

# DISTILLING DATASET INTO NEURAL FIELD

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Utilizing large-scale datasets is essential for training high-performance deep learning models, but it also comes with substantial computation and storage costs. To overcome these challenges, dataset distillation has emerged as a promising solution by compressing large-scale datasets into smaller synthetic versions that retain the essential information needed for training. This paper proposes a novel parameterization framework for dataset distillation, coined Distilling Dataset into Neural Field (DDiF), which leverages the neural field to store the necessary information of large-scale datasets. Due to the unique nature of the neural field, which takes coordinates as input and output quantity, DDiF effectively preserves the information and easily generates various shapes of data. Beyond the efficacy, DDiF has larger feature coverage than some previous literature if same budget is allowed, which is proved from the frequency domain perspective. Under the same budget setting, this larger coverage leads to a significant performance improvement in downstream tasks by providing more synthetic instances due to the coding efficiency. DDiF demonstrates both theoretical and empirical evidence of its ability to operate efficiently within a limited budget, while better preserving the information of the original dataset compared to conventional parameterization methods.

## 1 INTRODUCTION

High performances from recent deep learning models are largely driven by scaling laws (Bengio et al., 2007; Kaplan et al., 2020), which heavily rely on large-scale datasets. In spite of this performance gain, utilizing large-scale datasets incurs significant computation and storage costs. *Dataset distillation* has been proposed as a potential solution to address these challenges (Wang et al., 2018). The goal of dataset distillation is to synthesize a small-scale synthetic dataset that contains the essential information for training deep learning models. Naturally, one of the research directions in dataset distillation is defining the essential information and developing efficient methods to learn it. Many studies have been proposed to match some statistics, which is known to be critical in training neural networks, of large datasets with synthesized small ones.

In parallel, another crucial research direction is parameterizing a small-scale synthetic dataset under a limited storage budget. The naive parameterization method generates a synthetic instance in the same structure, i.e. data dimension, as the original instance. Due to issues of scalability and redundancy in this naive approach, various parameterizations enhance the efficiency under a limited storage budget. Specifically, parameterization methods commonly employ low-dimensional codes and decoding functions that transform a code in reduced dimensions into a data instance of the original input space. Conceptually, the decoding functions can be categorized into 1) static decoding (Kim et al., 2022; Shin et al., 2024); 2) parameterized decoding (Deng & Russakovsky, 2022; Sachdeva et al., 2023; Lee et al., 2022; Liu et al., 2022; Wei et al., 2024); and 3) deep generative prior (Cazenavette et al., 2023; Su et al., 2024a;b; Zhong et al., 2024). Although these methods have shown promising results, they have limitations in terms of expressiveness and robustness at varying resolutions. Additionally, there has been limited exploration of the theoretical foundations underlying their methods.

This paper introduces a new parameterization framework for dataset distillation that stores information into *synthetic neural fields* under a limited storage budget, coined **D**istilling **D**ataset into **N**eural **F**ield (DDiF). A *field* is a function that takes a coordinate as an input and returns a corresponding quantity, and a neural field parameterizes the field using a neural network. DDiF parameterizes a synthetic instance as a synthetic neural field. We noted that the neural field has a structural difference

054 compared to conventional decoding functions in dataset distillation, which map a low-dimensional  
 055 space to an instance-level space. Thanks to the nature of the neural field, DDiF effectively encodes  
 056 information from high-dimensional data, which is a crucial challenge in dataset distillation. More-  
 057 over, DDiF can easily decode data of extensive sizes, even if it has not seen it in the distillation  
 058 process. Beyond the efficiency, we provide the theoretical analysis of parameterization through the  
 059 expressiveness of synthetic instances. Based on this theoretical understanding, we prove that DDiF  
 060 has a larger function space that can represent than some previous literature, which implies higher  
 061 expressiveness. Throughout the various evaluation scenarios, DDiF consistently shows performance  
 062 improvements, generalization, robustness, and adaptiveness on vary modality datasets.

063 In summary, our contributions are as follows:

- 064 • We propose a new parameterization framework for dataset distillation, Distilling Dataset  
 065 into Neural Field, which employs the neural field to parameterized the synthetic instance.
- 066 • We provide theoretical analyses on DDIF by comparing the previous methods from the  
 067 perspective of feasible spaces covered by decoded synthetic instances.
- 068 • We empirically investigate the effect of DDiF in its performance gain, cross-architecture  
 069 generalization, and matching statistics generalization. Particularly, we present a new ex-  
 070 perimental design of generalization in *cross resolution* demonstrating the resolution adap-  
 071 tiveness of DDiF. The experimental datasets include images, video, audio, and 3D voxel.

## 074 2 PRELIMINARY

### 076 2.1 PROBLEM FORMULATION

078 This paper focuses on dataset distillation for classification tasks. We define a given large dataset  
 079 that needs to be distilled as  $\mathcal{T} := (X_{\mathcal{T}}, Y_{\mathcal{T}}) = \{(x_i, y_i)\}_{i=1}^{|\mathcal{T}|}$ , where  $X_{\mathcal{T}} := \{x_i\}_{i=1}^{|\mathcal{T}|}$  denotes a set  
 080 of  $D$ -dimensional input data  $x_i \in \mathcal{X} \subseteq \mathbb{R}^D$ , and  $Y_{\mathcal{T}} := \{y_i\}_{i=1}^{|\mathcal{T}|}$  denotes a set of corresponding  
 081 labels among  $C$ -classes  $y_i \in \mathcal{Y} = \{1, \dots, C\}$ . Let a classifier  $f_{\theta} : \mathcal{X} \rightarrow \mathcal{Y}$  be a neural network  
 082 parameterized by  $\theta \in \Theta$ . We also define a loss function  $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ .

084 The main goal of dataset distillation is to synthesize a small dataset such that a model trained on  
 085 this synthetic dataset can generalize well to a large dataset. Formally, given a synthetic dataset  
 086  $\mathcal{S} := (X_{\mathcal{S}}, Y_{\mathcal{S}}) = \{(\tilde{x}_j, \tilde{y}_j)\}_{j=1}^{|\mathcal{S}|}$  where  $X_{\mathcal{S}} := \{\tilde{x}_j\}_{j=1}^{|\mathcal{S}|}$ ,  $Y_{\mathcal{S}} := \{\tilde{y}_j\}_{j=1}^{|\mathcal{S}|}$ , and  $|\mathcal{S}| \ll |\mathcal{T}|$ , the  
 087 objective of dataset distillation is formulated as follows:

$$088 \min_{\mathcal{S}} \mathbb{E}_{(x,y) \sim p(x,y)} [\ell(f_{\theta_{\mathcal{S}}}(x), y)] \quad \text{where } \theta_{\mathcal{S}} = \arg \min_{\theta} \frac{1}{|\mathcal{S}|} \sum_{(\tilde{x}_j, \tilde{y}_j) \in \mathcal{S}} \ell(f_{\theta}(\tilde{x}_j), \tilde{y}_j) \quad (1)$$

091 Nonetheless, the optimization of Eq. (1) is both computationally intensive and lacks scalability, as it  
 092 entails a bi-level optimization problem for both  $\theta$  and  $\mathcal{S}$  (Zhao et al., 2020; Borsos et al., 2020). To  
 093 overcome these issues, several studies have suggested the surrogate objectives to effectively capture  
 094 essential information needed for training the neural network on  $\mathcal{T}$ , such as matching gradient (Zhao  
 095 et al., 2020), distribution (Zhao & Bilen, 2023), and trajectory (Cazenavette et al., 2022). In this  
 096 paper, we denote these objectives as  $\mathcal{L}(\mathcal{T}, \mathcal{S})$ .

### 098 2.2 PARAMETERIZATION OF DATASET DISTILLATION

099 The basic parameterization of a synthetic instance is to configure it in the same format as a real  
 100 instance. For example, if a real image instance  $x$  has a dimension of  $C \times H \times W$ , then the param-  
 101 eterized synthetic instance  $\tilde{x}$  also becomes a tensor of  $C \times H \times W$ . This input-sized parameterization  
 102 suffers from scalability as the dimension of a given instance increases. Also, input-sized parame-  
 103 terization does not utilize the storage budget efficiently because it contains redundant or irrelevant  
 104 information (Lei & Tao, 2023; Yu et al., 2023; Sachdeva & McAuley, 2023).

106 Addressing these concerns, several studies have proposed new parameterization methods to enhance  
 107 the efficiency and representation ability of synthetic dataset  $\mathcal{S}$ . In general, the parameterization  
 method utilizes a pair of compressed codes  $Z$  and decompressing function  $g_{\phi}$  to construct a synthetic

dataset: 1)  $Z := \{z_j\}_{j=1}^{|Z|}$  where  $z_j \in \mathbb{R}^{D'}$  and 2) decoding function  $g_\phi : \mathbb{R}^{D'} \rightarrow \mathbb{R}^D$ .<sup>1</sup> Under this framework, a decoded synthetic instance is represented by a combination of code and decoding function i.e.  $\tilde{x}_j = g_\phi(z_j)$ . Also, a set of decoded synthetic instances become  $X_S = \{g_\phi(z) | z \in Z\}$ .

Based on the form of the decoding function, they can be broadly categorized into 1) static decoding, 2) parameterized decoding, and 3) deep generative prior. Static decoding employs a non-parameterized decoding function  $g$ , such as resizing (Kim et al., 2022) and frequency transform (Shin et al., 2024). These methods are fast, easy to apply, and do not require a budget for the decoding function. However, since this decoding function is fixed, the structure of code  $z$  becomes limited without the ability to adaptively transform  $g$ . Also, using a static decoding function inevitably reduces expressiveness from a generalization perspective, leading to information loss.

Parameterized decoding utilizes a **linear combination** or a lightweight trainable neural network as a decoding function  $g_\phi$ . They employ a **linear combination with learnable coefficients** (Deng & Ruskovskiy, 2022; Sachdeva et al., 2023), decoder (Lee et al., 2022; Zheng et al., 2023), autoencoder (Liu et al., 2022; Duan et al., 2023), or transformer structure (Wei et al., 2024). Although flexible decoding functions are used, the parameters of the learned decoding function must also be stored within a limited budget, necessitating the use of a **simple structure (i.e. linear combination)** or a neural network with a small number of parameters. For this reason, its effectiveness has not been demonstrated for complex data types (e.g., 3D, medical images).

Deep generative prior leverages a pretrained deep generative model without additional training, focusing only on learning the latent vector (Cazenavette et al., 2023; Su et al., 2024a;b; Zhong et al., 2024). This framework encourages better generalization to unseen architecture and scale to high-dimensional datasets. However, it assumes easy access to a well-trained generative model, which can restrict the range of applications. Also, since the deep generative model contains a large number of parameters, the decoding process and backward update process take a long time.

### 2.3 NEURAL FIELD REPRESENTATION

In physics and mathematics, a *field*  $F$  is a function that takes a coordinate in space and outputs a quantity (Xie et al., 2022). If we apply the concept of field to data modeling, a dataset in grid-based representation can be regarded as a field. For example, an RGB image is a function that maps from pixel locations to pixel intensity. Similarly, a video datatype is a function that additionally takes time as input, and a 3D datatype is a function that outputs occupancy value on 3D coordinate system.

According to the universal approximation theorem (Cybenko, 1989), any field can be parameterized by a neural network, which is referred to as a *neural field*  $F_\psi$ . It implies that grid-based data can be expressed as a neural network. Let  $\mathcal{C} := \{c_i\}_{i \in \mathcal{I}}$  be a set of coordinates of grid-based data and  $\mathcal{Q} := \{q_i\}_{i \in \mathcal{I}}$  be a set of corresponding quantities. To encode a grid-based data point using a neural field  $F_\psi$ , we minimize a distortion measure, such as squared error, over all given coordinates as:

$$\min_{\psi} \sum_{i \in \mathcal{I}} \|F_\psi(c_i) - q_i\|_2^2 \quad (2)$$

Recently, the neural field has been adopted to many applications, such as representation learning (Park et al., 2019; Mildenhall et al., 2021), generative modeling (Skorokhodov et al., 2021; Dupont et al., 2021), medical imaging (Shen et al., 2022; Zang et al., 2021), and robotics (Li et al., 2022).

## 3 METHODOLOGY

This section proposes a new parameterization framework for dataset distillation that stores information of a real dataset in *synthetic neural fields* under a limited storage budget, coined Distilling Dataset into Neural Field (DDiF). The core idea of this paper is to store the distilled information in the synthetic function. Although there are several candidates for the form of synthetic function, we primarily focus on (neural) *field*. Figure 1 illustrates the overview of DDiF. In the following, we begin by explaining the reasons for choosing neural field as the form of synthetic function. Then, we introduce our framework, DDiF, which parameterizes a synthetic instance using a neural field. Finally, we provide a theoretical analysis for a better understanding of parameterization and DDiF.

<sup>1</sup>Although some studies (Deng & Ruskovskiy, 2022; Moser et al., 2024) use both  $\tilde{y}$  and  $z$  as inputs for  $g_\phi$ , we expressed it as  $g_\phi(z)$  for the sake of uniformity.

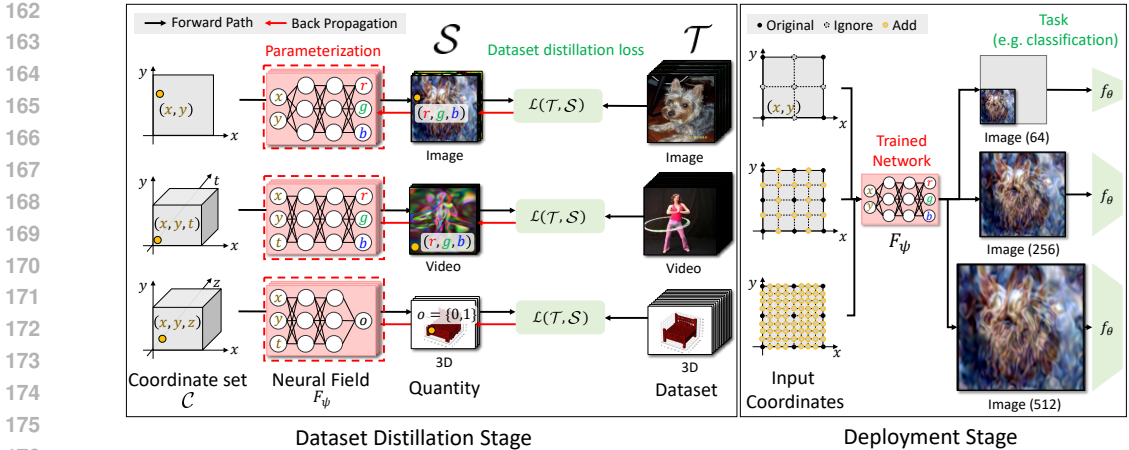


Figure 1: Overview of the proposed framework, DDiF. Following previous dataset distillation studies, the overall process consists of two stages. In the dataset distillation stage (left), each decoded synthetic instance is constructed by the output of each synthetic neural field  $F_\psi$  by inputting coordinate set  $\mathcal{C}$ . DDiF optimizes only the parameters  $\psi$ , as coordinate set  $\mathcal{C}$  does not require optimization or storage. Also, DDiF enables the encoding of various complex datasets. In the deployment stage (right), DDiF easily decodes data with sizes that were not seen during the dataset distillation stage.

### 3.1 WHY NEURAL FIELD IN DATASET DISTILLATION?

Even if the neural field is widely adopted, no research has been conducted on integrating the neural field into dataset distillation. Herein, we provide several properties, which are beneficial to dataset distillation due to its structure.

**Coding Efficiency.** The neural field effectively encodes information from high-dimensional data. Distilling high-dimensional datasets is a crucial challenge in expanding the applicability of dataset distillation (Lei & Tao, 2023). Input-sized parameterization typically scales poorly with resolution due to the *curse of discretization* (Mescheder, 2020). Utilizing a decoding function  $g_\phi : \mathbb{R}^{D'} \rightarrow \mathbb{R}^D$  might improve scalability, but ultimately, the output of the decoding function  $g_\phi$  depends on the data dimension  $D$ . It implies the need for a more complex decoding function as  $D$  increases, which becomes problematic with a limited storage budget. On the other hand, the neural field stores information regardless of data dimension, as it only requires inputting over more coordinates. Furthermore, the output dimension of the neural field is not the same as the data dimension; instead, it relies on the dimension of quantity. Figure 2 illustrates the structural differences between the existing decoding function and neural field.

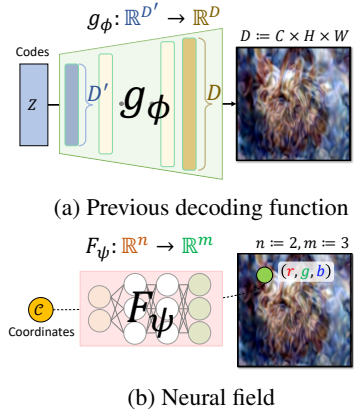


Figure 2: Illustration of structural difference between conventional decoding function and neural field.

**Resolution Adaptiveness.** The neural field is robust to diverse resolution. Consider there is a dataset consisting of data at multiple resolutions. Existing methods in dataset distillation require either creating separate synthetic datasets for each resolution or training at a specific resolution followed by resizing. These approaches need more storage space or result in insufficient information. On the contrary, the neural field enables to store the information from multiple resolutions without any additional storage. Furthermore, in the real world, it is often necessary to resize data depending on various downstream tasks (Wang et al., 2020b; Shorten & Khoshgoftaar, 2019). Existing parameterization methods can only apply postprocessing on optimized synthetic datasets, leading to insufficient information or information distortion. Meanwhile, the neural field easily obtains the various sizes of data by adjusting the coordinate set due to the continuous nature of the neural field. Furthermore, since the neural field is a continuous function, it provides more accurate values for unseen coordinates. Please refer to Section 4.2 for the empirical evidence of this claim.



### 3.2 DDiF: DISTILLING DATASET INTO NEURAL FIELD

As aforementioned, the neural field possesses beneficial features to dataset distillation. We introduce a basic framework for integrating the neural field into dataset distillation. In detail, DDiF parameterizes a synthetic instance  $\tilde{x}_j$  as a neural field  $F_{\psi_j}$ . DDiF consists of two main components: 1) *Coordinate set  $\mathcal{C}$*  and 2) *Synthetic neural fields  $F_{\Psi}$* .

