

CONCORD: CONCEPT-INFORMED DIFFUSION FOR DATASET DISTILLATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Dataset distillation has witnessed significant progress in synthesizing small-scale datasets that encapsulate rich information from large-scale original ones. Particularly, methods based on generative priors show promising performance, while maintaining computational efficiency and cross-architecture generalization. However, the generation process lacks explicit controllability for each sample. Previous distillation methods primarily match the real distribution from the perspective of the entire dataset, whereas overlooking conceptual completeness at the instance level. This oversight can result in missing or incorrectly represented object details and compromised dataset quality. To this end, we propose to incorporate the conceptual understanding of large language models (LLMs) to perform a CONCEPT-INFORMED Diffusion process for dataset distillation, in short as CONCORD. Specifically, distinguishable and fine-grained concepts are retrieved based on category labels to explicitly inform the denoising process and refine essential object details. By integrating these concepts, the proposed method significantly enhances both the controllability and interpretability of the distilled image generation, without relying on pre-trained classifiers. We demonstrate the efficacy of CONCORD by achieving state-of-the-art performance on ImageNet-1K and its subsets. It further advances the practical application of dataset distillation methods. The code implementation is attached in the supplementary material.

1 INTRODUCTION

In the current digital era, vast volumes of data are produced and disseminated across online platforms on a daily basis. The abundance of data boosts the training of robust neural network models, which often outperform human experts in a variety of domains (He et al., 2016; Dosovitskiy et al., 2022; Brown et al., 2020; Deng et al., 2009; Devlin et al., 2018). However, the heavy dependence on data also causes unbearable burden on the storage space and computational consumption. Strong neural networks often demand days or even months of training on high-capacity hardware, and this issue is exacerbated for more complex foundation models (Radford et al., 2021; He et al., 2022; Touvron et al., 2023; Bai et al., 2023). While pre-trained models are mostly available for general use, developing new networks from scratch remains necessary for certain specialized domains, and would be particularly challenging for resource-constrained research teams. In this context, Dataset Distillation (DD) emerges as a solution to condense rich information from original large-scale datasets into much smaller surrogate datasets (Wang et al., 2018; Zhao et al., 2021; Yu et al., 2023; Sachdeva & McAuley, 2023). With substantially reduced training time, the surrogate datasets aim to restore the performance levels of the original data for practical applications.

Typical DD methods incorporate meta-learning or metric matching to condense rich information into surrogate sets, and have achieved considerable performance on various benchmarks (Wang et al., 2018; Zhao et al., 2021; Nguyen et al., 2021b; Loo et al., 2022; Kim et al., 2022b). However, the distillation phase itself often demands even longer time compared with the training process on the original dataset (Cui et al., 2023; Sun et al., 2024). It would still be impractical for individual researchers to perform distillation on personalized datasets. Besides, these methods are easily biased towards the architecture adopted in the distillation phase, necessitating specialized designs to mitigate cross-architecture generalization challenges (Zhou et al., 2023; Wang et al., 2023a).

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

Figure 1: Comparison on example generated images with and without our proposed CONCORD method.^c indicates that CONCORD is applied. Incorporating rich knowledge from LLMs, CONCORD re nes instance-level conceptual completeness, and enhances the overall dataset quality.

Recently, a series of methods integrate generative models to synthesize training data (Cazenavette et al., 2023; Gu et al., 2024a; Su et al., 2024; Moser et al., 2024). The pre-acquired generative prior within these models contributes to better cross-architecture generalization as well as signi cantly lower distillation consumption. While the synthetic images yield state-of-the-art performance, the distillation process lacks explicit controllability for each sample. Most existing approaches condense information by mimicking the distribution of real data at the dataset level. On the one hand, the lack of instance-level control might result in conceptual incompleteness, where essential object details may be missing or inaccurately represented in the generated images. Due to the constrained storage budget typical of DD benchmarks, this information loss cannot be suf ciently compensated. On the other hand, the distribution imitation is dif cult to interpret, as the dataset quality can only be measured indirectly through training performance. It also raises a question: merely imitating the real distribution suf ce for generating effective surrogate datasets?

