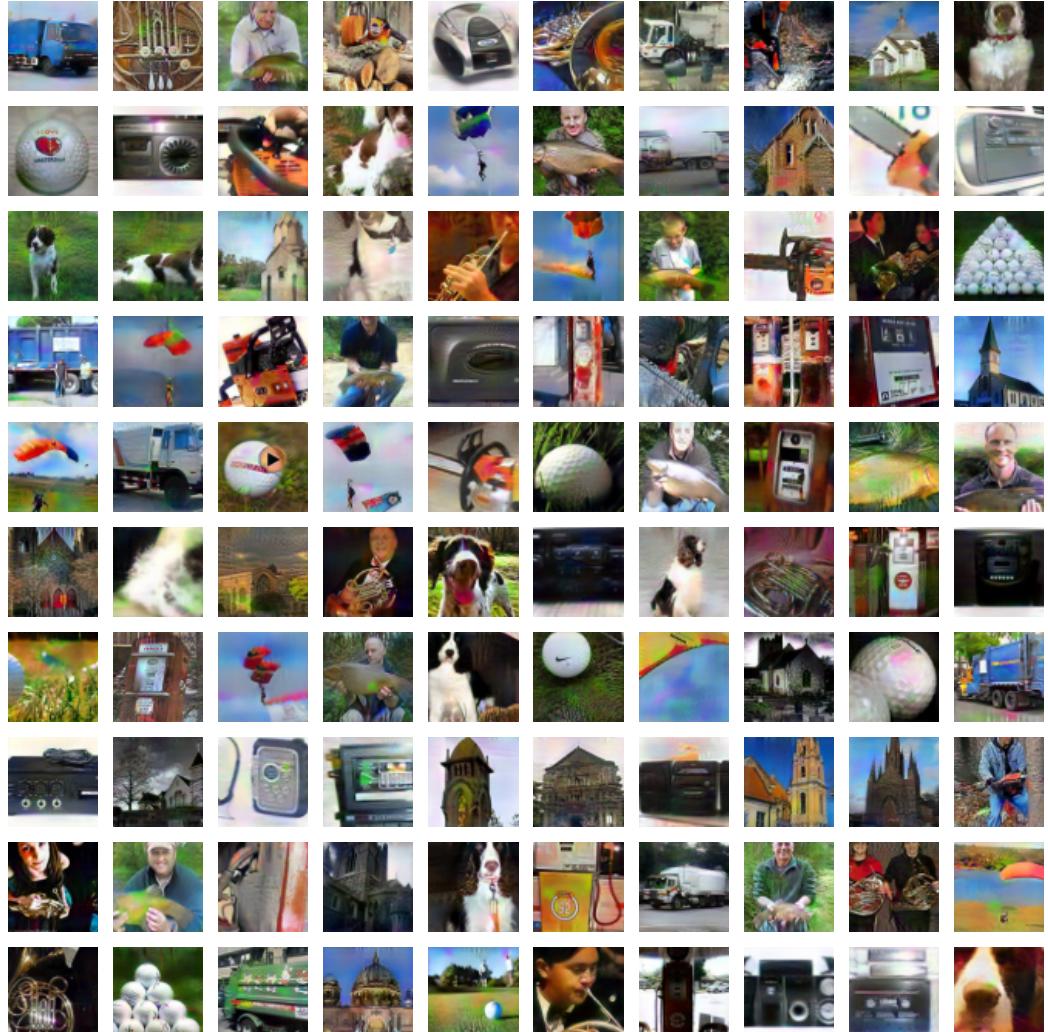


Figure 11: Visualization of synthetic data on ImageNette generated by **FADRM**.