**Coordinate Set  $\mathcal{C}$ .** To define the function, it is first necessary to define the input space of the function. Fundamentally, a decoded synthetic instance by the parameterization method is the same shape as the real instance. Therefore, our synthetic function must be defined in the space where the information of the real instances is stored. Suppose that the given real instance  $x \in X_{\mathcal{T}}$  is  $n$ -dimensional grid representation with resolution  $N_k, k = 1, \dots, n$ , and each element contains  $m$  values, i.e.  $m = 3$  of RGB values. Formally, the real instance  $x$  is element in  $\mathbb{R}^{m \times N_1 \times \dots \times N_n}$ . Then, the coordinate set  $\mathcal{C}$ , where the values of  $x$  are stored, are defined by a bounded set of lattice points:

$$\mathcal{C} := \left\{ (i_1, i_2, \dots, i_n) \mid i_k \in \{0, 1, \dots, N_k\}, \forall k = 1, \dots, n \right\}$$

Note that there are several properties of coordinate set  $\mathcal{C}$ , which become advantages in dataset distillation. First, since every real instance  $x \in X_{\mathcal{T}}$  is defined on the same coordinate set  $\mathcal{C}$  (the only difference is the value on each coordinate), we do not need to consider the coordinate individually. Also,  $\mathcal{C}$  is easily obtained if only the shape of the decoded instance is defined, without any additional information. For instance, assume that we want to get  $N \times N$ -shaped data instances. Then,  $\mathcal{C}$  is a set of lattice points in  $[0, N] \times [0, N]$ . Due to these static and bounded characteristics, DDiF does not need to optimize or store the coordinate set  $\mathcal{C}$ . From this property, DDiF is a new framework that leverages flexible decoding function  $g_{\phi}$  without inferring codes  $Z$ , and the traditional compressed codes become the parameters  $\psi$  of the synthetic neural field  $F_{\psi}$  that will be stored for distillation.

**Synthetic Neural Fields  $F_{\Psi}$ .** DDiF utilizes neural field  $F_{\psi} : \mathbb{R}^n \rightarrow \mathbb{R}^m$  to obtain the decoded synthetic instance  $\tilde{x}$  by inputting the coordinate set  $\mathcal{C}$ .<sup>2</sup> Specifically, given a coordinate set  $\mathcal{C}$ , the decoded synthetic instance by DDiF is  $\tilde{x} = [F_{\psi}(c)]_{c \in \mathcal{C}}$ . In short, we denote the decoded synthetic instance as  $F_{\psi}(\mathcal{C})$  in the DDiF framework. In DDiF, since a decoded synthetic instance  $\tilde{x}$  and a synthetic neural field  $F_{\psi}$  have one-to-one correspondence, obtaining  $K$  decoded synthetic instances requires  $K$  synthetic neural fields. We denote the parameter set of synthetic neural fields as  $\Psi := \{\psi_j\}_{j=1}^{|\Psi|}$  and the set of synthetic neural fields as  $F_{\Psi} := \{F_{\psi_j}\}_{j=1}^{|\Psi|}$ . For the structure of a synthetic neural field  $F_{\psi}$ , we follow the tradition of the neural field (Mildenhall et al., 2021; Tancik et al., 2020), which utilizes a simple  $L$ -layer neural network:

$$F_{\psi}(c) = W^{(L)}(h^{(L-1)} \circ \dots \circ h^{(0)})(c) + b^{(L)}, \quad h^{(l)}(c) = \sigma^{(l)}(W^{(l)}c + b^{(l)})$$

where  $W^{(l)} \in \mathbb{R}^{d_l \times d_{l-1}}$ ,  $b^{(l)} \in \mathbb{R}^{d_l}$  are weights and bias at layer  $l$ .  $\sigma^{(l)}$  denotes a nonlinear activation function. Under this formulation,  $\psi$  becomes  $\{W^{(l)}, b^{(l)}\}_{l=0}^L$ .

The nonlinear activation function  $\sigma^{(l)}$  plays an important role in the expressiveness of coordinate-based neural networks. Although the most widely used activation function is ReLU  $\sigma(x) = \max(0, x)$ , it limits the expressiveness due to the spectral bias (Rahaman et al., 2019; Xu et al., 2019) which indicates that ReLU-based MLP learns high frequencies very slowly. Several approaches have been proposed to mitigate the spectral bias, such as random Fourier features (Tancik et al., 2020) and sine activation functions (Sitzmann et al., 2020). We provide an ablation study to compare the effectiveness of these approaches in Section 4.3.

**Learning Framework.** Given a coordinate set  $\mathcal{C}$  and a parameter set of synthetic neural fields  $\Psi$ , a set of decoded synthetic instances is represented by  $X_{\mathcal{S}} = \{F_{\psi_j}(\mathcal{C})\}_{j=1}^{|\Psi|}$ . Under the arbitrary dataset distillation loss  $\mathcal{L}(\mathcal{T}, \mathcal{S})$ , the overall optimization of DDiF is formulated as follows:

$$\min_{\Psi} \mathcal{L}(\mathcal{T}, \mathcal{S}) \quad \text{where } \mathcal{S} = \{(F_{\psi_j}(\mathcal{C}), \tilde{y}_j)\}_{j=1}^{|\Psi|} \quad (3)$$

DDiF has no limitations in applying a learnable soft label (Sucholutsky & Schonlau, 2021; Cui et al., 2023), but we utilize predefined one-hot labels to ensure a fair comparison with previous parameterization i.e.  $\tilde{y}_j = y_j \in Y_{\mathcal{T}}$ . In practice, dataset distillation commonly utilizes randomly sampled

<sup>2</sup>We defined  $\mathcal{C}$  for dataset distillation training, but the domain of  $F_{\psi}$  is the entire  $n$ -dimensional space  $\mathbb{R}^n$ .

real instance  $x \in X_{\mathcal{T}}$  as the initialization of synthetic instance  $\tilde{x}$ . In the same context, we conduct the warm-up training for synthetic neural fields  $F_{\Psi}$ . Concretely, we train  $F_{\Psi}$  using Eq. (2) with randomly selected  $|\Psi|/C$  samples for each class. Algorithm 1 specifies a training procedure for our method DDiF.

**Budget calculation.** As Eq. (3) shows, the optimization target of DDiF is  $\Psi = \{\psi_j\}_{j=1}^{|\Psi|}$ . Note that each synthetic neural field  $F_{\psi}$  utilize  $d_0(n+1) + \sum_{l=1}^{L-1} d_l(d_{l-1}+1) + m(d_{L-1}+1) =: b$  parameters. We highlight that  $b$  does not depend on the size of the real instance  $D$ , so the high resolution would not necessarily increase the budget in  $F_{\psi}$ . Given a budget is  $B$ , we set a structure of  $F_{\psi}$ , such as width  $d_l$  and number of layers  $L$ , to satisfy  $|\Psi| \times b \leq B$ .

---

**Algorithm 1:** Training procedure of DDiF

---

**Input:** Original real dataset  $\mathcal{T}$ ; Dataset distillation loss  $\mathcal{L}$ ; Initialized  $\Psi = \{\psi_j\}_{j=1}^{|\Psi|}$ ; Learning rate  $\eta$

**Output:** Parameterized synthetic dataset  $\{(\psi_j, \tilde{y}_j)\}_{j=1}^{|\Psi|}$

---

```

1 Initialize coordinate set  $\mathcal{C}$  from  $x \in \mathcal{T}$ 
2 for  $j = 1$  to  $|\Psi|$  do
3   Sample a real instance  $(x, y) \sim \mathcal{T}$ 
4    $\tilde{y}_j \leftarrow y$ 
5   Optimize  $\psi_j$  with Eq. (2)
6 repeat
7   Sample a real mini-batch  $\mathcal{B}_{\mathcal{T}} \sim \mathcal{T}$ 
8    $\mathcal{B}_{\mathcal{S}} \leftarrow \{F_{\psi}(\mathcal{C})|_{\psi \in \Psi_{\mathcal{B}}}\}$  from  $\Psi_{\mathcal{B}} \sim \Psi$ 
9    $\Psi \leftarrow \Psi - \eta \nabla_{\Psi} \mathcal{L}(\mathcal{B}_{\mathcal{T}}, \mathcal{B}_{\mathcal{S}})$ 
10 until convergence;

```

---

### 3.3 THEORETICAL ANALYSIS

Even if several parameterization methods for dataset distillation are proposed, there has been little discussion regarding the theoretical understanding of their methods. Herein, we provide a simple but intuitive theoretical analysis of parameterization through the expressiveness of synthetic instances. Then, we investigate the expressiveness of DDiF when using the sine activation function. Lastly, we compare DDiF with a previous work, FreD (Shin et al., 2024).

Let’s say that the expressiveness of a synthetic instance is equal to the coverage of its corresponding data space. Then, the optimization of the synthetic instance limits the coverage by its optimization feasibility. Given this conceptual assumption, we observe the relationship between the feasible space of synthetic instances and the optimal value of the dataset distillation objective. Herein, we assume only synthetic inputs  $X_{\mathcal{S}}$  as the optimization variable. Consequently, dataset distillation loss  $\mathcal{L}(\mathcal{T}, \mathcal{S})$  is simply expressed as a function of  $X_{\mathcal{S}}$  i.e.  $\mathcal{L}(\mathcal{T}, \mathcal{S})$  is represented by  $\mathcal{L}(X_{\mathcal{S}})$ .

**Observation 1.** Let two matrix variables  $X_1 := [x_{11}, \dots, x_{1M}]$  and  $X_2 := [x_{21}, \dots, x_{2M}]$  consisting of columns  $x_{ij} \in \mathcal{X}_i \subseteq \mathbb{R}^d$  for  $i = 1, 2$  and  $j = 1, \dots, M$ . If  $\mathcal{X}_1 \subseteq \mathcal{X}_2$ , then  $\min_{X_1 \in \mathcal{X}_1^M} \mathcal{L}(X_1) \geq \min_{X_2 \in \mathcal{X}_2^M} \mathcal{L}(X_2)$ .

Please refer to Appendix A.1 for proof. Observation 1 claims that if a particular synthetic instance’s feasible space contains another synthetic instance’s feasible space, the former’s optimal value is lower under the same number of synthetic instances. This result is intuitive when considering situations where either optimizable dimension or value is constrained. Please refer to Appendix C.5 for the empirical evidence of Observation 1.

As an extension of Observation 1, which considers input-sized parameterization, Observation 2 investigates the relationship between the feasible space of decoded synthetic instances and the optimal value of the dataset distillation objective:

**Observation 2.** Consider two functions  $g_1, g_2$  where  $g_i : \mathcal{Z}_i \rightarrow \mathbb{R}^d$  for  $i = 1, 2$ . Also, consider two matrix variables  $Z_i := [z_{i1}, \dots, z_{iM}]$  where their columns  $z_{ij} \in \mathcal{Z}_i$  for  $i = 1, 2$  and  $j = 1, \dots, M$ . We denotes  $\hat{g}_i(Z_i) := [g_i(z_{i1}), \dots, g_i(z_{iM})]$  for  $i = 1, 2$ . Set  $\mathcal{G}_i := \{g|g : \mathcal{Z}_i \rightarrow \mathbb{R}^d\}$  for  $i = 1, 2$ . If  $g_1(\mathcal{Z}_1) \subseteq g_2(\mathcal{Z}_2)$  for any  $g_1 \in \mathcal{G}_1$  and  $g_2 \in \mathcal{G}_2$ , then  $\min_{g_1 \in \mathcal{G}_1, Z_1 \in \mathcal{Z}_1^M} \mathcal{L}(\hat{g}_1(Z_1)) \geq \min_{g_2 \in \mathcal{G}_2, Z_2 \in \mathcal{Z}_2^M} \mathcal{L}(\hat{g}_2(Z_2))$ .

Please refer to Appendix A.2 for proof. Observation 2 claims a similar statement: the optimal value of dataset distillation loss becomes smaller as the feasible space of decoded synthetic instances becomes larger. The larger feasible space enables a more complex combination, or *expressiveness*, of codes  $Z$  and decoding function  $g$ , and the combination enhances the performance. The experimental finding in Shin et al. (2024) that an increase in frequency dimension leads to improved dataset distillation performance supports Observation 2.

Now, we focus on the expressiveness of DDiF. The following theorem states the representation ability of a two-layer sinusoidal neural network which is a special case of DDiF:

**Theorem 1.** (Novello, 2022) Consider a neural network  $F_\psi : \mathbb{R} \rightarrow \mathbb{R}$  with two hidden layers and width  $d$ . If  $F_\psi$  utilizes a sine activation function, then  $F_\psi$  represents a function which is the sum of sines and cosines:

$$F_\psi(x) = b^{(2)} + \sum_{k \in \mathbb{Z}^d} \alpha_k \cos(\omega_k x) + \beta_k \sin(\omega_k x) \quad (4)$$

where  $\omega_k := \langle k, W^{(0)} \rangle$ ,  $\varphi_{k,i} := \langle k, b^{(0)} \rangle + b_i^{(1)}$ ,  $\alpha_k := \sum_{i=1}^d A_{k,i} \sin(\varphi_{k,i})$  and  $\beta_k := \sum_{i=1}^d A_{k,i} \cos(\varphi_{k,i})$ . Also,  $A_{k,i} := W_i^{(2)} \lambda_k(W_i^{(1)})$  and  $\lambda_k(W_i^{(1)}) := \prod_{j=1}^d J_{k_j}(W_{ij}^{(1)})$  where  $J_{k_j}$  denotes Bessel function of the first kind of order  $k_j$ .

**Remark 1.** By using trigonometric identity, Eq. (4) is represented by the sum of cosines:

$$F_\psi(x) = b^{(2)} + \sum_{k \in \mathbb{Z}^d} \sum_{i=1}^d A_{k,i} \cos(\omega_k x + \varphi'_{k,i}) \quad \text{where } \varphi'_{k,i} = \varphi_{k,i} - \frac{\pi}{2} \quad (5)$$

According to Theorem 1 and Remark 1, DDiF enables to change the amplitudes  $A_{k,i}$ , frequencies  $\omega_k$ , phases  $\varphi'_{k,i}$ , and shift  $b^{(0)}$ . Eq. (5) has a similar form of expressiveness of previous work, FreD (Shin et al., 2024). FreD optimizes frequency coefficients, which are selected by the explained variance ratio. They utilize inverse discrete cosine transform (IDCT) to decode synthetic instances from the frequency domain. Suppose that FreD utilizes IDCT with  $N$  equidistant locations on  $\mathbb{R}$ . Also, when  $\mathcal{U} \subset \mathcal{C}_N := \{0, \dots, N-1\}$  is the index set of selected frequency dimension; the optimized frequency coefficients is denoted as  $\Gamma := \{\gamma_u | u \in \mathcal{U}\}$ . Then, the value of decoded synthetic instance by FreD is expressed in the form of a function over  $\mathcal{C}_N$ :

$$g(x; \Gamma) = \sum_{u \in \mathcal{U}} \gamma_u \cos\left(\frac{\pi u}{N} x + \frac{\pi u}{2N}\right) \quad \text{where } x \in \mathcal{C}_N \quad (6)$$

Thanks to the similarity, we provide the relationship between the feasible space of DDiF and FreD:

**Theorem 2.** For any  $x \in \mathcal{C}_N$ ,  $g(x; \Gamma) \subsetneq F_\psi(x)$ .

Please refer to Appendix A.4 for proof. Theorem 2 implies that the range of decoded values of DDiF is wider than FreD. From a dataset distillation perspective, if  $|\mathcal{U}| \geq 6$ , DDiF enables to construct a valid neural network i.e.  $d \geq 1$  with same budget. Consequently, according to Observation 2, DDiF achieves a smaller optimal value of dataset distillation loss than FreD. Please refer to Appendix D.1 for a detailed comparison between DDiF and FreD based on Eqs. (5) and (6).

## 4 EXPERIMENTS

### 4.1 EXPERIMENT SETTING

To evaluate the efficacy of DDiF, we primarily focus on the experiments on high-resolution image datasets, ImageNet-Subset (Cazenavette et al., 2022; 2023) with 128, 256 resolution, which is a more challenging task in dataset distillation. Please refer to Appendix C for the experiments on low-resolution datasets, such as CIFAR-10 and CIFAR-100 (Krizhevsky et al., 2009). Furthermore, we conduct several experiments to verify the applicability of DDiF on various modalities. We utilize miniUCF for video (Wang et al., 2024; Khurram, 2012), Mini Speech Commands for audio (Kim et al., 2022; Warden, 2018), and ModelNet (Wu et al., 2015), ShapeNet (Chang et al., 2015) for 3D. We enumerate the baselines and details of loss configurations in Appendix B.

### 4.2 EXPERIMENTAL RESULTS

**Main Performance.** Table 1 shows the overall performance for ImageNet-Subset with resolution 128 and 256 under IPC=1. We utilize trajectory matching (TM) for 128 resolution and distribution matching (DM) for 256 resolution. DDiF achieves the best performances in all experimental settings. Remarkably, DDiF shows a significant performance margin compared to the second-best performer: from 3.4%p to 6.1%p in 128 resolution, and from 5.8%p to 11.7%p in 256 resolution. Please check  $\Delta$  row of Table 1. We underline the performance gap from vanilla to DDiF is significant, proving the efficacy of our method as dataset parameterization. In Table 2, DDiF shows consistent improvement and highly competitive performance with baseline. These results demonstrate that the neural field positively impacts efficiency and performance enhancement, particularly when the budget is very limited. Please refer to Appendix C for the experiments on CIFAR-10 and CIFAR-100.

Table 1: Test accuracies (%) on ImageNet-Subset with regard to various dataset parameterizations under IPC=1. "IPC" denotes instances per class, which indicates the total allowed budget. We utilize trajectory matching (TM) for  $128 \times 128$  and distribution matching (DM) for  $256 \times 256$ . For Vanilla (TM) in  $128 \times 128$ , we report our implementation results, which are higher than the officially reported. Bold and Underline mean the best and second-best performance of each column, respectively. "-" indicates no reported results.  $\Delta$  represents the performance gap over Vanilla. The full table with standard deviations is in Appendix C.8.