To this end, we intend to explicitly enhance instance-level conceptual completeness during the diffusion process with the assistance of large language models (LLMs). LLMs have obtained extensive conceptual understanding across a variety of objects, which can be utilized to facilitate examining and re ning the defects and incorrect concept representations in the images. Our method involves initially retrieving distinguishable concepts speci c to the target categories, and subsequently performing the CONCEPT-infORMed Diffusion inference (in short as CONCORD) to supplement missing or incorrect details. The approach offers several advantages. Firstly, the ne-grained control exerted by the retrieved concepts allows for more accurate re nement of object details, which also enables higher levels of personalization. Secondly, the concepts provide explicit explanations why the generated images are better suited for model training. Additionally, we employ concepts from similar categories to construct negative samples, thereby ensuring more accurate and stabilized control over the generation process. By prioritizing the enhancement of crucial concepts in addition to distribution imitating, our method generates more effective distilled data for training models.

As shown in Fig. 1, the samples generated by Minimax (Gu et al., 2024a) often fail to include complete and correct concepts for their respective categories, such as unrealistic back legs of the beagle and a missing wing in the cabbage butterfly image. With the assistance of rich knowledge from LLMs, the proposed CONCORD method signi cantly improves the conceptual completeness, and reduces image defects. The proposed CONCORD method can be plugged into any diffusion-based generative pipelines for dataset distillation. We conduct extensive experiments on both Minimax and Stable Diffusion baselines (Ramesh et al., 2022) to illustrate the superiority of CONCORD, which achieves state-of-the-art performance on the full ImageNet-1K dataset and its subsets. Notably, the method only incorporates descriptive concepts to inform the diffusion process, eliminating the dependence on pre-trained classi ers. It reduces the required computational consumption, and thereby enhances the practicality of our approach for broader application possibilities.

2 RELATED WORK

Dataset Distillation Aiming at reducing the demanded storage and computational consumption for training neural networks, dataset distillation (DD) has been increasingly investigated in recent

years (Yu et al., 2023; Sachdeva & McAuley, 2023) and achieved broad applications (Gu et al., 2024b; Xiong et al., 2023; Maekawa et al., 2024; Wang et al., 2023b). DD synthesizes small-scale datasets reflecting rich information from the original large-scale ones and is firstly designed with meta-learning schemes (Wang et al., 2018; Nguyen et al., 2021b;a; Zhou et al., 2022; Loo et al., 2022; 2023). The optimization is conducted upon a meta loss where a neural network or estimation is built on the surrogate data and then evaluated on the real data. Other methods optimize the synthetic images by matching training characteristics with real images (Zhao et al., 2021; Zhao & Bilén, 2023; Liu et al., 2023; Vahidian et al., 2024; Cazenavette et al., 2022; Zhao et al., 2023). The imitation on real distribution effectively improves the information contained in small surrogate datasets. Data parametrization (Kim et al., 2022b; Liu et al., 2022; Wei et al., 2024) and generative prior (Cazenavette et al., 2023; Gu et al., 2024a; Wang et al., 2023a) are also considered for more efficient DD method construction. However, most of existing DD methods remain as black boxes, lacking the ability of explicitly controlling the distilling direction. As a result, the practicality of DD methods are still poor from real-world applications. In this work, we aim at enhancing both the interpretability and controllability of the dataset distillation process.

Diffusion Models Diffusion models have acquired substantial success in generating high-quality images (Ho et al., 2020; Dhariwal & Nichol, 2021; Kingma et al., 2021; Nichol & Dhariwal, 2021). There have also been a series of works focusing on diffusion-based image manipulating or editing. DiffusionCLIP incorporates a CLIP model into the diffusion model re-tuning to provide optimization guidance (Kim et al., 2022a). DiffuseIT, DiffEdit and Prompt-to-Prompt integrate the editing into manifold constraint, mask guidance and cross attention control, respectively (Kwon & Ye, 2023; Couairon et al., 2023; Hertz et al., 2023). However, most of them manipulate image instances following certain instructions. SDEdit proposes to control the training data generation, yet it requires the assistance of pre-trained models (Yeo et al., 2024). In this work, we design a training-free denoising guidance towards images suitable for model training.