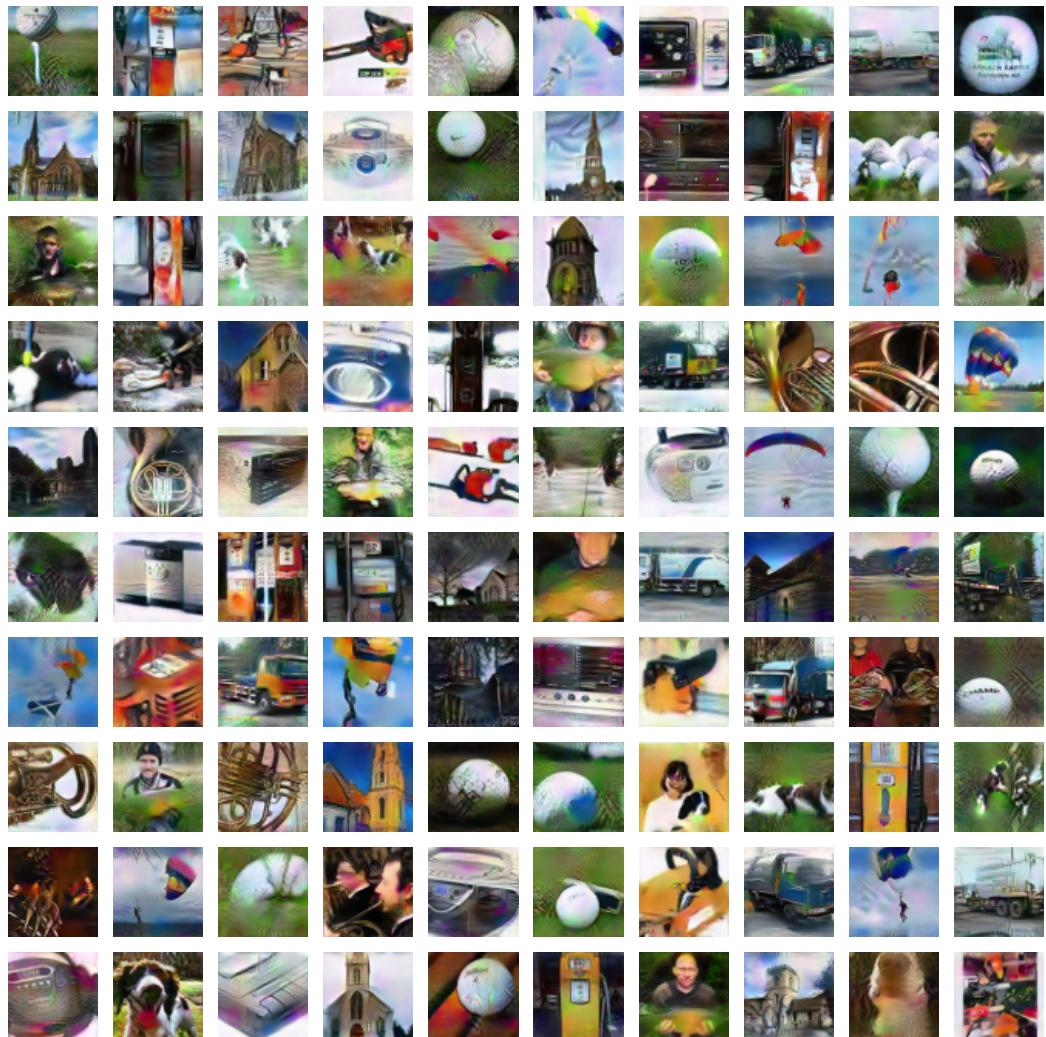


Figure 12: Visualization of synthetic data on ImageNette generated by **FADRM+**.