Resolution		128 × 128						256 × 256					
Subset		Nette	Woof	Fruit	Yellow	Meow	Squawk	Nette	Woof	Fruit	Yellow	Meow	Squawk
Input sized	Vanilla	<b>51.4</b>	<b>29.7</b>	<b>28.8</b>	<b>47.5</b>	<b>33.3</b>	<b>41.0</b>	32.1	20.0	19.5	33.4	21.2	27.6
	FRePo	48.1	29.7	-	-	-	-	-	-	-	-	-	-
Static	IDC	61.4	34.5	38.0	56.5	39.5	50.2	53.7	30.2	<u>33.1</u>	<u>52.2</u>	34.6	47.0
	FreD	66.8	<u>38.3</u>	<u>43.7</u>	<u>63.2</u>	43.2	57.0	54.2	<u>31.2</u>	32.5	49.1	34.0	43.1
Parameterized	HaBa	51.9	32.4	34.7	50.4	36.9	41.9	-	-	-	-	-	-
	SPEED	66.9	38.0	43.4	62.6	43.6	<u>60.9</u>	<u>57.7</u>	-	-	-	-	-
	LatentDD	-	-	-	-	-	-	56.1	28.0	30.7	-	<u>36.3</u>	<u>47.1</u>
	NSD	<u>68.6</u>	35.2	39.8	61.0	<u>45.2</u>	52.9	-	-	-	-	-	-
DGM Prior	GLaD	38.7	23.4	23.1	-	26.0	35.8	-	-	-	-	-	-
	H-GLaD	45.4	28.3	25.6	-	29.6	39.7	-	-	-	-	-	-
Function	DDiF	<b>72.0</b>	<b>42.9</b>	<b>48.2</b>	<b>69.0</b>	<b>47.4</b>	<b>67.0</b>	<b>67.8</b>	<b>39.6</b>	<b>43.2</b>	<b>63.1</b>	<b>44.8</b>	<b>67.0</b>
	$\Delta(\%)$	(40.1)	(44.4)	(67.4)	(45.3)	(42.3)	(63.4)	(111.2)	(96.0)	(121.5)	(88.9)	(107.4)	(142.8)
Entire dataset $\mathcal{T}$		87.4	67.0	63.9	84.4	66.7	87.5	92.5	80.1	70.2	90.5	72.2	93.2

Table 2: Test accuracies (%) on ImageNet-Subset under IPC=10.

Method	Nette	Woof	Fruit	Yellow	Meow	Squawk
TM	63.0	35.8	40.3	60.0	40.4	52.3
FRePo	66.5	42.2	-	-	-	-
IDC	70.8	39.8	46.4	68.7	47.9	65.4
FreD	72.0	41.3	47.0	69.2	48.6	67.3
HaBa	64.7	38.6	42.5	63.0	42.9	56.8
SPEED	<u>72.9</u>	<u>44.1</u>	<b>50.0</b>	<b>70.5</b>	<b>52.0</b>	<u>71.8</u>
DDiF	<b>74.6</b>	<b>44.9</b>	<u>49.8</u>	<b>70.5</b>	<u>50.6</u>	<b>72.3</b>
Entire $\mathcal{T}$	87.4	67.0	63.9	84.4	66.7	87.5

Table 3: Average test accuracies (%) on ImageNet-Subset ( $128 \times 128$ ) across AlexNet, VGG11, ResNet18, and ViT, under IPC=1.

Method	Nette	Woof	Fruit	Yellow	Meow	Squawk
TM	22.0	14.8	17.1	22.3	16.2	25.5
IDC	27.9	19.5	<u>23.9</u>	28.0	19.8	29.9
FreD	<u>36.2</u>	<u>23.7</u>	23.6	<u>31.2</u>	19.1	37.4
GLaD	30.4	17.1	21.1	-	19.6	28.2
H-GLaD	30.8	17.4	21.5	-	20.1	28.8
LD3M	32.0	19.9	21.4	-	<u>22.1</u>	<u>30.4</u>
DDiF	<b>59.3</b>	<b>34.1</b>	<b>39.3</b>	<b>51.1</b>	<b>33.8</b>	<b>54.0</b>

**Cross-architecture Generalization.** The network structure used for dataset distillation might differ from the one used for training with the distilled dataset. Accordingly, the parameterization methods should achieve consistent performance enhancement across various test network architectures. In this study, we utilize AlexNet (Krizhevsky et al., 2012b), VGG11 (Simonyan & Zisserman, 2014), ResNet18 (He et al., 2016), and ViT (Dosovitskiy, 2020) while ConvNetD5 is utilized in training. Table 3 presents that DDiF consistently outperforms. Remarkably, DDiF shows a significant performance gap to the second-best performer from 10.4%p to 23.1%p. These results indicate that DDiF effectively encodes important task-relevant information regardless of the training network.

**Robustness to the Dataset Distillation Loss.** Another important evaluation metric for parameterization is whether it constantly improves the performance across various dataset distillation losses. We utilize gradient matching (DC) and distribution matching (DM) to evaluate the robustness to the dataset distillation loss. In Table 4, DDiF achieves the best performances in both dataset distillation losses. These results confirm the robustness of DDiF in dataset distillation loss, representing its wide adaptiveness.

**Various Modality.** Since dataset distillation studies have mainly developed within the image domain and only a few studies have focused on the specific domain, the applicability of previous methods may be constrained. However, our method, DDiF, consistently achieves the best performance across different domains, suggesting its potential as a robust baseline for use across various modalities. First, Figure 3 shows the test performances on the video domain with regard to the required budget for running each method. Existing

Table 4: Test accuracies (%) on ImageNet-Subset ( $128 \times 128$ ) across DC and DM under IPC=1.

$\mathcal{L}$	Method	Nette	Woof	Fruit	Meow	Squawk
DC	Vanilla	34.2	22.5	21.0	22.0	32.0
	IDC	45.4	25.5	26.8	25.3	34.6
	FreD	<u>49.1</u>	<u>26.1</u>	<u>30.0</u>	<u>28.7</u>	<u>39.7</u>
	GLaD	35.4	22.3	20.7	22.6	33.8
	H-GLaD	36.9	24.0	22.4	24.1	35.3
	DDiF	<b>61.2</b>	<b>35.2</b>	<b>37.8</b>	<b>39.1</b>	<b>54.3</b>
DM	Vanilla	30.4	20.7	20.4	20.1	26.6
	IDC	48.3	27.0	29.9	30.5	38.8
	FreD	<u>56.2</u>	<u>31.0</u>	<u>33.4</u>	<u>33.3</u>	<u>42.7</u>
	GLaD	32.2	21.2	21.8	22.3	27.6
	H-GLaD	34.8	23.9	24.4	24.2	29.5
	DDiF	<b>69.2</b>	<b>42.0</b>	<b>45.3</b>	<b>45.8</b>	<b>64.6</b>

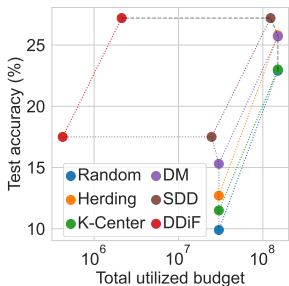


Figure 3: Test accuracies (%) on Video domain. Each black line denotes the same number of decoded instances per class, 1 and 5, respectively.

Table 5: Test accuracies (%) on Audio domain.

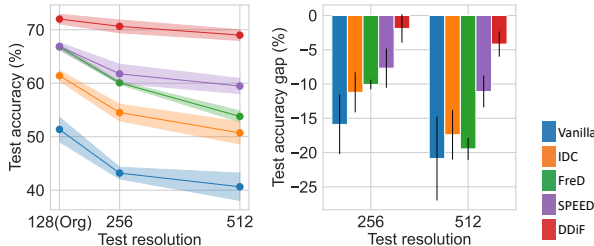
Spec. / class	10	20
Random	42.6	57.0
Herding	56.2	72.9
DSA	65.0	74.0
DM	69.1	77.2
<hr/>		
IDC-I	73.3	83.0
IDC-I + HaBa	74.5	84.3
IDC	<u>82.9</u>	<u>86.6</u>
DDiF	<b>90.5</b>	<b>92.7</b>
<hr/>		
Entire	93.4	

Table 6: Test accuracies on 3D domain under IPC=1.

$\mathcal{L}$	Method	ModelNet	ShapeNet
DC	Random	60.9	68.5
	Vanilla	56.3	48.6
	IDC	78.7	79.9
	FreD	<u>85.6</u>	<u>88.2</u>
	DDiF	<b>87.1</b>	<b>89.6</b>
DM	Vanilla	75.3	80.4
	IDC	85.6	85.3
	FreD	<u>87.3</u>	<u>90.6</u>
	DDiF	<b>88.4</b>	<b>93.1</b>
Entire		91.6	98.3

dataset distillation methods, DM and SDD (Wang et al., 2024) achieve higher performance than coreset selections, but they still require a large storage budget size. DDiF achieves competitive performance even with a budget equivalent to only 1.7% of SDD. Second, Table 5 and 6 shows the test performances of audio and 3D domain, respectively. Both tables indicate that DDiF achieves the highest performance. These results demonstrate the efficacy of DDiF in various modality datasets.

**Cross-resolution Generalization.** As mentioned in Section 3.1, previous studies can only perform post-resizing on optimized synthetic datasets, resulting in information distortion. In contrast, DDiF can easily generate data of various sizes by adjusting the coordinate set, thanks to the continuous nature of the neural field. To verify this property, we introduce an experiment when the size of data changes from the dataset distillation stage to the deployment stage, which is the first experimental design in the dataset distillation community. We define this setting as a *cross-resolution generalization*. We apply the interpolation techniques, such as nearest, bilinear, and bicubic, for the optimized synthetic data of existing parameterization.



(a) Accuracy (b) Accuracy degradation

Figure 4: (a) Test accuracies (%) with different image resolutions. The original resolution is 128x128. (b) Test accuracy gap (%) from original. We use bilinear interpolation for previous studies.

Figure 4a shows test accuracies on each test resolution when utilizing the corresponding network architecture, ConvNetD6 for 256 and ConvNetD7 for 512. DDiF shows the best performance over all resolutions. Figure 4b shows the amount of the test performance decrease percentage with the resolution change i.e.  $(ACC_{org} - ACC_{test})/ACC_{org}$ . DDiF shows the least decrease with regard to resolution difference, being evidently robust to resolution change. Its robustness opens a new adaptability of dataset distillation methods to more wide-range situations. Please refer to Appendix C for the experimental results on different resizing techniques and the same architecture.

### 4.3 MORE ANALYSES

**Fixed Number of Decoded Instances.** To explore coding efficiency and expressiveness, we fixed the number of decoded synthetic instances and examined the performance by varying the budget required for each synthetic instance. Figure 5 shows the performance curve when the number of decoded synthetic instances is fixed as 1 per class. We observe that DDiF envelops the performance curve of baselines. Specifically, it is empirical evidence of Theorem 2 that the performance curve of DDiF envelops the performance curve of FreD. We also emphasize that DDiF maintains high performance even with a small budget, whereas the baselines show significant performance degradation as the allocated budget for each instance decreases. These results demonstrate that DDiF exhibits higher coding efficiency and expressiveness compared to previous studies.

486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539

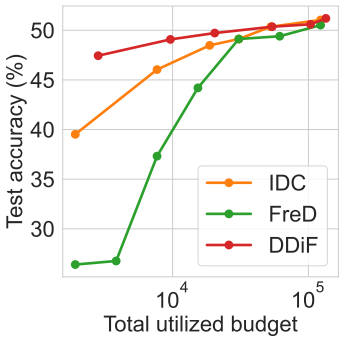


Figure 5: Test accuracies (%) under one decoded instance per class.

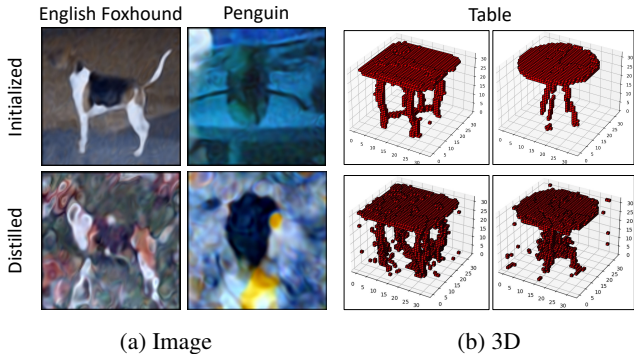


Figure 6: Visualization of the decoded synthetic instances from DDiF on various modalities.

**Qualitative Analysis.** Figure 6 visualizes the decoded synthetic instances by DDiF on various modalities. As shown in the first row in Figure 6, DDiF effectively decodes high-dimensional data even with a very small budget, regardless of the modality. For instance, each synthetic neural field utilizes a budget that is 1.96% of the original data size for images and 2.87% for 3D. After the distillation stage, each decoded synthetic instance involves class-discriminative features, even with significant budget reductions (See the second row). Notably, since the synthetic neural field is a continuous function, the quantity changes smoothly as the position changes. Please refer to Appendix C for more visualization results.

**Ablation Studies.** DDiF employs a neural field, which is a neural network, to store the distilled information. We conducted several ablation studies to investigate the effect of different components of the neural network. Figures 7a and 7b visualize the performance variation as the number of layers and the network width change, respectively. As the number of layers and the network width increases, the number of decoded synthetic instances decreases because the number of parameters in each synthetic neural field also increases. By comparing the layer and width, the width has fewer decoded synthetic instances, but it has a more gradual performance change than the layer. To explain the rationale, as shown in Eq. (5), width  $d$  directly affects  $k$ , which is involved in the number of basis functions, while  $L$  only affects other factors. Thus, increasing  $d$  leads to a modest change, as the increase in expressiveness offsets the reduction in quantity. In contrast, increasing  $L$  results in a relatively larger change, as the expressiveness remains similar while the quantity decreases. Even when the structure of the neural field is changed as FFN (Tancik et al., 2020), DDiF consistently shows performance improvement (see Figure 7c).

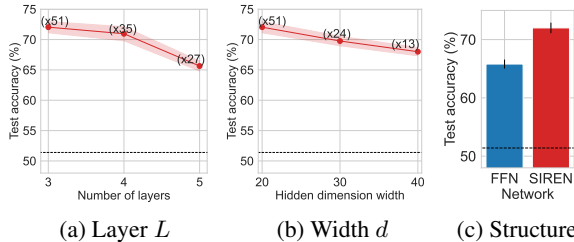


Figure 7: Ablation studies on (a) layer  $L$ , (b) width  $d$ , and (c) structure of neural field. The bottom black dashed line indicates the performance of Vanilla.

## 5 CONCLUSION

This paper introduces DDiF, a novel parameterization framework for dataset distillation that encodes information from large-scale datasets into synthetic neural fields under a constrained storage budget. By utilizing neural fields, DDiF efficiently captures distilled information and can easily decode data of varying sizes. We provide a theoretical analysis of dataset distillation parameterization, focusing on the feasible space of decoded synthetic instances. Additionally, we demonstrate that DDiF possesses a larger feasible space compared to previous methods, indicating greater expressiveness. In various evaluation scenarios, DDiF consistently exhibits improvements in performance, generalization, robustness, and adaptability across diverse modality datasets. Please refer to Appendix E for the limitation of DDiF.



## REFERENCES

- 540  
541  
542 Yoshua Bengio, Yann LeCun, et al. Scaling learning algorithms towards ai. *Large-scale kernel*  
543 *machines*, 34(5):1–41, 2007.
- 544  
545 Zalán Borsos, Mojmír Mutný, and Andreas Krause. Coresets via bilevel optimization for continual  
546 learning and streaming. *Advances in neural information processing systems*, 33:14879–14890,  
547 2020.
- 548  
549 George Cazenavette, Tongzhou Wang, Antonio Torralba, Alexei A Efros, and Jun-Yan Zhu. Dataset  
550 distillation by matching training trajectories. In *Proceedings of the IEEE/CVF Conference on*  
*Computer Vision and Pattern Recognition*, pp. 4750–4759, 2022.
- 551  
552 George Cazenavette, Tongzhou Wang, Antonio Torralba, Alexei A Efros, and Jun-Yan Zhu. Gener-  
553 alizing dataset distillation via deep generative prior. In *Proceedings of the IEEE/CVF Conference*  
*on Computer Vision and Pattern Recognition*, pp. 3739–3748, 2023.
- 554  
555 Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li,  
556 Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, et al. Shapenet: An information-rich 3d  
557 model repository. *arXiv preprint arXiv:1512.03012*, 2015.
- 558  
559 Justin Cui, Ruochen Wang, Si Si, and Cho-Jui Hsieh. Scaling up dataset distillation to imagenet-  
560 1k with constant memory. In *International Conference on Machine Learning*, pp. 6565–6590.  
PMLR, 2023.
- 561  
562 George Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics of control,*  
*signals and systems*, 2(4):303–314, 1989.
- 563  
564 Zhiwei Deng and Olga Russakovsky. Remember the past: Distilling datasets into addressable mem-  
565 ories for neural networks. *Advances in Neural Information Processing Systems*, 35:34391–34404,  
566 2022.
- 567  
568 Alexey Dosovitskiy. An image is worth 16x16 words: Transformers for image recognition at scale.  
569 *arXiv preprint arXiv:2010.11929*, 2020.
- 570  
571 Yuxuan Duan, Jianfu Zhang, and Liqing Zhang. Dataset distillation in latent space. *arXiv preprint*  
*arXiv:2311.15547*, 2023.
- 572  
573 Rakesh Dugad and Narendra Ahuja. A fast scheme for image size change in the compressed domain.  
574 *IEEE Transactions on Circuits and Systems for Video Technology*, 11(4):461–474, 2001.
- 575  
576 Emilien Dupont, Yee Whye Teh, and Arnaud Doucet. Generative models as distributions of func-  
577 tions. *arXiv preprint arXiv:2102.04776*, 2021.
- 578  
579 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-  
580 nition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp.  
770–778, 2016.
- 581  
582 Jeremy Howard. A smaller subset of 10 easily classified classes from imagenet and a little more  
583 french. URL <https://github.com/fastai/imagenette>, 2019.
- 584  
585 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child,  
586 Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language  
models. *arXiv preprint arXiv:2001.08361*, 2020.
- 587  
588 Soomro Khurram. Ucf101: A dataset of 101 human actions classes from videos in the wild. *arXiv*  
*preprint arXiv: 1212.0402*, 2012.
- 589  
590 Jang-Hyun Kim, Jinuk Kim, Seong Joon Oh, Sangdoon Yun, Hwanjun Song, Joonhyun Jeong, Jung-  
591 Woo Ha, and Hyun Oh Song. Dataset condensation via efficient synthetic-data parameterization.  
In *International Conference on Machine Learning*, pp. 11102–11118. PMLR, 2022.
- 592  
593 Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017. URL  
<https://arxiv.org/abs/1412.6980>.