3 METHOD

In this section, we demonstrate the detailed modules of our proposed CONCEPT-INFORMED Diffusion method (CONCORD). Firstly, we present the preliminary knowledge on dataset distillation and the possibility of performing concept-informed diffusion in Sec. 3.1. Subsequently we illustrate the design of concept acquirement and objective design in Sec. 3.2 and Sec. 3.3, respectively.

3.1 CONCEPT-INFORMED DIFFUSION

Given a target real dataset $\mathcal{T} = \{f(x_i; y_i)g_{i=1}^{jTj}\}$, the aim of dataset distillation is to generate a small surrogate dataset $\mathcal{S} = \{f(x_i; y_i)g_{i=1}^{jSj}\}$, where $jSj \ll jTj$, such that training a network on \mathcal{S} approximates as closely as possible the performance attained when training on \mathcal{T} . Typical methods incorporate meta-learning or metric matching to condense the information from real data into the surrogate dataset. However, the dependence on bi-level optimization often leads to excessive computation demands and bias towards specific adopted architectures (Sun et al., 2024; Zhou et al., 2023). Recently, methods utilizing the generative priors of diffusion models emerge as solutions for more efficient dataset distillation (Gu et al., 2024a; Su et al., 2024; Moser et al., 2024).

Diffusion for Distillation Diffusion-based generative models learn data distributions via denoising. Firstly, a forward process is defined by obtaining $x^{(T)}$ from clean data $x^{(0)} \sim q(x^{(0)})$ as a Markov chain of gradually adding Gaussian noise at time steps (Ho et al., 2020):

$$q(x^{(1:T)}|x^{(0)}) := \prod_{t=1}^T q(x^{(t)}|x^{(t-1)}); \text{ where } q(x^{(t)}|x^{(t-1)}) := N(x^{(t)}; \sqrt{1-\alpha_t}x^{(t-1)}; \alpha_t); \quad (1)$$

where $\alpha_t \in (0, 1)$ is a variance schedule. Denoting $\beta_t := 1 - \alpha_t$ and $\beta_t := \prod_{s=0}^{t-1} \alpha_s$, $x^{(t)}$ at an arbitrary time step can be directly sampled with a Gaussian noise $\epsilon \sim N(0; 1)$:

$$x^{(t)} = \sqrt{\beta_t}x^{(0)} + \sqrt{1-\beta_t}\epsilon; \quad (2)$$

Denoising diffusion probabilistic models (DDPMs) approximate the data distribution with a network:

$$p(x^{(t-1)}|x^{(t)}) = N(x^{(t-1)}; \mu(x^{(t)}; t); \sigma(x^{(t)}; t)); \quad (3)$$

162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215

Figure 2: The pipeline of the proposed CONCORD method. Descriptive concepts are retrieved and utilized to inform the diffusion denoising process. The samples with better instance-level concept completeness help to construct a surrogate dataset with better overall quality.

where

$$(x^{(t)}; t) := \frac{1}{p-t} x^{(t)} - \frac{1}{1-t} (x^{(t)}; t); \quad (4)$$

and the $(x^{(t)}; t)$ is the predicted noise, wheris optimized by:

$$\min_{t; x^{(0)} \sim q(x^{(0)}); N(0,1)} \mathbb{E} \left[\left(\frac{1}{p-t} x^{(0)} + \frac{1}{1-t} (x^{(t)}; t) \right)^2 \right]; \quad (5)$$

By denoising a fixed number of random noises, a surrogate dataset can be generated that encapsulates the distribution of original data. Gu et al. (2024a) introduce additional minimax criteria to distill more representative and diverse samples from real data, improving the quality of the generated surrogate datasets. However, the distilling process primarily focuses on imitating dataset-level concept distributions, while overlooking the instance-level conceptual completeness at the inference stage. Since the final distilled samples are directly derived from random noise, without explicit control over the content, essential object details might be missing or incorrectly represented in the generated images. Moreover, the constrained storage budget typical of dataset distillation benchmarks limits the ability to compensate the instance-level information loss by increasing data scale, further compromising the quality of the distilled dataset. Thus, there is an urgent demand for techniques that allow for explicit control during the denoising process, enhancing both the conceptual completeness and the overall quality of the surrogate dataset.