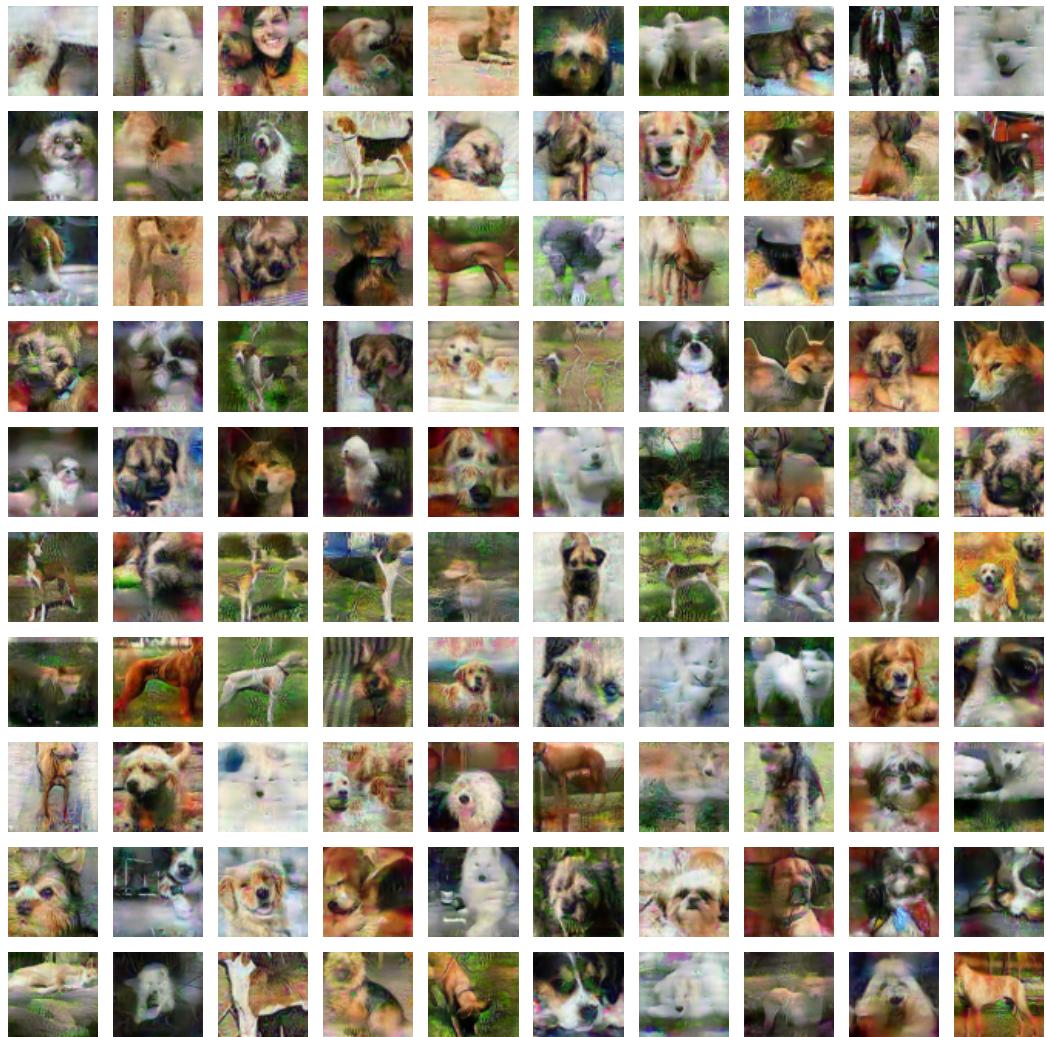


Figure 13: Visualization of synthetic data on ImageWoof generated by **FADRM**.