- 594 Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.  
595 2009.  
596
- 597 Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep  
598 convolutional neural networks. In F. Pereira, C.J. Burges, L. Bottou, and K.Q. Weinberger  
599 (eds.), *Advances in Neural Information Processing Systems*, volume 25. Curran Associates, Inc.,  
600 2012a. URL [https://proceedings.neurips.cc/paper\\_files/paper/2012/  
601 file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2012/file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf).
- 602 Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convo-  
603 lutional neural networks. *Advances in neural information processing systems*, 25, 2012b.  
604
- 605 Hae Beom Lee, Dong Bok Lee, and Sung Ju Hwang. Dataset condensation with latent space knowl-  
606 edge factorization and sharing. *arXiv preprint arXiv:2208.10494*, 2022.
- 607 Shiye Lei and Dacheng Tao. A comprehensive survey of dataset distillation. *IEEE Transactions on  
608 Pattern Analysis and Machine Intelligence*, 2023.  
609
- 610 Yunzhu Li, Shuang Li, Vincent Sitzmann, Pulkit Agrawal, and Antonio Torralba. 3d neural scene  
611 representations for visuomotor control. In *Conference on Robot Learning*, pp. 112–123. PMLR,  
612 2022.
- 613 Songhua Liu, Kai Wang, Xingyi Yang, Jingwen Ye, and Xinchao Wang. Dataset distillation via  
614 factorization. *arXiv preprint arXiv:2210.16774*, 2022.
- 615 Lars Morten Mescheder. *Stability and expressiveness of deep generative models*. PhD thesis, Uni-  
616 versität Tübingen, 2020.  
617
- 618 Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and  
619 Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications  
620 of the ACM*, 65(1):99–106, 2021.
- 621 Brian B Moser, Federico Raue, Sebastian Palacio, Stanislav Frolov, and Andreas Dengel. Latent  
622 dataset distillation with diffusion models. *arXiv preprint arXiv:2403.03881*, 2024.  
623
- 624 Tiago Novello. Understanding sinusoidal neural networks. *arXiv preprint arXiv:2212.01833*, 2022.  
625
- 626 Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove.  
627 DeepSDF: Learning continuous signed distance functions for shape representation. In *Proceedings  
628 of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 165–174, 2019.
- 629 Nasim Rahaman, Aristide Baratin, Devansh Arpit, Felix Draxler, Min Lin, Fred Hamprecht, Yoshua  
630 Bengio, and Aaron Courville. On the spectral bias of neural networks. In *International Confer-  
631 ence on Machine Learning*, pp. 5301–5310. PMLR, 2019.
- 632 Edgar Riba, Dmytro Mishkin, Daniel Ponsa, Ethan Rublee, and Gary Bradski. Kornia: an open  
633 source differentiable computer vision library for pytorch. In *Proceedings of the IEEE/CVF Winter  
634 Conference on Applications of Computer Vision*, pp. 3674–3683, 2020.  
635
- 636 Noveen Sachdeva and Julian McAuley. Data distillation: A survey. *arXiv preprint  
637 arXiv:2301.04272*, 2023.
- 638 Noveen Sachdeva, Zexue He, Wang-Cheng Kang, Jianmo Ni, Derek Zhiyuan Cheng, and Julian  
639 McAuley. Farzi data: Autoregressive data distillation. *arXiv preprint arXiv:2310.09983*, 2023.  
640
- 641 Liyue Shen, John Pauly, and Lei Xing. Nerp: implicit neural representation learning with prior  
642 embedding for sparsely sampled image reconstruction. *IEEE Transactions on Neural Networks  
643 and Learning Systems*, 35(1):770–782, 2022.
- 644 Donghyeok Shin, Seungjae Shin, and Il-Chul Moon. Frequency domain-based dataset distillation.  
645 *Advances in Neural Information Processing Systems*, 36, 2024.  
646
- 647 Connor Shorten and Taghi M Khoshgoftaar. A survey on image data augmentation for deep learning.  
*Journal of big data*, 6(1):1–48, 2019.

- 648 Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image  
649 recognition. *arXiv preprint arXiv:1409.1556*, 2014.  
650
- 651 Vincent Sitzmann, Julien Martel, Alexander Bergman, David Lindell, and Gordon Wetzstein. Im-  
652 plicit neural representations with periodic activation functions. *Advances in neural information  
653 processing systems*, 33:7462–7473, 2020.
- 654 Ivan Skorokhodov, Savva Ignatyev, and Mohamed Elhoseiny. Adversarial generation of continuous  
655 images. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*,  
656 pp. 10753–10764, 2021.  
657
- 658 K Soomro. Ucf101: A dataset of 101 human actions classes from videos in the wild. *arXiv preprint  
659 arXiv:1212.0402*, 2012.
- 660 Duo Su, Junjie Hou, Weizhi Gao, Yingjie Tian, and Bowen Tang.  $D^4$ : Dataset distillation via  
661 disentangled diffusion model. In *Proceedings of the IEEE/CVF Conference on Computer Vision  
662 and Pattern Recognition*, pp. 5809–5818, 2024a.  
663
- 664 Duo Su, Junjie Hou, Guang Li, Ren Togo, Rui Song, Takahiro Ogawa, and Miki Haseyama. Gener-  
665 ative dataset distillation based on diffusion model. *arXiv preprint arXiv:2408.08610*, 2024b.  
666
- 667 Iliia Sucholutsky and Matthias Schonlau. Soft-label dataset distillation and text dataset distillation.  
668 In *2021 International Joint Conference on Neural Networks (IJCNN)*, pp. 1–8. IEEE, 2021.
- 669 Matthew Tancik, Pratul Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh  
670 Singhal, Ravi Ramamoorthi, Jonathan Barron, and Ren Ng. Fourier features let networks learn  
671 high frequency functions in low dimensional domains. *Advances in neural information processing  
672 systems*, 33:7537–7547, 2020.
- 673 Haohan Wang, Xindi Wu, Zeyi Huang, and Eric P Xing. High-frequency component helps explain  
674 the generalization of convolutional neural networks. In *Proceedings of the IEEE/CVF conference  
675 on computer vision and pattern recognition*, pp. 8684–8694, 2020a.  
676
- 677 Kai Wang, Jianyang Gu, Daquan Zhou, Zheng Zhu, Wei Jiang, and Yang You. Dim: Distilling  
678 dataset into generative model. *arXiv preprint arXiv:2303.04707*, 2023.
- 679 Tongzhou Wang, Jun-Yan Zhu, Antonio Torralba, and Alexei A Efros. Dataset distillation. *arXiv  
680 preprint arXiv:1811.10959*, 2018.  
681
- 682 Zhihao Wang, Jian Chen, and Steven CH Hoi. Deep learning for image super-resolution: A survey.  
683 *IEEE transactions on pattern analysis and machine intelligence*, 43(10):3365–3387, 2020b.  
684
- 685 Ziyu Wang, Yue Xu, Cewu Lu, and Yong-Lu Li. Dancing with still images: Video distillation  
686 via static-dynamic disentanglement. In *Proceedings of the IEEE/CVF Conference on Computer  
687 Vision and Pattern Recognition*, pp. 6296–6304, 2024.
- 688 Pete Warden. Speech commands: A dataset for limited-vocabulary speech recognition. *arXiv  
689 preprint arXiv:1804.03209*, 2018.
- 690 Xing Wei, Anjia Cao, Funing Yang, and Zhiheng Ma. Sparse parameterization for epitomic dataset  
691 distillation. *Advances in Neural Information Processing Systems*, 36, 2024.  
692
- 693 Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and Jianxiong  
694 Xiao. 3d shapenets: A deep representation for volumetric shapes. In *Proceedings of the IEEE  
695 conference on computer vision and pattern recognition*, pp. 1912–1920, 2015.
- 696 Yiheng Xie, Towaki Takikawa, Shunsuke Saito, Or Litany, Shiqin Yan, Numair Khan, Federico  
697 Tombari, James Tompkin, Vincent Sitzmann, and Srinath Sridhar. Neural fields in visual comput-  
698 ing and beyond. In *Computer Graphics Forum*, volume 41, pp. 641–676. Wiley Online Library,  
699 2022.
- 700
- 701 Zhi-Qin John Xu, Yaoyu Zhang, Tao Luo, Yanyang Xiao, and Zheng Ma. Frequency principle:  
Fourier analysis sheds light on deep neural networks. *arXiv preprint arXiv:1901.06523*, 2019.

702 Shaolei Yang, Shen Cheng, Mingbo Hong, Haoqiang Fan, Xing Wei, and Shuaicheng Liu. Neural  
703 spectral decomposition for dataset distillation. *arXiv preprint arXiv:2408.16236*, 2024.  
704

705 Ruonan Yu, Songhua Liu, and Xinchao Wang. Dataset distillation: A comprehensive review. *IEEE*  
706 *Transactions on Pattern Analysis and Machine Intelligence*, 2023.

707 Guangming Zang, Ramzi Idoughi, Rui Li, Peter Wonka, and Wolfgang Heidrich. Intratomo: self-  
708 supervised learning-based tomography via sinogram synthesis and prediction. In *Proceedings of*  
709 *the IEEE/CVF International Conference on Computer Vision*, pp. 1960–1970, 2021.  
710

711 Bo Zhao and Hakan Bilen. Dataset condensation with differentiable siamese augmentation. In  
712 *International Conference on Machine Learning*, pp. 12674–12685. PMLR, 2021.

713 Bo Zhao and Hakan Bilen. Dataset condensation with distribution matching. In *Proceedings of the*  
714 *IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 6514–6523, 2023.  
715

716 Bo Zhao, Konda Reddy Mopuri, and Hakan Bilen. Dataset condensation with gradient matching.  
717 *arXiv preprint arXiv:2006.05929*, 2020.

718 Haizhong Zheng, Jiachen Sun, Shutong Wu, Bhavya Kailkhura, Zhuoqing Mao, Chaowei Xiao, and  
719 Atul Prakash. Leveraging hierarchical feature sharing for efficient dataset condensation. *arXiv*  
720 *preprint arXiv:2310.07506*, 2023.  
721

722 Xinhao Zhong, Hao Fang, Bin Chen, Xulin Gu, Tao Dai, Meikang Qiu, and Shu-Tao Xia. Hierar-  
723 chical features matter: A deep exploration of gan priors for improved dataset distillation. *arXiv*  
724 *preprint arXiv:2406.05704*, 2024.  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

## A PROOFS

### A.1 PROOF OF OBSERVATION 1

**Observation 1.** Let two matrix variables  $X_1 := [x_{11}, \dots, x_{1M}]$  and  $X_2 := [x_{21}, \dots, x_{2M}]$  consisting of columns  $x_{ij} \in \mathcal{X}_i \subseteq \mathbb{R}^d$  for  $i = 1, 2$  and  $j = 1, \dots, M$ . If  $\mathcal{X}_1 \subseteq \mathcal{X}_2$ , then  $\min_{X_1 \in \mathcal{X}_1^M} \mathcal{L}(X_1) \geq \min_{X_2 \in \mathcal{X}_2^M} \mathcal{L}(X_2)$ .

*Proof.* For  $i = 1, 2$ , let  $X_i^* = \arg \min_{X_i \in \mathcal{X}_i^M} \mathcal{L}(X_i)$ . Note that  $\mathcal{X}_1^M \subseteq \mathcal{X}_2^M$  since  $\mathcal{X}_1 \subseteq \mathcal{X}_2$ . By the definition of  $X_2^*$ , for any  $X_2 \in \mathcal{X}_2^M$ ,  $\mathcal{L}(X_2^*) \leq \mathcal{L}(X_2)$ . Since  $X_1^* \in \mathcal{X}_1^M \subseteq \mathcal{X}_2^M$ ,  $\mathcal{L}(X_2^*) \leq \mathcal{L}(X_1^*)$ . In conclusion,  $\min_{X_1 \in \mathcal{X}_1^M} \mathcal{L}(X_1) \geq \min_{X_2 \in \mathcal{X}_2^M} \mathcal{L}(X_2)$ .  $\square$

### A.2 PROOF OF OBSERVATION 2

**Observation 2.** Consider two functions  $g_1, g_2$  where  $g_i : \mathcal{Z}_i \rightarrow \mathbb{R}^d$  for  $i = 1, 2$ . Also, consider two matrix variables  $Z_i := [z_{i1}, \dots, z_{iM}]$  where their columns  $z_{ij} \in \mathcal{Z}_i$  for  $i = 1, 2$  and  $j = 1, \dots, M$ . We denote  $\widehat{g}_i(Z_i) := [g_i(z_{i1}), \dots, g_i(z_{iM})]$  for  $i = 1, 2$ . Set  $\mathcal{G}_i := \{g | g : \mathcal{Z}_i \rightarrow \mathbb{R}^d\}$  for  $i = 1, 2$ . If  $g_1(\mathcal{Z}_1) \subseteq g_2(\mathcal{Z}_2)$  for any  $g_1 \in \mathcal{G}_1$  and  $g_2 \in \mathcal{G}_2$ , then  $\min_{g_1 \in \mathcal{G}_1, Z_1 \in \mathcal{Z}_1^M} \mathcal{L}(\widehat{g}_1(Z_1)) \geq \min_{g_2 \in \mathcal{G}_2, Z_2 \in \mathcal{Z}_2^M} \mathcal{L}(\widehat{g}_2(Z_2))$ .

*Proof.* For  $i = 1, 2$ , let  $g_i^*, Z_i^* = \arg \min_{g_i \in \mathcal{G}_i, Z_i \in \mathcal{Z}_i^M} \mathcal{L}(\widehat{g}_i(Z_i))$ . Note that  $g_1^*(z_{1j}^*) \in g_1^*(\mathcal{Z}_1) \subseteq g_2(\mathcal{Z}_2)$  for  $j = 1, \dots, M$  and any  $g_2 \in \mathcal{G}_2$ . It implies that there exists some  $g_2 \in \mathcal{G}_2, Z_2 \in \mathcal{Z}_2^M$  such that  $\widehat{g}_2(Z_2) = \widehat{g}_1^*(Z_1^*)$ . By the definition of  $g_2^*, Z_2^*$ , for any  $g_2 \in \mathcal{G}_2, Z_2 \in \mathcal{Z}_2^M$ ,  $\mathcal{L}(g_2^*(Z_2^*)) \leq \mathcal{L}(\widehat{g}_2(Z_2))$ . Therefore,  $\mathcal{L}(g_2^*(Z_2^*)) \leq \mathcal{L}(\widehat{g}_1^*(Z_1^*))$ . In conclusion,  $\min_{g_1 \in \mathcal{G}_1, Z_1 \in \mathcal{Z}_1^M} \mathcal{L}(\widehat{g}_1(Z_1)) \geq \min_{g_2 \in \mathcal{G}_2, Z_2 \in \mathcal{Z}_2^M} \mathcal{L}(\widehat{g}_2(Z_2))$ .  $\square$

### A.3 DERIVATION OF THE EQ. (5)

Recall that  $\alpha_k := \sum_{i=1}^d A_{k,i} \sin(\varphi_{k,i})$  and  $\beta_k := \sum_{i=1}^d A_{k,i} \cos(\varphi_{k,i})$ .

$$\begin{aligned}
 F_\psi(x) &= b^{(2)} + \sum_{k \in \mathbb{Z}^d} \alpha_k \cos(\omega_k x) + \beta_k \sin(\omega_k x) \\
 &= b^{(2)} + \sum_{k \in \mathbb{Z}^d} \left\{ \sum_{i=1}^d A_{k,i} \sin(\varphi_{k,i}) \cos(\omega_k x) + \sum_{i=1}^d A_{k,i} \cos(\varphi_{k,i}) \sin(\omega_k x) \right\} \\
 &= b^{(2)} + \sum_{k \in \mathbb{Z}^d} \sum_{i=1}^d A_{k,i} \sin(\omega_k x + \varphi_{k,i}) \\
 &= b^{(2)} + \sum_{k \in \mathbb{Z}^d} \sum_{i=1}^d A_{k,i} \cos(\omega_k x + \varphi'_{k,i}) \quad \text{where } \varphi'_{k,i} = \varphi_{k,i} - \frac{\pi}{2}
 \end{aligned}$$

### A.4 PROOF OF THEOREM 2

**Theorem 2.** For any  $x \in \mathcal{C}_N$ ,  $g(x; \Gamma) \subsetneq F_\psi(x)$ .

*Proof.* To prove  $g(x; \Gamma) \subseteq F_\psi(x)$ , it is sufficient to show that there exist neural network parameters  $\psi = \{W^{(j)}, b^{(j)}\}_{j=0}^2$  which satisfy  $g(x; \Gamma) = F_\psi(x)$  for any  $\Gamma$  and  $x \in \mathcal{C}_N = \{0, \dots, N-1\}$ . First, we decompose Eq. (5) into a sum over  $\mathcal{K} \subseteq \mathbb{Z}^d$  and the other terms:

$$\begin{aligned}
 F_\psi(x) &= b^{(2)} + \sum_{k \in \mathbb{Z}^d} \sum_{i=1}^d A_{k,i} \cos(\omega_k x + \varphi'_{k,i}) \\
 &= b^{(2)} + \sum_{k \in \mathcal{K}} \sum_{i=1}^d A_{k,i} \cos(\omega_k x + \varphi'_{k,i}) + \sum_{k \in \mathbb{Z}^d \setminus \mathcal{K}} \sum_{i=1}^d A_{k,i} \cos(\omega_k x + \varphi'_{k,i})
 \end{aligned}$$