Concept Informing Dhariwal & Nichol (2021) introduce classifier guidance with the gradients of a classifier network $\nabla_{x^{(t)}} \log p(y|x^{(t)}; t)$ during the diffusion process. However, when the classifier can acquire activations from a broad set of possible details to make predictions, the conceptual completeness associated with the specific category can remain insufficient. Therefore, we propose to explicitly inform the diffusion process with fine-grained and distinguishable concepts tied to the category (e.g. attributes). The concept-informed diffusion offers several advantages: firstly, various concepts of a category provide more detailed information compared with using the category alone, allowing for explicit reasoning and refinement during the generation process. Secondly, in circumstances where classifiers are difficult to obtain, concepts remain viable given the category. We define the set of concepts associated with the category label of the current samples $C = \{a_j g_j^{[C]}\}$, where $|C|$ is a pre-defined number of concepts. Subsequently, an objective $O(x^{(t)}; C)$ can be derived reflecting the semantic similarity between the generated sample and these concepts. The informed denoising process can be represented with the objective as:

$$p(x^{(t-1)}|x^{(t)}) = N(x^{(t-1)}; (x^{(t)}; t) + (x^{(t)}; t)r_{x^{(t)}} O(x^{(t)}; C); (x^{(t)}; t)); \quad (6)$$

Song et al. (2020) introduce another form of denoising diffusion implicit models (DDIMs) that construct a deterministic non-Markovian inference process as:

$$x^{(t-1)} = \frac{p}{t-1} x^{(0)} + \frac{p}{1-t-1} (x^{(t)}; t); \quad (7)$$

where the estimated observation $x^{(0)}$ of clean original data $x^{(0)}$ can be obtained by computing the posterior expectation with $x^{(t)}$ (Robbins, 1992):

$$x^{(0)} := \frac{x^{(t)}}{p-t} - \frac{1}{p-t} (x^{(t)}; t); \quad (8)$$

Algorithm 1: Concept-Informed Diffusion

Input: diffusion model \mathcal{D} , original dataset \mathcal{T} , concept set \mathcal{C} , required sample number N_s
Output: surrogate dataset \mathcal{S}
Initialize the surrogate dataset $\mathcal{S} = \emptyset$
for index in $1..N_s$ do
 Obtain a random noisy sample $x^{(T)}$ and a category label c
 Retrieve the positive and negative concepts \mathcal{P} and \mathcal{N} from \mathcal{C} according to label c
 for time step in $T..1$ do
 Predict the noise $(x^{(t)}; t)$
 Calculate the concept matching object $O(x^{(t)}; \mathcal{C}; \mathcal{C})$ according to Eq. 12
 Update the predicted noise $\hat{x}^{(t)}$ according to Eq. 9
 Conduct denoising step to obtain $x^{(t-1)}$ according to Eq. 7
 end
 Add the predicted clean sample to the surrogate dataset \mathcal{S} set $x^{(0)}$
end

Similarly, we can apply the concept informing through:

$$\hat{x}^{(t)} := (x^{(t)}; t) \oplus \lambda \frac{1}{\tau} \text{tr}_{x^{(t)}} O(x^{(t)}; \mathcal{C}); \quad (9)$$

where λ is the informing weight adjustable for control extent. The updated $\hat{x}^{(t)}$ is subsequently used for the above reverse diffusion process. When the informing can be applied to both frameworks, in this work, we mainly incorporate DDIM for developing our algorithm.

3.2 CONCEPT ACQUIREMENT

Based on the aim of enhancing the discriminative details and mitigating conceptual incompleteness in the generated images, we intend to explicitly inform the diffusion process with distinguishable concepts. While concluding or manually designing visual concepts being infeasible for a large number of categories, large language models (LLMs) offer a valuable solution with rich conceptual understanding acquired during the training process. Inspired by this, we design prompts for the corresponding categories in the target dataset to elicit fine-grained attributes from LLMs, which are used as concepts to inform the diffusion process. Menon & Vondrick (2023) design prompts to retrieve descriptions used for zero-shot image classification. While our task shares certain similarity, it differs primarily in the nature of the required descriptions. Descriptions used for classification should comprehensively reflect various aspects of the corresponding category. In comparison, those used for constructing surrogate datasets are supposed to be distinguishable across different categories to ensure that the generated data can facilitate model training. Therefore, we design an example prompt shown in Fig. 2, where distinction from other classes is emphasized for retrieving descriptions.