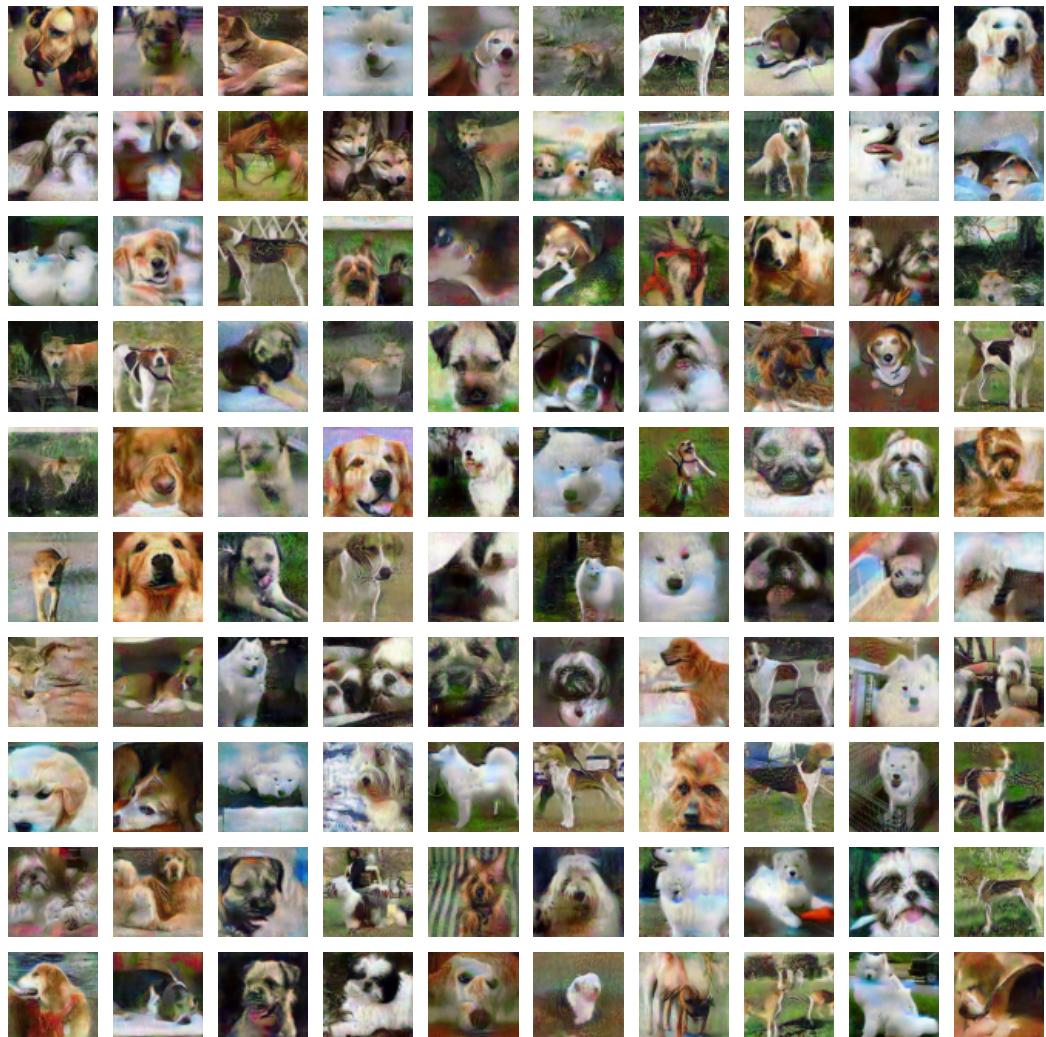


Figure 14: Visualization of synthetic data on ImageWoof generated by **FADRM+**.