Let  $\mathcal{K} = \{u\mathbf{e}_1 | u \in \mathcal{U}\}$  where  $\mathbf{e}_1 = [1, 0, \dots, 0]^T$  is the standard basis vector in  $\mathbb{R}^d$ . Set  $W^{(0)} = \frac{\pi}{N}\mathbf{e}_1$ ,  $b^{(0)} = \frac{\pi}{2N}\mathbf{e}_1$ ,  $b^{(1)} = \frac{\pi}{2}\mathbf{e}_1$ , and  $b^{(2)} = -\sum_{k \in \mathbb{Z}^d \setminus \mathcal{K}} \sum_{i=1}^d A_{k,i} \cos(\omega_k x + \varphi'_{k,i})$ . Since the absolute value of the Bessel function of the first kind has an upper bound  $|J_p(r)| < \infty$  for any  $p$  and  $r$ ,  $W^{(1)}$  and  $W^{(2)}$  can be configured to satisfy  $A_{k,i} = W_i^{(2)} \prod_{j=1}^d J_{k_j}(W_{ij}^{(1)}) = \gamma_u$ . Under these parameters,  $F_\psi(x)$  is same as  $g(x; \Gamma)$ :

$$\begin{aligned} F_\psi(x) &= b^{(2)} + \sum_{k \in \mathcal{K}} \sum_{i=1}^d A_{k,i} \cos(\omega_k x + \varphi'_{k,i}) + \sum_{k \in \mathbb{Z}^d \setminus \mathcal{K}} \sum_{i=1}^d A_{k,i} \cos(\omega_k x + \varphi'_{k,i}) \\ &= \sum_{u \in \mathcal{U}} \gamma_u \cos\left(\frac{\pi u}{N}x + \frac{\pi u}{2N}\right) = g(x; \Gamma) \end{aligned}$$

Next, we prove that there is no  $\Gamma$  such that  $g(x; \Gamma) = F_\psi(x)$  for some  $\psi$  and  $x \in \mathcal{C}_N$ . Note that it is possible to choose  $k$  and  $W^{(0)}$  which satisfy  $\omega_k = \langle k, W^{(0)} \rangle \neq \frac{\pi u}{N}$  for any  $u \in \mathcal{U}$ . Due to the orthogonality of cosine functions having different frequencies, there is no  $\Gamma$  that represents  $\cos(\omega_k x + \varphi'_{k,i})$  terms. It implies  $g(x; \Gamma)$  cannot express  $F_\psi(x)$  with some  $\psi$  parameters.  $\square$

## B EXPERIMENTAL DETAILS

### B.1 DATASETS

**Image Domain.** We evaluate DDiF on a various benchmark image datasets. 1) ImageNet-Subset (Howard, 2019; Cazenavette et al., 2022) is dataset consists of a subset of similar characteristics in the ImageNet. In the experiment, we consider various types of subsets by following Cazenavette et al. (2022): ImageNette (various objects), ImageWoof (dog breeds), ImageFruit (fruit category), ImageMeow (cats), ImageSquawk (birds), ImageYellow (yellowish objects). Each subset has 10 classes and more than 10,000 instances. We utilize two types of resolution:  $128 \times 128$  and  $256 \times 256$ . 2) CIFAR-10 (Krizhevsky et al., 2009) consists of 60,000 RGB images in 10 classes. Each image has a  $32 \times 32$  size. Each class contains 5,000 images for training and 1,000 images for testing. 3) CIFAR-100 (Krizhevsky et al., 2009) consists of 60,000  $32 \times 32$  RGB images of 100 categories. Each class is split into 500 for training and 100 for testing.

**Video Domain.** We utilize MiniUCF (Wang et al., 2024), a subset of UCF101 (Soomro, 2012) which includes 50 classes. The videos are sampled to 16 frames, and the frames are cropped and resized to  $112 \times 112$ . Each data has  $16 \times 3 \times 112 \times 112$  size.

**Audio Domain.** We utilize Mini Speech Commands (Kim et al., 2022), a subset of the original Speech Commands dataset (Warden, 2018). We follow the data processing of Kim et al. (2022). The dataset consists of 8 classes, and each class has 875/125 data for training/testing. Each data is  $64 \times 64$  log-scale magnitude spectrograms by short time Fourier transform (STFT).

**3D Domain.** We utilize a core version of ModelNet-10 (Wu et al., 2015) and ShapeNet (Chang et al., 2015), which are widely used in 3D. They includes 10 classes and 16 classes, respectively. Each 3D point cloud data is converted into  $32 \times 32 \times 32$  voxel.

### B.2 NETWORK ARCHITECTURES.

**ConvNet.** By following previous studies, we leverage the ConvNetDn as a default network architecture for both distillation and evaluation of synthetic datasets. The ConvNetDn is a convolutional neural network with  $n$  duplicate blocks. Each  $n$  blocks consist of a convolution layer with  $3 \times 3$ -shape 128 filters, an instance normalization layer, ReLU, and an average pooling with  $2 \times 2$  kernel size with stride 2. Lastly, it contains a linear classifier, which outputs the logits. Depending on the resolution of real dataset, we utilize different depth  $n$ . Specifically, ConvNetD3 for  $32 \times 32$  CIFAR-10 and CIFAR-100, ConvNetD4 for  $64 \times 64$  Audio spectrograms, ConvNetD5 for  $128 \times 128$  ImageNet-Subset, ConvNetD6 for  $256 \times 256$  ImageNet-Subset, and ConvNetD7 for  $512 \times 512$  ImageNet-Subset.



**AlexNet.** AlexNet is a basic convolutional neural network architecture suggested in (Krizhevsky et al., 2012a). It consists of 5 convolution layers, 3 max-pooling layers, 2 Normalized layers, 2 fully connected layers and 1 SoftMax layer. In each convolution layer, ReLU activation function is utilized. We adopt this network to evaluate cross-architecture performance of DDiF.

**VGG11.** VGG11 (Simonyan & Zisserman, 2014) is also applied for evaluation, which attributes to 11 weighted layers. It consists of 8 convolution layers and 3 fully connected layers. Its design is straightforward yet powerful, providing a balance between depth and computational efficiency. The number of trainable parameters is around 132M, making it larger than earlier models but still suitable for medium-scale tasks. We adopt this network to evaluate the cross-architecture performance of DDiF.

**ResNet18.** ResNet18 (He et al., 2016) introduces residual connections, which help mitigate the vanishing gradient problem in deep networks by allowing gradients to bypass certain layers. It consists of 18 layers with 4 residual blocks, each composed of two convolutional layers followed by activation and normalization with around 11M trainable parameters. We utilize ResNet18 as one of the architecture for evaluating synthetic datasets.

**ViT.** Vision Transformer (Dosovitskiy, 2020) utilizes the transformer architecture, initially designed for sequence modeling tasks in NLP. For image classification, it divides images into non-overlapping patches and processes them as a sequence using self-attention mechanisms. ViT has around 10M trainable parameters in its base form and offers a competitive alternative to CNNs, demonstrating the effectiveness of transformers in vision tasks. We selected ViT as the final network to evaluate synthetic image datasets.

**Conv3DNet.** For 3D domain, we utilize Conv3DNet (Shin et al., 2024), a 3D version of ConvNet. Conv3DNet consists of three repeated blocks, each containing a  $3 \times 3 \times 3$  convolutional layer with 64 filters, 3D instance normalization, ReLU activation, and 3D average pooling with a  $2 \times 2 \times 2$  filter and a stride of 2. Lastly, it contains a linear classifier.

### B.3 BASELINES.

Since our main focus lies on the parameterization of dataset distillation, we compare DDiF with 1) static decoding, which are IDC (Kim et al., 2022) and FreD (Shin et al., 2024); 2) parameterized decoding, which are RTP (Deng & Russakovsky, 2022), HaBa (Liu et al., 2022), SPEED (Wei et al., 2024), LatentDD (Duan et al., 2023), and NSD (Yang et al., 2024); and 3) deep generative prior, which include GLaD (Cazenavette et al., 2023), H-GLaD (Zhong et al., 2024), and LD3M (Moser et al., 2024). We also demonstrate the performance improvement of DDiF compared to input-sized parameterization, denoted as Vanilla.

### B.4 IMPLEMENTATION SETTINGS.

Although any loss can be adapted to DDiF, we utilize TM (Cazenavette et al., 2022) for  $\mathcal{L}$  as a default unless specified. Following previous studies, we use DSA (Zhao & Bilen, 2021), which consists of color jittering, cropping, cutout, flipping, scaling, and rotation. We adopt ZCA whitening on CIFAR-10 (IPC=1, 10) and CIFAR-100 (IPC=1) with the Kornia (Riba et al., 2020) implementation. We adopt SIREN (Sitzmann et al., 2020) for synthetic field  $F_\psi$  as a default. SIREN is a multilayer perceptron with a sinusoidal activation function, and it is widely used in the neural field area due to its simple structure. We use the same width across all layers in a synthetic neural field i.e.  $d_l = d$  for all  $l$ . We utilize normalized coordinates defined on  $[-1, 1]^n$  for  $n$ -dimension data to enhance stability (Sitzmann et al., 2020), rather than using integer coordinates, which have a wide range. For cross-resolution experiments, we utilize the coordinate set  $C$ , which consists of evenly spaced points within the interval  $[-1, 1]$  according to the target resolution. We provide the detailed configuration of the synthetic neural field, the resulting size of each neural field, the number of synthetic instances per class, and the total number of neural fields in Table 7. We use Adam optimizer (Kingma & Ba, 2017) for all experiments. We fix the iteration number and learning rate for warm-up initialization of synthetic neural field as 5000 and 0.0005. Without any description to distillation loss, we generally use matching training trajectory (TM) objective for dataset distillation loss  $\mathcal{L}$ . Following the

Table 7: Configurations of the synthetic neural field. In the case of Video, there is no increment of decoded instances because we experimented with the fixed number of decoded instances.

Modality	Dataset	IPC	$n$	$L$	$d$	$m$	size( $\psi$ )	Increment of decoded instances	
Image	CIFAR10	1	2	2	6	3	81	$\times 37$	
		10	2	2	6	3	81	$\times 37.9$	
		50	2	2	20	3	543	$\times 5.64$	
	CIFAR100	1	2	2	10	3	173	$\times 17$	
		10	2	2	15	3	333	$\times 9.2$	
		50	2	2	30	3	1113	$\times 2.76$	
	ImageNet-Subset (128)	1	2	3	20	3	963	$\times 51$	
		10	2	3	20	3	963	$\times 51$	
		50	2	3	40	3	3523	$\times 13.94$	
	ImageNet-Subset (256)	1	2	3	40	3	3523	$\times 55$	
		Video	1	3	6	40	3	8483	–
			5	3	6	40	3	8483	–
Audio	Mini Speech Commands	10	2	3	10	1	261	$\times 15.6$	
		20	2	3	10	1	261	$\times 15.6$	
3D	ModelNet	1	3	3	20	1	941	$\times 34$	
	ShapeNet	1	3	3	20	1	941	$\times 34$	

previous studies, we utilize two types of default TM hyperparameters same as SPEED (Wei et al., 2024) and FreD (Shin et al., 2024). We run 15,000 iterations for TM and 20,000 iterations for DM. We use a mixture of RTX 3090, L40S, and Tesla A100 to run our experiments. We follow the conventional evaluation procedure of the previous studies: train 5 randomly initialized networks with an optimized synthetic dataset and evaluate the classification performance. We provide the detailed hyperparameter in Table 23.

## B.5 ALGORITHM FOR DECODING PROCESS OF SYNTHETIC INSTANCES

The main difference between DDiF and previous parameterization methods is the decoding process for synthetic instances. Basically, the neural field takes coordinates as input and output quantities. To generate a single synthetic instance, each coordinate  $c \in \mathcal{C}$  is input into the synthetic neural field  $F_\psi$ , and then the resulting value is assigned to the corresponding coordinate of the decoded synthetic instance  $\tilde{x}_j^{(c)}$ . Algorithm 2 illustrates this synthetic instance decoding process of DDiF. In this algorithm, we included a loop over coordinates for clarity. However, it should be noted that in the actual implementation, the coordinate set is input to the neural network in a full-batch manner.

### Algorithm 2: Decoding process of Synthetic instances in DDiF

**Input:** Set of parameters of synthetic neural fields  $\Psi$ ; Coordinate set  $\mathcal{C}$

**Output:** Set of decoded synthetic instances  $X_S$

```

1 Initialize  $X_S \leftarrow \emptyset$ 
2 for  $j = 1$  to  $|\Psi|$  do
3   Initialize  $\tilde{x}_j$ 
4   for  $c \in \mathcal{C}$  do
5      $\tilde{x}_j^{(c)} \leftarrow F_{\psi_j}(c)$ 
6    $X_S \leftarrow X_S \cup \tilde{x}_j$ 

```

## C ADDITIONAL EXPERIMENTAL RESULTS

### C.1 PERFORMANCE COMPARISON ON LOW-DIMENSIONAL DATASETS

To verify the wide-applicability of DDiF, we conducted experiments on low-dimensional datasets, such as CIFAR-10 and CIFAR-100. In Table 8, DDiF exhibits highly competitive performances with previous studies. These results demonstrates that DDiF also properly applicable to low-dimensional datasets, while DDiF shows significant performance improvement in high-dimensional datasets.

Table 8: Test accuracies (%) on CIFAR-10 and CIFAR-100. **Bold** and Underline means best and second-best performance of each column, respectively. "—" indicates no reported results.

Dataset	CIFAR10			CIFAR100			
	1	10	50	1	10	50	
Input sized	TM	46.3±0.8	65.3±0.7	<b>71.6±0.2</b>	24.3±0.3	40.1±0.4	<b>47.7±0.2</b>
	FRePo	46.8±0.7	65.5±0.4	<b>71.7±0.2</b>	28.7±0.1	42.5±0.2	<b>44.3±0.2</b>
Static	IDC	50.0±0.4	67.5±0.5	<b>74.5±0.2</b>	—	—	—
	FreD	60.6±0.8	70.3±0.3	<b>75.8±0.1</b>	34.6±0.4	42.7±0.2	<b>47.8±0.1</b>
Parameterized	HaBa	48.3±0.8	69.9±0.4	<b>74.0±0.2</b>	—	—	<b>47.0±0.2</b>
	RTP	66.4±0.4	71.2±0.4	<b>73.6±0.5</b>	34.0±0.4	42.9±0.7	—
	HMN	65.7±0.3	73.7±0.1	<b>76.9±0.2</b>	36.3±0.2	45.4±0.2	<b>48.5±0.2</b>
	SPEED	63.2±0.1	73.5±0.2	<b>77.7±0.4</b>	40.4±0.4	45.9±0.3	<b>49.1±0.2</b>
	NSD	<b>68.5±0.8</b>	73.4±0.2	<b>75.2±0.6</b>	36.5±0.3	<b>46.1±0.2</b>	—
Function	DDiF	<u>66.5±0.4</u>	<b>74.0±0.4</b>	<b>77.5±0.3</b>	<b>42.1±0.2</b>	<u>46.0±0.2</u>	<b>49.9±0.2</b>

## C.2 ADDITIONAL PERFORMANCE COMPARISON ON HIGH-DIMENSIONAL DATASET

**Comparison under Large budget.** In general, high-dimensional dataset distillation under a large budget is rarely addressed in previous studies due to their significant computational cost. To demonstrate the efficacy of DDiF even in a large budget setting, we conducted experiments on ImageNette (128) under IPC=50. DDiF with TM achieves  $75.2\% \pm 1.3\%$  while vanilla with TM achieves  $72.8\% \pm 0.8\%$ . It means that DDiF effectively improves the dataset distillation performance even with a larger storage budget for high resolution.

**Comparison with Combination-based parameterization.** Building on the observation that there are common representations across classes, several studies have proposed utilizing class-shared bases to improve budget efficiency and reduce spatial redundancies. In particular, RTP (Deng & Russakovsky, 2022) generates synthetic instances by linearly combining class-shared bases, coefficients, and label vectors. Due to the advantage of combination, RTP has demonstrated high performance across various datasets. We believe it is necessary to compare RTP and DDiF on high-dimensional datasets, as the combination-based parameterization can also decode a large number of synthetic instances within the same budget.

Table 9: Test accuracies (%) on ImageNet-Subset ( $128 \times 128$ ) under IPC=1.

	Increment of decoded instances	Nette	Woof	Fruit
TM	×1	51.4 ± 2.3	29.7 ± 0.9	28.8 ± 1.2
TM+RTP	×64	69.6 ± 0.4	38.8 ± 1.1	45.2 ± 1.7
TM+DDiF	×51	<b>72.0 ± 0.9</b>	<b>42.9 ± 0.7</b>	<b>48.2 ± 1.2</b>

Table 9 presents the results of RTP and DDiF on the  $128 \times 128$  resolution ImageNet-Subset. For a fair comparison, we use the same dataset distillation loss (TM) for both. While RTP consistently shows performance improvements when applied to TM, DDiF achieves higher performance improvement even though DDiF shows smaller increments of decoded instances. These results suggest that DDiF enhances both the quantity and quality of synthetic instances.

## C.3 PERFORMANCE COMPARISON UNDER THE FIXED NUMBER OF DECODED INSTANCES

In the main paper, we primarily conducted performance comparisons under the same storage budget, which is the most general setting in dataset distillation. Herein, we investigate the expressiveness of parameterization methods by conducting performance comparisons under the fixed number of decoded instances. Through these experiments, we demonstrate that DDiF’s superiority stems not only from the diversity improvement but also from the quality improvement.

**Comparison with Input-sized parameterization.** First, we compare with input-sized parameterization. Under the IPC=1 setting, DDiF can decode 51 and 24 synthetic instances per class on 128 and 256 resolution, respectively. Therefore, we conducted experiments in two scenarios where the number of decoded synthetic instances was determined based on either 1) the number of Vanilla’s or 2) the number of DDiF’s. Table 10 presents the results as follows:

- When the number of decoded instances is small, there are some performance drop cases. However, we emphasize that 1) this performance drop is not larger than 2.3%p and 2) DDiF utilizes a much smaller budget. Moreover, it is interesting that DDiF achieves higher performance than vanilla in some cases. We believe that it is related to previous findings that input-sized parameterization includes superfluous or irrelevant information (Lei & Tao, 2023; Yu et al., 2023; Sachdeva & McAuley, 2023).
- When the number of decoded instances is large, DDiF utilizes much less budget but achieves performance comparable to input-sized parameterization.

These experimental results support that DDiF involves sufficient representational power while using a much smaller budget compared to the input-sized parameterization.

Table 10: Test accuracies (%) on ImageNet-Subset with input-sized parameterization (Vanilla) and DDiF. We utilize TM for  $128 \times 128$  resolution and DM for  $256 \times 256$  resolution. "DIPC" denotes the number of decoded instances per class.

Resolution	DIPC	Methods	Utilized Budget	Nette	Woof	Fruit
128	1	Vanilla	491,520	$51.4 \pm 2.3$	$29.7 \pm 0.9$	$28.8 \pm 1.2$
		DDiF	9,630	$49.1 \pm 2.0$	$29.4 \pm 0.7$	$27.3 \pm 1.3$
	51	Vanilla	$51 \times 491,520$	$73.0 \pm 0.7$	$42.8 \pm 0.7$	$48.2 \pm 0.7$
		DDiF	491,520	$72.0 \pm 0.9$	$42.9 \pm 0.7$	$48.2 \pm 1.2$
256	1	Vanilla	1,966,080	32.1	20.0	19.5
		DDiF	79,530	$31.2 \pm 0.8$	$21.2 \pm 0.9$	$21.3 \pm 1.5$
	55	Vanilla	$55 \times 1,966,080$	$70.1 \pm 1.0$	$37.5 \pm 1.2$	$41.3 \pm 0.8$
		DDiF	1,966,080	$67.8 \pm 1.0$	$39.6 \pm 1.6$	$43.2 \pm 1.7$

**Comparison with Parameterization methods.**

In addition, we conducted a performance comparison with FreD and SPEED, which show strong performance on high-resolution image datasets. Under the IPC=1 setting, FreD and SPEED can decode 8 and 15 synthetic instances per class on 128 resolution, respectively. In Table 11, under the same number of decoded instances, DDiF achieves higher performance while using less budget compared to FreD and SPEED. Furthermore, under the same budget, DDiF generates more decoded instances and achieves superior performance with a significant margin. These results repeatedly support the claim that DDiF exhibits high coding efficiency and expressiveness.

Table 11: Test accuracies (%) on ImageNet-Subset ( $128 \times 128$ ) with SPEED and DDiF. We utilize TM.

	DIPC	Utilized budget	Nette
FreD	8	491,520	$66.8 \pm 0.4$
SPEED	15	491,520	$66.9 \pm 0.7$
DDiF	8	77,040	$67.1 \pm 0.4$
	15	144,450	$68.3 \pm 1.1$
	51	491,520	$72.0 \pm 0.9$

C.4 ADDITIONAL RESULTS ON CROSS-RESOLUTION GENERALIZATION

**Detailed experimental results.** We apply the spatial-domain upsampling methods, such as nearest, bilinear, and bicubic, for the optimized synthetic instances of previous parameterization methods. Particularly, FreD can also utilize frequency-domain upsampling since it stores masked frequency coefficients. The most widely used for frequency-domain upsampling is zero-padding the frequency coefficients before inverse frequency transform, which means assigning zeros to the high-frequency components (Dugad & Ahuja, 2001). For example, in the case of DCT, which is the default setting of FreD, the process of upsampling an  $N$ -resolution frequency coefficient  $F$  to an  $M(> N)$ -resolution image can be described as

$\text{IDCT} \left( \begin{bmatrix} \lambda \times F & 0_{M-N} \\ 0_{M-N} & 0_{M-N} \end{bmatrix} \right)$  where  $\lambda = \left(\frac{M}{N}\right)^2$  is the scaling factor. We refer to this frequency-domain upsampling as "zero-padding". For DDiF, we constructed a coordinate set suitable for the test resolution, and then input it into the optimized synthetic neural field to generate the upsampled synthetic instance. We refer to this method as "coordinate interpolation".

Table 12 presents the detailed experimental results. We observe that the previous parameterizations show drastic performance degradation regardless of the interpolation method used. Whereas, DDiF still achieves the highest cross-resolution generalization performance. Also, Figure 8 presents the cross-resolution generalization performance under the same network architecture in the distillation stage. As seen in Figure 4, DDiF achieves the best performance and shows the least performance degradation over all resolutions. These extensive results consistently demonstrate that DDiF is robust to resolution field change, and this robustness is largely driven by the continuous nature of the synthetic neural field.

**Comparison with Full dataset downsampling.** The most intuitive and straightforward way to reduce the budget of a dataset is by downsampling, which reduces the budget size of each instance. It means that it is possible to downsample the full dataset to a specific resolution for storage and then upsample it to the test resolution. One may doubt that this simple method can show high cross-resolution generalization performance. To verify the superiority of DDiF on cross-resolution generalization, we additionally conducted performance comparison experiments with full dataset downsampling.

Table 13 shows the experiment results where the full dataset was downsampled and then upsampled to the original resolution during the test phase. When the budget is similar (when the downsampled resolution is  $4 \times 4$ ), DDiF outperforms the full dataset downsampling method. We also experimented when the downsampled resolution was the same as the resolution of the decoded synthetic instance. In this case, the downsampled dataset achieved better performance, as expected. However, this setting requires 1,289 times more budget than DDiF since the number of instances is not reduced. Such a setup deviates from the core purpose of dataset distillation, which aims to optimize performance under strict budget constraints.

Furthermore, we emphasize that the full dataset downsampling has two limitations on cross-resolution generalization. First, it requires a slightly different assumption. Cross-resolution generalization experiment, which we proposed, is modeled to evaluate the dataset distillation ability to generalize to higher resolution settings. This involves training on a low-resolution ( $128 \times 128$ ) synthetic dataset and testing on high-resolution ( $256 \times 256$ ) data. However, this full dataset downsampling experiment diverges from this cross-resolution setting since it assumes the availability of a high-resolution dataset ( $256 \times 256$ ). Second, the training time during the test phase increases since the number of data points remains the same as the full dataset. While the memory budget may

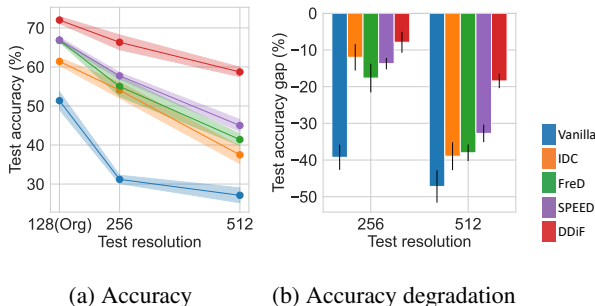


Figure 8: (a) Test accuracies (%) with different image resolutions. (b) Test accuracy gap (%) from original. We use bilinear interpolation for previous studies. We utilize ConvNetD5, which is the same architecture in the distillation stage.

Table 12: Test accuracies (%) with different resolutions and networks, we denote the difference between ordinary and cross-resolution performance as  $\text{Diff}(\%) = ACC_{org} - ACC_{test}$  and relative ratio as  $\text{Ratio} = \frac{ACC_{org} - ACC_{test}}{ACC_{org}}$ .

Test resolution	Test network	Method	Upsampling	Accuracy ( $\uparrow$ )	Diff ( $\downarrow$ )	Ratio ( $\downarrow$ )
256	ConvNetD5	Vanilla	nearest	29.5 $\pm$ 1.6	21.9	0.43
			bilinear	31.2 $\pm$ 1.1	20.2	0.39
			bicubic	30.7 $\pm$ 2.0	20.7	0.40
		IDC	nearest	54.8 $\pm$ 1.2	6.6	0.11
			bilinear	54.0 $\pm$ 2.0	7.4	0.12
			bicubic	55.0 $\pm$ 1.6	6.4	0.10
		SPEED	nearest	58.8 $\pm$ 1.4	8.1	0.12
			bilinear	57.7 $\pm$ 0.8	9.2	0.14
			bicubic	58.1 $\pm$ 1.0	8.8	0.13
		FreD	nearest	55.2 $\pm$ 2.2	11.6	0.17
			bilinear	55.0 $\pm$ 2.6	11.8	0.16
			bicubic	56.4 $\pm$ 1.4	10.4	0.16
			zero-padding	53.8 $\pm$ 1.4	13.0	0.19
		DDiF	coord. interpolation	66.3 $\pm$ 1.9	5.7	0.08
	ConvNetD6	Vanilla	nearest	44.0 $\pm$ 1.7	7.3	0.14
			bilinear	43.2 $\pm$ 1.1	8.2	0.16
			bicubic	43.9 $\pm$ 1.8	7.4	0.14
		IDC	nearest	54.7 $\pm$ 1.6	6.7	0.11
			bilinear	54.5 $\pm$ 1.6	6.9	0.11
			bicubic	55.4 $\pm$ 1.3	6.0	0.10
		SPEED	nearest	62.0 $\pm$ 1.0	4.9	0.07
			bilinear	61.8 $\pm$ 1.8	5.1	0.08
			bicubic	62.6 $\pm$ 1.1	4.3	0.06
		FreD	nearest	60.9 $\pm$ 0.8	5.9	0.09
bilinear			60.1 $\pm$ 0.3	6.7	0.10	
bicubic			61.4 $\pm$ 0.8	5.8	0.09	
		zero-padding	61.8 $\pm$ 1.0	5.0	0.07	
	DDiF	coord. interpolation	70.6 $\pm$ 1.2	1.4	0.02	
512	ConvNetD5	Vanilla	nearest	27.4 $\pm$ 1.4	24.0	0.47
			bilinear	27.1 $\pm$ 1.9	24.2	0.47
			bicubic	27.1 $\pm$ 1.0	24.3	0.47
		IDC	nearest	38.6 $\pm$ 2.7	22.8	0.37
			bilinear	37.5 $\pm$ 2.3	23.9	0.39
			bicubic	39.5 $\pm$ 2.1	21.9	0.36
		SPEED	nearest	43.3 $\pm$ 2.2	23.6	0.35
			bilinear	45.0 $\pm$ 1.5	21.9	0.33
			bicubic	44.8 $\pm$ 3.0	22.1	0.33
		FreD	nearest	42.5 $\pm$ 2.5	24.3	0.36
			bilinear	41.4 $\pm$ 1.5	25.4	0.38
			bicubic	41.6 $\pm$ 1.3	25.2	0.38
			zero-padding	42.9 $\pm$ 1.5	23.9	0.36
		DDiF	coord. interpolation	58.7 $\pm$ 1.2	13.3	0.18
	ConvNetD7	Vanilla	nearest	41.2 $\pm$ 1.5	10.1	0.20
			bilinear	40.6 $\pm$ 2.6	10.7	0.21
			bicubic	40.4 $\pm$ 1.9	10.9	0.21
		IDC	nearest	51.5 $\pm$ 1.8	9.9	0.16
			bilinear	50.7 $\pm$ 2.1	10.7	0.17
			bicubic	51.2 $\pm$ 2.8	10.7	0.17
		SPEED	nearest	59.6 $\pm$ 2.0	7.3	0.11
			bilinear	59.5 $\pm$ 2.0	7.4	0.11
			bicubic	60.1 $\pm$ 1.7	6.8	0.10
		FreD	nearest	55.1 $\pm$ 1.2	11.7	0.18
bilinear			53.8 $\pm$ 1.0	13.0	0.19	
bicubic			54.4 $\pm$ 0.9	12.4	0.19	
		zero-padding	56.3 $\pm$ 0.8	10.5	0.16	
	DDiF	coord. interpolation	69.0 $\pm$ 1.0	3.0	0.04	



be comparable, this increased time cost is undesirable, especially in scenarios where efficiency is critical.

Table 13: Performance comparison when the test resolution is  $256 \times 256$ . We utilize bicubic interpolation for full dataset resizing. The relative budget ratio indicates the ratio of full dataset downsampling over DDiF.

Test network	Method	Original resolution	Downsampled resolution	Relative budget ratio	Accuracy
ConvNetD5	Downsample	$256 \times 256$ $256 \times 256$	$4 \times 4$ $128 \times 128$	1.3 1,289.4	$45.2 \pm 2.3$ $91.3 \pm 0.5$
	DDiF	$128 \times 128$	–	1.0	$66.3 \pm 1.9$
ConvNetD6	Downsample	$256 \times 256$ $256 \times 256$	$4 \times 4$ $128 \times 128$	1.3 1,289.4	$44.7 \pm 0.4$ $91.2 \pm 0.0$
	DDiF	$128 \times 128$	–	1.0	$70.6 \pm 1.2$

### C.5 EMPIRICAL EVIDENCE ON OBSERVATION 1

The theoretical analysis in this paper begins with Observation 1, which states that as the feasible space of synthetic instances in input-sized parameterization increases, the performance of dataset distillation also improves. We have provided the proof for Observation 1 (see Appendix A); however, we believe that it is also necessary to examine experimental results for an intuitive understanding. Under the fixed number of synthetic instances, we investigate the performance changes of input-sized parameterization while imposing constraints on the feasible space of synthetic instances. We consider two types of constraints: 1) dimension masking and 2) value clipping. As shown in Tables 14 and 15, when the feasible space becomes smaller (i.e., as the restrictions are enforced more strongly), the dataset distillation performances decrease. These results serve as direct evidence of observations in theoretical analysis.

Table 14: Dimension masking

Masking (%)	0	25	50	75
Accuracy (%)	$51.4 \pm 2.3$	$50.2 \pm 1.7$	$48.7 \pm 3.2$	$41.1 \pm 1.9$

Table 15: Value clipping into  $[-R, R]$

$R$	$\infty$	2.0	1.0	0.5
Accuracy (%)	$51.4 \pm 2.3$	$43.2 \pm 0.9$	$25.8 \pm 2.8$	$21.0 \pm 1.3$

### C.6 ROBUSTNESS TO CORRUPTION

We further investigate the robustness against corruption of DDiF in the trained synthetic datasets by evaluating on ImageNet-Subset-C. This subset is designed specifically to assess robustness across varying corruption types and severity levels. We report the average test accuracies over 15 corruption types, each evaluated across 5 levels of severity, for each class in ImageNet-Subset-C. Table 16 summarizes the performance of DDiF and baseline models. The results in Table 16 clearly demonstrate that DDiF has the same substantial robustness to resolution change and corruption.

Table 16: Test accuracies (%) on ImageNet-Subset-C under IPC=1. Test accuracy on ImageSquawk-C of TM is not reported on previous works because the default size of the dataset does not fit to  $128 \times 128$

Method	ImageNette-C	ImageWoof-C	ImageFruit-C	ImageYellow-C	ImageMeow-C	ImageSquawk-C
TM	$38.0 \pm 1.6$	$23.8 \pm 1.0$	$22.7 \pm 1.1$	$35.6 \pm 1.7$	$23.2 \pm 1.1$	-
IDC	$34.5 \pm 0.6$	$18.7 \pm 0.4$	$28.5 \pm 0.9$	$36.8 \pm 1.4$	$22.2 \pm 1.2$	$26.8 \pm 0.5$
FreD	$51.2 \pm 0.6$	$31.0 \pm 0.9$	$32.3 \pm 1.4$	$48.2 \pm 1.0$	$30.3 \pm 0.3$	$45.9 \pm 0.6$
DDiF	$54.5 \pm 0.6$	$34.0 \pm 0.4$	$36.6 \pm 0.4$	$47.2 \pm 0.7$	$30.3 \pm 0.8$	$53.8 \pm 0.5$

### C.7 TIME COMPLEXITY

As mentioned earlier, the neural field takes coordinates as input and produces quantities as output. This distinct characteristic of the neural field offers the advantage of being resolution-invariant but may raise concerns regarding the decoding process time. We admit that the decoding time of DDiF indeed increases as resolution grows due to the need to forward a larger number of coordinates through the neural field. To address this concern, we conducted an experiment measuring the wall-clock time for decoding a single instance.

Table 17 shows the results with 128 image resolution demonstrate that while the time cost of DDiF is larger than methods relying on non-parameterized decoding functions, such as IDC and FreD, it remains comparable to methods that use parameterized decoding functions, such as HaBa and SPEED, and exhibits a lower time cost than methods that utilize pre-trained generative models, such as GLaD. As the resolution increases to 256, the decoding time of DDiF also increases and it is slightly larger than non-parameterized decoding functions. In conclusion, although the decoding process time of DDiF increases as the resolution increases, it does not differ significantly from that of conventional parameterization methods. We attribute this to 1) the use of a small neural network structure for the synthetic neural field and 2) the full-batch forwarding of the coordinate set in the implementation.

Table 17: Wall-clock time (ms) of the decoding process for a single synthetic instance. "ms" indicates the millisecond.

	128×128	256×256
Vanilla	0.31	0.31
IDC	0.40	2.83
FreD	0.46	1.20
HaBa	2.81	–
SPEED	2.20	–
GLaD	31.33	–
DDiF	2.49	3.25

## C.8 FULL TABLE WITH STANDARD DEVIATION AND ADDITIONAL VISUALIZATION

For improved layout, we have positioned full tables with standard deviation and additional example figures at the end of the paper. Please refer to Table 18 for ImageNet-Subset (128 × 128) under IPC=1; Table 19 for ImageNet-Subset (256 × 256) under IPC=1; Table 20 for ImageNet-Subset under (128 × 128) IPC=10; Table 21 for cross-architecture; and Table 22 for robustness to the loss. In addition, please refer to Figures 9 to 14 for Image domain; Figure 15 for 3D domain; and Figure 16 for video domain.

## D ADDITIONAL DISCUSSIONS

### D.1 MORE COMPARISON WITH FRED

As seen in Eqs. (5) and (6), the functions represented by DDiF and FreD are similar, both being the sum of cosine functions. Therefore, we can perform a term-by-term comparison of both equations.

- DDiF enables to change the amplitudes  $A_{k,i}$ , frequencies  $\omega_k$ , phases  $\varphi'_{k,i}$ , and shift  $b^{(0)}$ , while FreD only allows the amplitudes. It indicates that DDiF has higher representation ability than FreD.
- Although FreD is a finite sum of cosine functions with a fixed frequency, DDiF represents an infinite sum of cosine functions with tunable frequency. It means that DDiF can cover a wide range of frequencies from low to high by selecting various  $k$ . According to the empirical findings in Wang et al. (2020a), it has been demonstrated that datasets with a larger number of frequencies exhibit improved generalization performance.
- DDiF is a continuous function, whereas FreD is a discrete function. Due to this characteristic, FreD cannot encode information for coordinates that were not provided during the distillation stage. In contrast, since DDiF operates over a continuous domain, it inherently stores information for coordinates that were not supplied during the distillation stage.