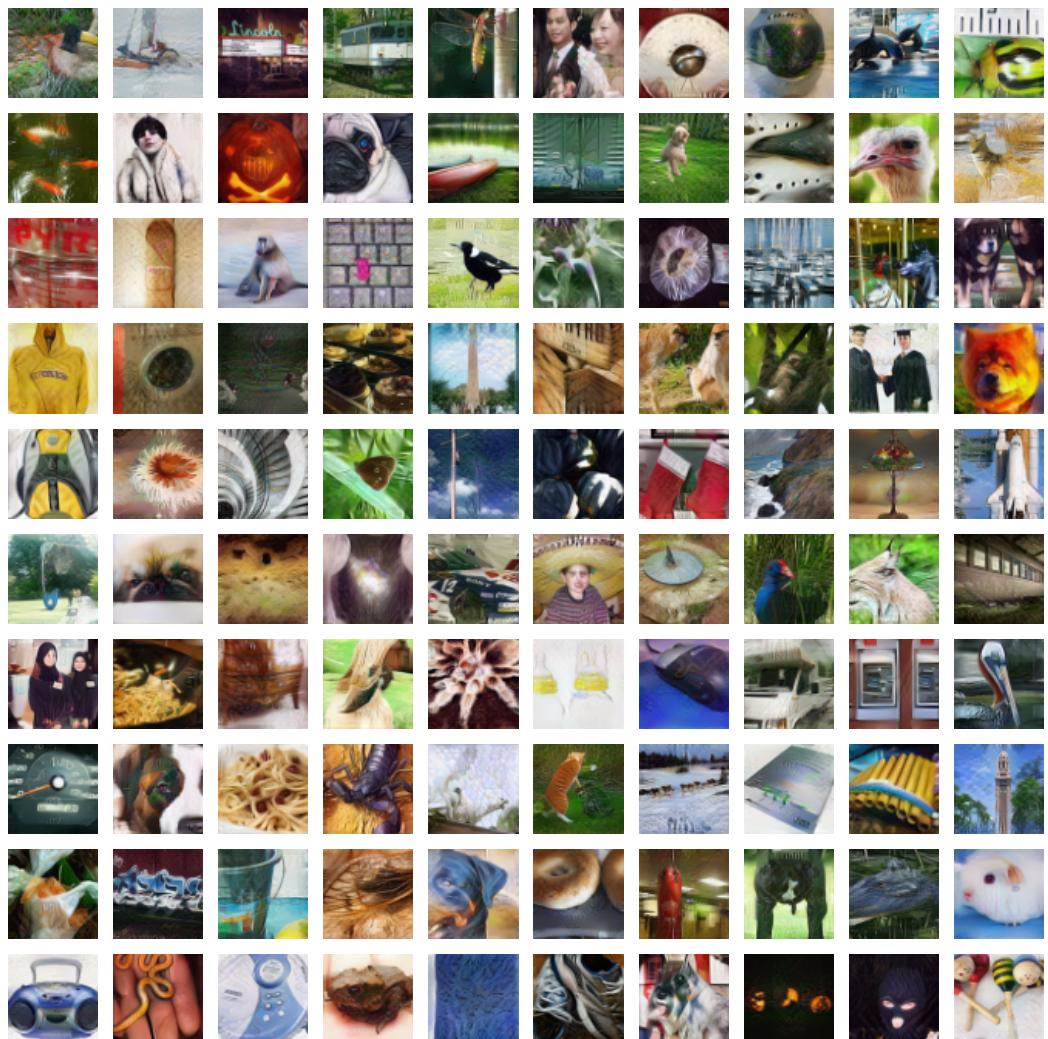


Figure 15: Visualization of synthetic data on ImageNet-1k generated by **FADRM**.

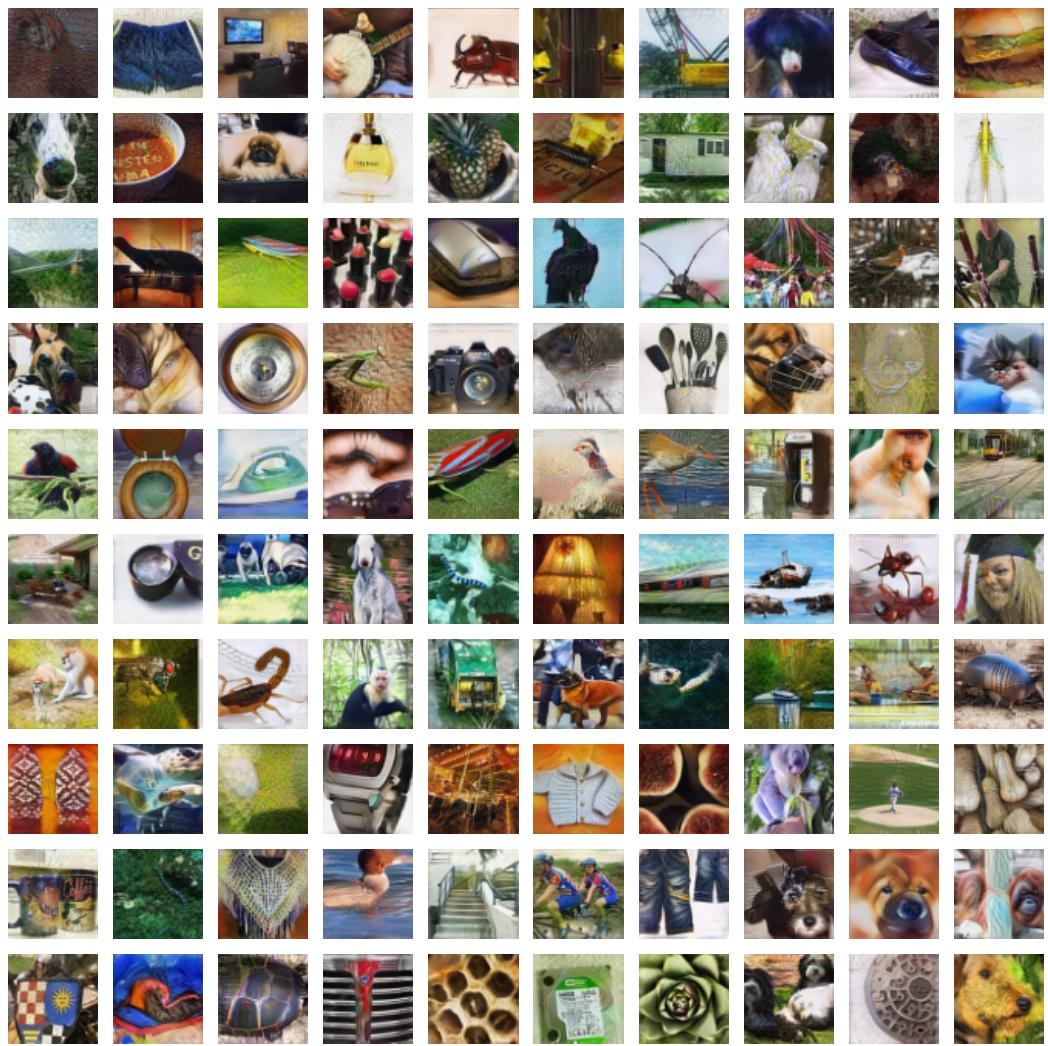


Figure 16: Visualization of synthetic data on ImageNet-1k generated by **FADRM+**.