In summary, DDiF has a more flexible and expressive function than FreD.

### D.2 DISCUSSION ABOUT THEORETICAL ANALYSIS

We provide the theoretical analysis to propose a framework for comparing parameterization methods, which have traditionally been evaluated solely based on performance, through expressive-

ness—specifically, the size of the feasible space. Observations 1 and 2 indicate that a larger feasible space for decoded synthetic instances through input-sized parameterization or parameterization results in lower dataset distillation loss. Remark 1 and Theorem 2 to show that DDiF has higher expressiveness than prior work (FreD). We believe that the proposed theoretical analysis framework will serve as a cornerstone for future theoretical comparisons of parameterization methods in dataset distillation areas.

However, this theoretical analysis still has room for further improvement. In the theoretical analysis, we consider the fixed number of decoded instance scenarios, not the fixed storage budget. We experimentally demonstrated the superiority of DDiF not only in the fixed storage budget scenario but also in various situations with a fixed number of decoded instances (see Figure 5 and Appendix C.3). In particular, we demonstrated the high efficiency and expressiveness of DDiF by showing that, even with a smaller budget in the fixed number of decoded instances, DDiF achieved competitive or higher performance than previous studies. Even though the proposed theoretical analysis in this paper is experimentally verified through extensive results, this framework has the limitation of not primarily focusing on the fixed storage budget scenario, which is the widely used setting of dataset distillation. We believe that constructing a theoretical framework to compare the expressiveness of parameterization methods under a fixed storage budget is necessary, and the proposed theoretical analysis in this paper can serve as a foundational background for such efforts.

### D.3 COMPARISON WITH DiM

As mentioned in Section 3, the synthetic function has several possible forms. DiM (Wang et al., 2023) employs a probability density function as a synthetic function and utilizes a deep generative model to parameterize it. Specifically, the decoded synthetic instance of DiM is the sampled output of a deep generative model by inputting random noise:

$$\min_{\phi} \mathcal{L}(\mathcal{T}, \mathcal{S}) \text{ where } \mathcal{S} = g_{\phi}(\mathcal{Z}), \mathcal{Z} \sim \mathcal{N}(0, I)$$

DDiF and DiM have in common that they store the distilled information in the synthetic function and only store the parameters of the function without any additional codes. However, there are several structural differences.

- The output of DiM still depends on the data dimension. As aforementioned, this type of decoding function requires a more complicated structure and storage budget as the data dimension grows larger. Actually, DiM has not been extensively tested on high-resolution datasets. On the contrary, DDiF stores information regardless of data dimension, which indicates broader applicability across various resolutions.
- The decoding process of DiM is stochastic. Due to the stochasticity, DiM can sample the diverse decoded synthetic instances and save the redeployment cost. However, at the same time, DiM carries the risk of generating less informative synthetic instances. Consequently, it leads to instability in training on downstream tasks. Furthermore, DiM might suffer from redundant sampling due to mode collapse, a well-known issue of the generative model. Meanwhile, since the decoding process of DDiF is deterministic, DDiF has an advantage in stability.

## E LIMITATION

**Less efficiency on Low-dimensional datasets.** One of the main ideas in parameterization is expanding the trade-off curve between quantity and quality: reducing the utilized budget of each synthetic instance, while maximally preserving the expressiveness of it. While DDiF has an extensive feature coverage theoretically and shows competitive performances with previous studies experimentally, it might be less efficient for some low-dimensional datasets due to the structural features of the neural field. This is because, given the low-dimensional instance, it can be difficult to design a neural field that is sufficiently expressive with fewer parameters. In spite of this issue, we repeatedly highlight that DDiF has a significant performance improvement on high-dimensional datasets. Furthermore, we speculate that a deeper analysis of the neural field structure could be an interesting direction for future research in dataset distillation.

**Individual Parameterization.** Several studies have claimed that storing intra- and inter-class information in a shared component is effective. From this perspective, our proposed method, which has a one-to-one correspondence between synthetic instances and synthetic neural fields, does not have a component to store shared information. We believe that it can be extended by adding modulation or conditional code, and this paper may serve as a good starting point.

Table 18: Test accuracies (%) on ImageNet-Subset ( $128 \times 128$ ) with regard to various dataset parameterization methods under IPC=1.

	Subset	Nette	Woof	Fruit	Yellow	Meow	Squawk
Input sized	TM	51.4 $\pm$ 2.3	29.7 $\pm$ 0.9	28.8 $\pm$ 1.2	47.5 $\pm$ 1.5	33.3 $\pm$ 0.7	41.0 $\pm$ 1.5
	FRePo	48.1 $\pm$ 0.7	29.7 $\pm$ 0.6	-	-	-	-
Static	IDC	61.4 $\pm$ 1.0	34.5 $\pm$ 1.1	38.0 $\pm$ 1.1	56.5 $\pm$ 1.8	39.5 $\pm$ 1.5	50.2 $\pm$ 1.5
	FreD	66.8 $\pm$ 0.4	38.3 $\pm$ 1.5	43.7 $\pm$ 1.6	63.2 $\pm$ 1.0	43.2 $\pm$ 0.8	57.0 $\pm$ 0.8
Parameterized	HaBa	51.9 $\pm$ 1.7	32.4 $\pm$ 0.7	34.7 $\pm$ 1.1	50.4 $\pm$ 1.6	36.9 $\pm$ 0.9	41.9 $\pm$ 1.4
	SPEED	66.9 $\pm$ 0.7	38.0 $\pm$ 0.9	43.4 $\pm$ 0.6	62.6 $\pm$ 1.3	43.6 $\pm$ 0.7	60.9 $\pm$ 1.0
	NSD	68.6 $\pm$ 0.8	35.2 $\pm$ 0.4	39.8 $\pm$ 0.2	61.0 $\pm$ 0.5	45.2 $\pm$ 0.1	52.9 $\pm$ 0.7
DGM Prior	GLaD	38.7 $\pm$ 1.6	23.4 $\pm$ 1.1	23.1 $\pm$ 0.4	-	26.0 $\pm$ 1.1	35.8 $\pm$ 1.4
	H-GLaD	45.4 $\pm$ 1.1	28.3 $\pm$ 0.2	25.6 $\pm$ 0.7	-	29.6 $\pm$ 1.0	39.7 $\pm$ 0.8
Function	DDiF	72.0 $\pm$ 0.9	42.9 $\pm$ 0.7	48.2 $\pm$ 1.2	69.0 $\pm$ 0.8	47.4 $\pm$ 1.3	67.0 $\pm$ 1.3

Table 19: Test accuracies (%) on ImageNet-Subset ( $256 \times 256$ ) resolution under IPC=1. "DM" in this table means it utilize distribution matching objective for distillation.

	Subset	Nette	Woof	Fruit	Yellow	Meow	Squawk
Input sized	DM	32.1	20.0	19.5	33.4	21.2	27.6
Static	IDC	53.7 $\pm$ 1.2	30.2 $\pm$ 1.5	33.1 $\pm$ 1.5	52.2 $\pm$ 1.4	34.6 $\pm$ 1.8	47.0 $\pm$ 1.5
	FreD	54.2 $\pm$ 1.1	31.2 $\pm$ 0.9	32.5 $\pm$ 1.9	49.1 $\pm$ 0.4	34.0 $\pm$ 1.2	43.1 $\pm$ 1.5
Parameterized	SPEED (DM)	57.7 $\pm$ 0.9	-	-	-	-	-
DGM Prior	LatentDM (DM)	56.1	28.0	30.7	-	36.3	47.1
Function	DDiF	67.8 $\pm$ 1.0	39.6 $\pm$ 1.6	43.2 $\pm$ 1.7	63.1 $\pm$ 0.8	44.8 $\pm$ 1.1	67.0 $\pm$ 0.9

Table 20: Test accuracies (%) on ImageNet-Subset ( $128 \times 128$ ) with regard to various dataset parameterization methods under IPC=10.

	Subset	Nette	Woof	Fruit	Yellow	Meow	Squawk
Input sized	TM	63.0 $\pm$ 1.3	35.8 $\pm$ 1.8	40.3 $\pm$ 1.3	60.0 $\pm$ 1.5	40.4 $\pm$ 2.2	52.3 $\pm$ 1.0
	FRePo	66.5 $\pm$ 0.8	42.2 $\pm$ 0.9	-	-	-	-
Static	IDC	70.8 $\pm$ 0.5	39.8 $\pm$ 0.9	46.4 $\pm$ 1.4	68.7 $\pm$ 0.8	47.9 $\pm$ 1.4	65.4 $\pm$ 1.2
	FreD	72.0 $\pm$ 0.8	41.3 $\pm$ 1.2	47.0 $\pm$ 1.1	69.2 $\pm$ 0.6	48.6 $\pm$ 0.4	67.3 $\pm$ 0.8
Parameterized	HaBa	64.7 $\pm$ 1.6	38.6 $\pm$ 1.3	42.5 $\pm$ 1.6	63.0 $\pm$ 1.6	42.9 $\pm$ 0.9	56.8 $\pm$ 1.0
	SPEED	72.9 $\pm$ 1.5	44.1 $\pm$ 1.4	50.0 $\pm$ 0.8	70.5 $\pm$ 1.5	52.0 $\pm$ 1.3	71.8 $\pm$ 1.3
Function	DDiF	74.6 $\pm$ 0.7	44.9 $\pm$ 0.5	49.8 $\pm$ 0.8	70.5 $\pm$ 1.8	50.6 $\pm$ 1.1	72.3 $\pm$ 1.3

Table 21: Test accuracies (%) on Cross Architecture networks with ImageNet Subsets ( $128 \times 128$ ), IPC=1

Test Net	Method	Dataset					
		Nette	Woof	Fruit	Yellow	Meow	Squawk
AlexNet	TM	13.2 $\pm$ 0.6	10.0 $\pm$ 0.0	10.0 $\pm$ 0.0	11.0 $\pm$ 0.2	9.8 $\pm$ 0.0	–
	IDC	17.4 $\pm$ 0.9	16.5 $\pm$ 0.7	17.9 $\pm$ 0.7	20.6 $\pm$ 0.9	16.8 $\pm$ 0.5	20.7 $\pm$ 1.0
	FreD	35.7 $\pm$ 0.4	23.9 $\pm$ 0.7	15.8 $\pm$ 0.7	19.8 $\pm$ 1.2	14.4 $\pm$ 0.5	36.3 $\pm$ 0.3
	DDiF	<b>60.7</b> $\pm$ 2.3	<b>36.4</b> $\pm$ 2.3	<b>41.8</b> $\pm$ 0.6	<b>56.2</b> $\pm$ 0.8	<b>40.3</b> $\pm$ 1.9	<b>60.5</b> $\pm$ 0.4
VGG11	TM	17.4 $\pm$ 2.1	12.6 $\pm$ 1.8	11.8 $\pm$ 1.0	16.9 $\pm$ 1.1	13.8 $\pm$ 1.3	–
	IDC	19.6 $\pm$ 1.5	16.0 $\pm$ 2.1	13.8 $\pm$ 1.3	16.8 $\pm$ 3.5	13.1 $\pm$ 2.0	19.1 $\pm$ 1.2
	FreD	21.8 $\pm$ 2.9	17.1 $\pm$ 1.7	12.6 $\pm$ 2.6	18.2 $\pm$ 1.1	13.2 $\pm$ 1.9	18.6 $\pm$ 2.3
	DDiF	<b>53.6</b> $\pm$ 1.5	<b>29.9</b> $\pm$ 1.9	<b>33.8</b> $\pm$ 1.9	<b>44.2</b> $\pm$ 1.7	<b>32.0</b> $\pm$ 1.8	<b>37.9</b> $\pm$ 1.5
ResNet18	TM	34.9 $\pm$ 2.3	20.7 $\pm$ 1.0	23.1 $\pm$ 1.5	43.4 $\pm$ 1.1	22.8 $\pm$ 2.2	–
	IDC	43.6 $\pm$ 1.3	23.2 $\pm$ 0.8	32.9 $\pm$ 2.8	44.2 $\pm$ 3.5	28.2 $\pm$ 0.5	47.8 $\pm$ 1.9
	FreD	48.8 $\pm$ 1.8	28.4 $\pm$ 0.6	34.0 $\pm$ 1.9	49.3 $\pm$ 1.1	29.0 $\pm$ 1.8	50.2 $\pm$ 0.8
	DDiF	<b>63.8</b> $\pm$ 1.8	<b>37.5</b> $\pm$ 1.9	<b>42.0</b> $\pm$ 1.9	<b>55.9</b> $\pm$ 1.0	<b>35.8</b> $\pm$ 1.8	<b>62.6</b> $\pm$ 1.5
ViT	TM	22.6 $\pm$ 1.1	15.9 $\pm$ 0.4	23.3 $\pm$ 0.4	18.1 $\pm$ 1.3	18.6 $\pm$ 0.9	–
	IDC	31.0 $\pm$ 0.6	22.4 $\pm$ 0.8	31.1 $\pm$ 0.8	30.3 $\pm$ 0.6	21.4 $\pm$ 0.7	32.2 $\pm$ 1.2
	FreD	38.4 $\pm$ 0.7	25.4 $\pm$ 1.7	31.9 $\pm$ 1.4	37.6 $\pm$ 2.0	19.7 $\pm$ 0.8	44.4 $\pm$ 1.0
	DDiF	<b>59.0</b> $\pm$ 0.4	<b>32.8</b> $\pm$ 0.8	<b>39.4</b> $\pm$ 0.8	<b>47.9</b> $\pm$ 0.6	<b>27.0</b> $\pm$ 0.6	<b>54.8</b> $\pm$ 1.1

Table 22: Compatibility on different dataset distillation loss with ImageNet Subsets ( $128 \times 128$ ), IPC=1

$\mathcal{L}_{DD}$	Method	Dataset				
		Nette	Woof	Fruit	Meow	Squawk
DC	Vanilla	34.2 $\pm$ 1.7	22.5 $\pm$ 1.0	21.0 $\pm$ 0.9	22.0 $\pm$ 0.6	32.0 $\pm$ 1.5
	IDC	45.4 $\pm$ 0.7	25.5 $\pm$ 0.7	26.8 $\pm$ 0.4	25.3 $\pm$ 0.6	34.6 $\pm$ 0.5
	FreD	49.1 $\pm$ 0.8	26.1 $\pm$ 1.1	30.0 $\pm$ 0.7	28.7 $\pm$ 1.0	39.7 $\pm$ 0.7
	GLaD	35.4 $\pm$ 1.2	22.3 $\pm$ 1.1	20.7 $\pm$ 1.1	22.6 $\pm$ 0.8	33.8 $\pm$ 0.9
	H-GLaD	36.9 $\pm$ 0.8	24.0 $\pm$ 0.8	22.4 $\pm$ 1.1	24.1 $\pm$ 0.9	35.3 $\pm$ 1.0
	DDiF	<b>61.2</b> $\pm$ 1.0	<b>35.2</b> $\pm$ 1.7	<b>37.8</b> $\pm$ 1.1	<b>39.1</b> $\pm$ 1.3	<b>54.3</b> $\pm$ 1.0
DM	Vanilla	30.4 $\pm$ 2.7	20.7 $\pm$ 1.0	20.4 $\pm$ 1.9	20.1 $\pm$ 1.2	26.6 $\pm$ 2.6
	IDC	48.3 $\pm$ 1.3	27.0 $\pm$ 1.0	29.9 $\pm$ 0.7	30.5 $\pm$ 1.0	38.8 $\pm$ 1.4
	FreD	56.2 $\pm$ 1.0	31.0 $\pm$ 1.2	33.4 $\pm$ 0.5	33.3 $\pm$ 0.6	42.7 $\pm$ 0.8
	GLaD	32.2 $\pm$ 1.7	21.2 $\pm$ 1.5	21.8 $\pm$ 1.8	22.3 $\pm$ 1.6	27.6 $\pm$ 1.9
	H-GLaD	34.8 $\pm$ 1.0	23.9 $\pm$ 1.9	24.4 $\pm$ 2.1	24.2 $\pm$ 1.1	29.5 $\pm$ 1.5
	DDiF	<b>69.2</b> $\pm$ 1.0	<b>42.0</b> $\pm$ 0.4	<b>45.3</b> $\pm$ 1.8	<b>45.8</b> $\pm$ 1.1	<b>64.6</b> $\pm$ 1.1

1458  
 1459  
 1460  
 1461  
 1462  
 1463  
 1464  
 1465  
 1466  
 1467  
 1468  
 1469  
 1470  
 1471  
 1472  
 1473  
 1474  
 1475  
 1476  
 1477  
 1478  
 1479  
 1480  
 1481  
 1482  
 1483  
 1484  
 1485  
 1486  
 1487  
 1488  
 1489  
 1490  
 1491  
 1492  
 1493  
 1494  
 1495  
 1496  
 1497  
 1498  
 1499  
 1500  
 1501  
 1502  
 1503  
 1504  
 1505  
 1506  
 1507  
 1508  
 1509  
 1510  
 1511

Table 23: Configuration of hyperparameters.

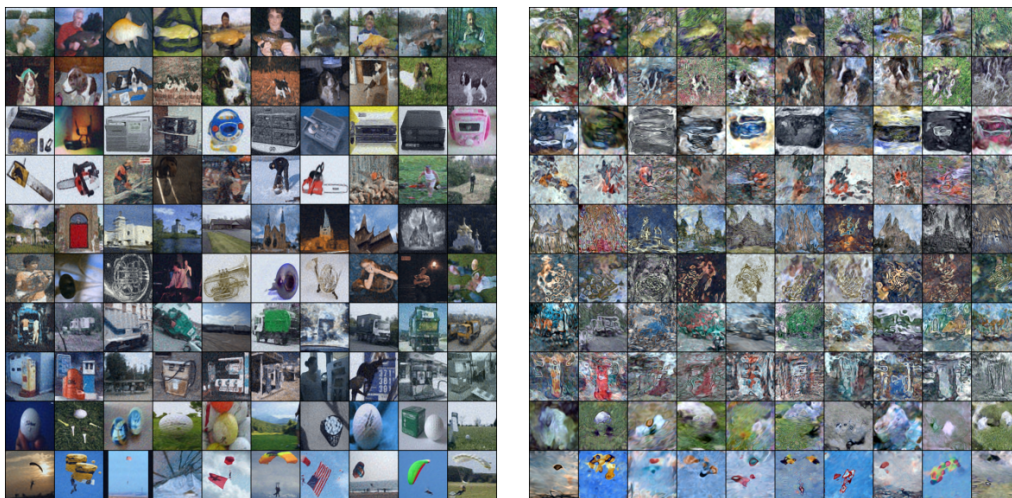
(a) Gradient matching (DC)				(b) Distribution matching (DM)			
Dataset	IPC	Synthetic batch size	Learning rate (Neural field)	Dataset	IPC	Synthetic batch size	Learning rate (Neural field)
ImageNet-Subset (128 × 128)	1	-	$5 \times 10^{-5}$	ImageNet-Subset (128 × 128)	1	-	$5 \times 10^{-5}$
Mini Speech Commands	10	64	$10^{-5}$	ImageNet-Subset (256 × 256)	1	-	$10^{-5}$
	20	64	$10^{-4}$				
ModelNet	1	-	$10^{-4}$	MiniUCF	1	-	$10^{-4}$
ShapeNet	1	-	$10^{-4}$		5	-	$10^{-4}$
				ModelNet	1	-	$10^{-4}$
				ShapeNet	1	-	$10^{-4}$

(c) Trajectory matching (TM)								
Dataset	IPC	Synthetic steps	Expert epochs	Max start epoch	Synthetic batch size	Learning rate (Neural field)	Learning rate (Step size)	Learning rate (Teacher)
CIFAR-10	1	60	2	10	74	$10^{-3}$	$10^{-5}$	$10^{-2}$
	10	60	2	10	256	$10^{-3}$	$10^{-5}$	$10^{-2}$
	50	60	2	40	235	$10^{-4}$	$10^{-5}$	$10^{-2}$
CIFAR-100	1	60	2	40	170	$10^{-3}$	$10^{-5}$	$10^{-2}$
	10	60	2	40	230	$10^{-3}$	$10^{-5}$	$10^{-2}$
	50	60	2	40	276	$10^{-3}$	$10^{-5}$	$10^{-2}$
ImageNet-Subset (128 × 128)	1	20	2	10	102	$10^{-4}$	$10^{-6}$	$10^{-2}$
	10	40	2	20	30	$10^{-4}$	$10^{-5}$	$10^{-2}$
	50	40	2	20	30	$10^{-4}$	$10^{-5}$	$10^{-2}$



1512  
1513  
1514  
1515  
1516  
1517  
1518  
1519  
1520  
1521  
1522  
1523  
1524  
1525  
1526  
1527  
1528  
1529  
1530  
1531  
1532  
1533  
1534  
1535  
1536  
1537  
1538  
1539  
1540  
1541  
1542  
1543  
1544  
1545  
1546  
1547  
1548  
1549  
1550  
1551  
1552  
1553  
1554  
1555  
1556  
1557  
1558  
1559  
1560  
1561  
1562  
1563  
1564  
1565



(a) Initialization

(b) Synthesized

Figure 9: (a) Warm-up initialized images on ImageNette with DDiF, (b) Best-performed synthetic dataset represented by DDiF. We visualize the first 10 images, while DDiF constructs 51 images per class under the same budget.



(a) Initialization

(b) Synthesized

Figure 10: (a) Warm-up initialized images on ImageWoof with DDiF, (b) Best-performed synthetic dataset represented by DDiF. We visualize the first 10 images, while DDiF constructs 51 images per class under the same budget.



1566  
 1567  
 1568  
 1569  
 1570  
 1571  
 1572  
 1573  
 1574  
 1575  
 1576  
 1577  
 1578  
 1579  
 1580  
 1581  
 1582  
 1583  
 1584  
 1585  
 1586  
 1587  
 1588  
 1589  
 1590  
 1591  
 1592  
 1593  
 1594  
 1595  
 1596  
 1597  
 1598  
 1599  
 1600  
 1601  
 1602  
 1603  
 1604  
 1605  
 1606  
 1607  
 1608  
 1609  
 1610  
 1611  
 1612  
 1613  
 1614  
 1615  
 1616  
 1617  
 1618  
 1619

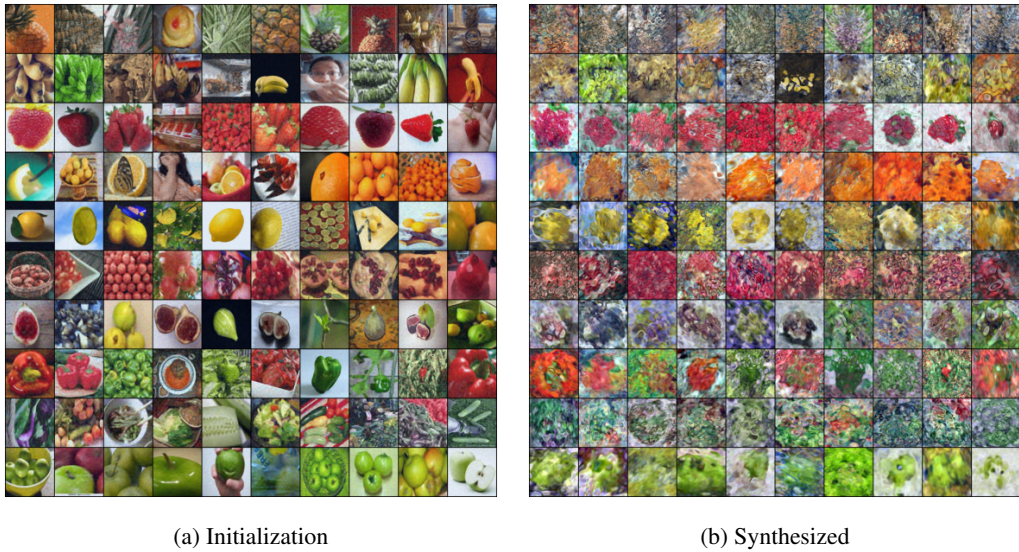


Figure 11: (a) Warm-up initialized images on ImageFruit with DDiF, (b) Best-performed synthetic dataset represented by DDiF. We visualize the first 10 images, while DDiF constructs 51 images per class under the same budget.

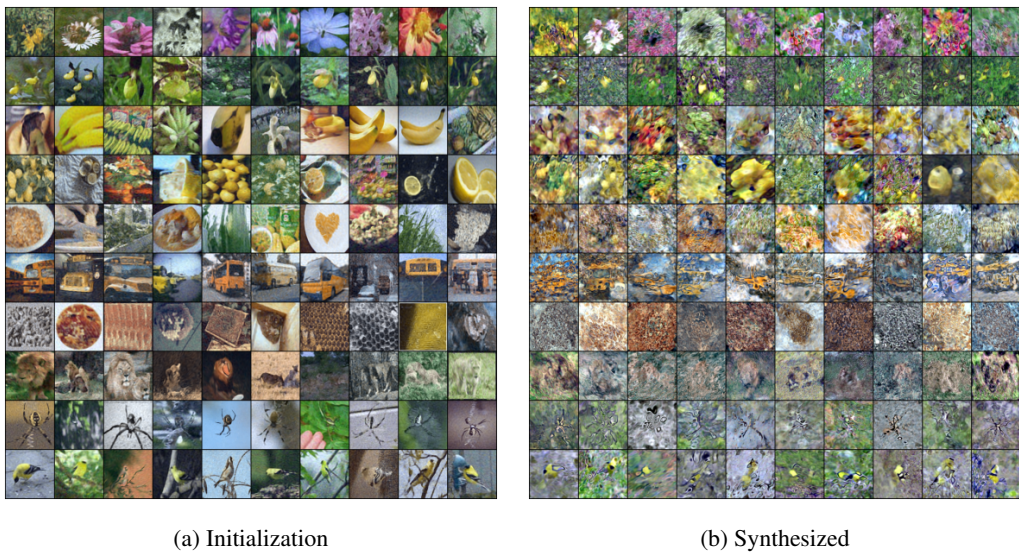
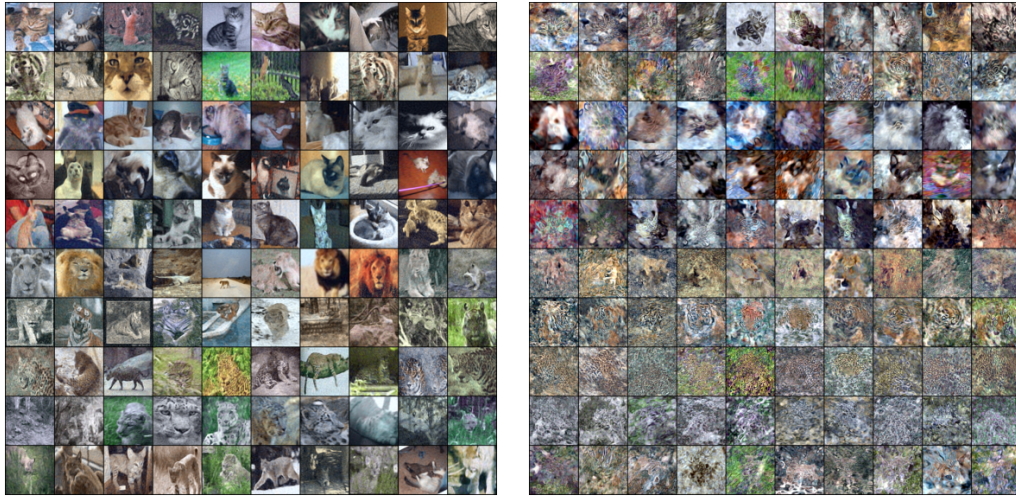


Figure 12: (a) Warm-up initialized images on ImageYellow with DDiF, (b) Best-performed synthetic dataset represented by DDiF. We visualize the first 10 images, while DDiF constructs 51 images per class under the same budget.



1620  
1621  
1622  
1623  
1624  
1625  
1626  
1627  
1628  
1629  
1630  
1631  
1632  
1633  
1634  
1635  
1636  
1637  
1638  
1639  
1640

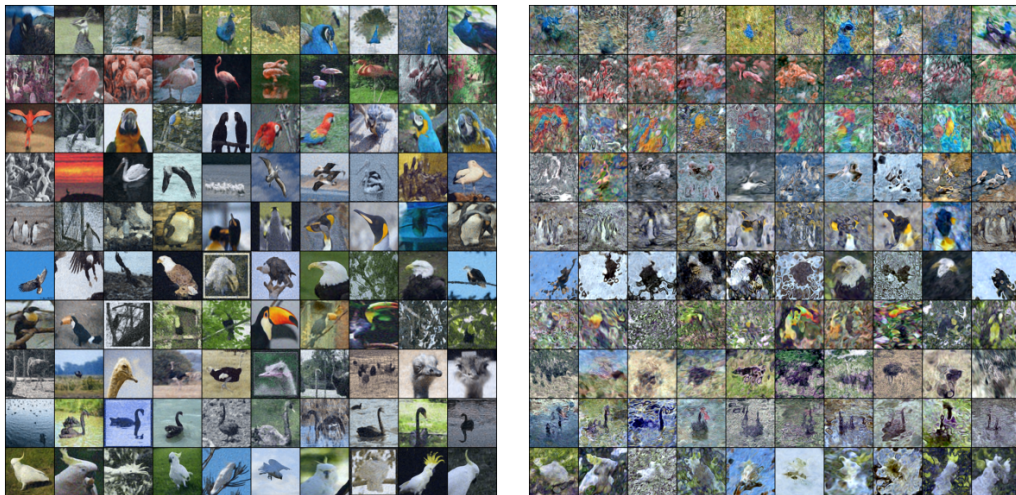


(a) Initialization

(b) Synthesized

1641 Figure 13: (a) Warm-up initialized images on ImageMeow with DDiF, (b) Best-performed synthetic  
1642 dataset represented by DDiF. We visualize the first 10 images, while DDiF constructs 51 images per  
1643 class under the same budget.  
1644

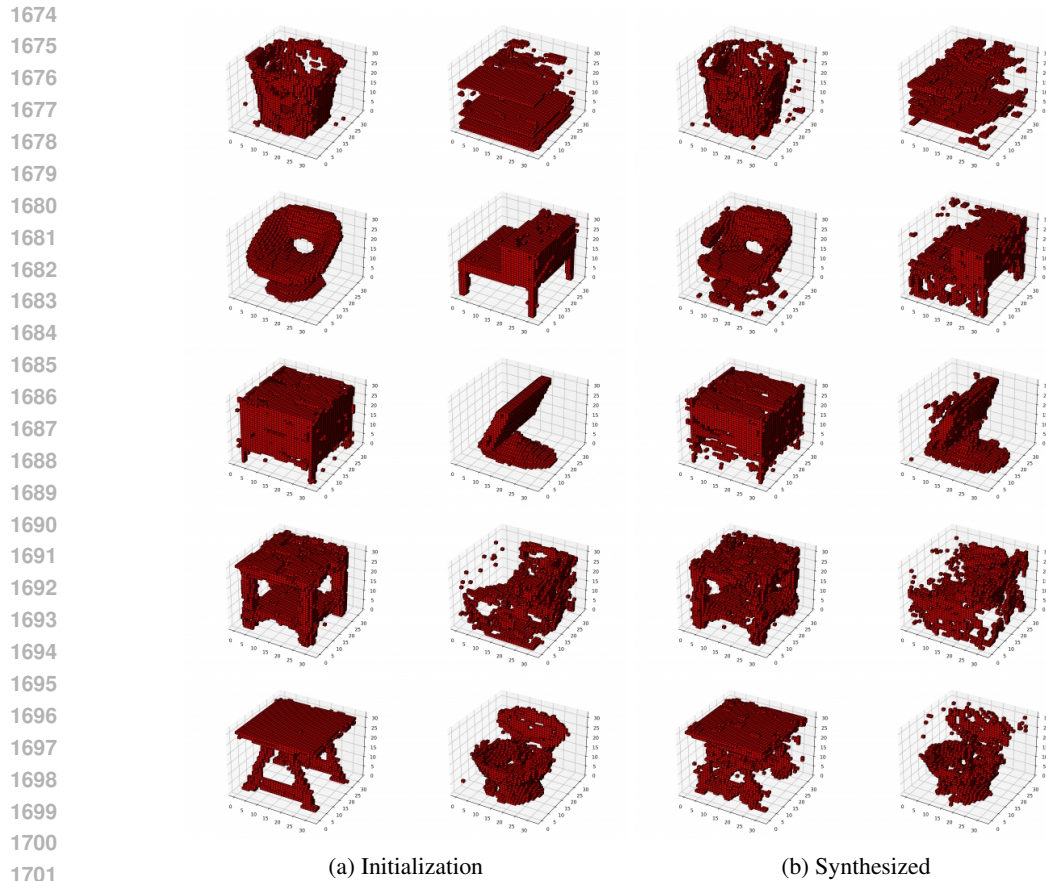
1645  
1646  
1647  
1648  
1649  
1650  
1651  
1652  
1653  
1654  
1655  
1656  
1657  
1658  
1659  
1660  
1661  
1662  
1663  
1664  
1665  
1666  
1667  
1668



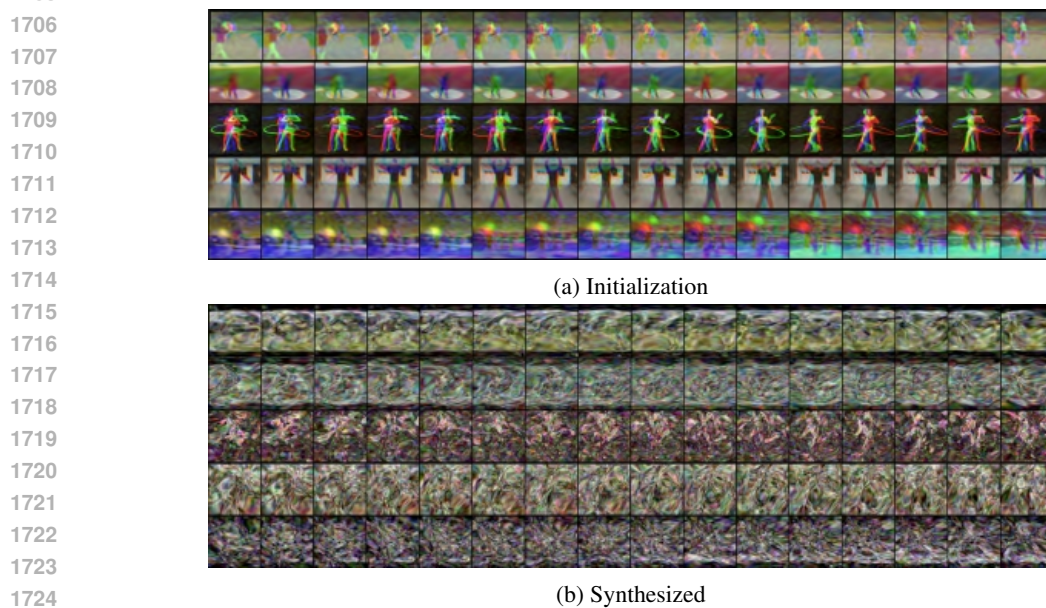
(a) Initialization

(b) Synthesized

1669 Figure 14: (a) Warm-up initialized images on ImageSquawk with DDiF, (b) Best-performed synthetic  
1670 dataset represented by DDiF. We visualize the first 10 images, while DDiF constructs 51  
1671 images per class under the same budget.  
1672  
1673



1702 Figure 15: (a) Warm-up initialization images and (b) Synthesized images of ModelNet-10. Start  
 1703 reading labels from the top and continue to right: 1) bathtub, 2) bed, 3) chair, 4) desk, 5) dresser, 6)  
 1704 monitor, 7) nightstand, 8) sofa, 9) table, 10) toilet



1725 Figure 16: (a) Warm-up initialization images and (b) Synthesized images of MiniUCF. Start reading  
 1726 labels from the top: 1) FrisbeeCatch, 2) HammerThrow, 3) HulaHoop, 4) JumpingJack, 5) Parallel-  
 1727 Bars