Concept Validity Evaluation Once a set of concepts is obtained, it is crucial to evaluate their validity on actual data, as some concepts may not be well-represented in the real data due to biases in data collection. Thus, before the concept matching process, we first retrieve an over-abundant amount of concepts, and then filter them through a validity evaluation process to identify those with the strongest relevance to the real data. For this purpose, we utilize a CLIP model to extract embedded features from both real images and textual descriptions. For a category activation

A of a text description c on the images $\{x_i\}$, $y_i = \text{lg}_{i=1}^{ij}$ can be calculated by:

$$A = \frac{1}{|\mathcal{I}|} \sum_{i=1}^{|\mathcal{I}|} \cos(\phi(x_i); c); \quad (10)$$

where $\phi(\cdot)$ denotes the embedded feature extraction function of the CLIP model, and computes the cosine similarity. We select a pre-defined number C of descriptions for each category with the highest activation scores for further use in the informing process. This ensures that the selected concepts retain the integrity of the knowledge distilled from the real dataset, making them more representative and relevant for improving instance-level conceptual completeness in generated data.

3.3 CONCEPTMATCHING

With the distinguishable concepts acquired, a straightforward approach to measure the relationship between the generated sample $x^{(t)}$ and the corresponding concept c_j is to compute their cosine similarity:

$$O(x^{(t)}; C) = \frac{1}{|C|} \sum_{j=1}^{|C|} \langle x^{(t)}; c_j \rangle \quad (11)$$

which is similar to the concept validity evaluation process. We argue that beyond the positive informing from the concepts associated with the corresponding category, it is equally important to adjust the diffusion control by considering the overall dataset distribution. Therefore, we employ concepts from other categories as negative samples to provide more stable diffusion guidance.

Contrastive Matching Inspired by the contrastive loss adopted in CLIP training (Radford et al., 2021; Patel et al., 2023), we propose a similar strategy to incorporate negative concepts. Since multiple positive concepts should work together to provide adequate guidance, we modify the supervised contrastive loss (Khosla et al., 2020) into an image-text version:

$$O(x^{(t)}; C, \bar{C}) = \frac{1}{|C|} \sum_{i=1}^{|C|} \log \frac{\exp(\langle x^{(t)}; c_i \rangle)}{\exp(\langle x^{(t)}; c_i \rangle) + \sum_{c_j \in \bar{C}} \exp(\langle x^{(t)}; c_j \rangle)} \quad (12)$$

where \bar{C} denotes the set of negative concepts.

Negative Concept Selection With a large number of potential negative categories, selecting appropriate negative concepts is essential for effectively informing the diffusion process. We first compute the cosine similarity between the category labels, and use the similarity as sampling weight for negative category selection. This approach ensures that categories with higher similarity to the target category are prioritized as negative samples. Compared with random selection, the similarity-based approach offers more precise control over the diffusion process. Additionally, compared with a fixed range of negative categories, the dynamic sampling allows for more diverse denoising control.

4 EXPERIMENTS

4.1 IMPLEMENTATION DETAILS

We adopt Minimax (Gu et al., 2024a) and Stable Diffusion unCLIP Img2Img (Ramesh et al., 2022) as baselines to evaluate our proposed training-free approach, which is applied at the inference stage. The informing weight α in Eq. 9 is set as 1. For each category, 5 descriptive attributes are selected with the highest activation scores, as detailed in Sec. 3.2. And 10 negative descriptions from different categories are used for contrastive loss calculation. A total denoising step number of 50 is adopted for the generation process, and the generated images are resized to 256 subsequent validation. The validation protocol follows RDED (Sun et al., 2024), where soft label is adopted to obtain better performance. The model training lasts for 300 epochs. All reported results are based on 3 random runs, with the averaged accuracy and the variance included. All the experiments are conducted on a single NVIDIA A100 GPU. Further implementation details are provided in Sec. B.

We believe that DD for small-resolution datasets has been well solved by previous methods. Thus, the main experiments in this work are conducted on ImageNet-1K (Deng et al., 2009) and its sub-sets including ImageNet-100 and ImageWoof (Fastai). Additionally, we incorporate Food-101 (Bossard et al., 2014) as another benchmark to evaluate the effectiveness of the proposed CONCORD method.

4.2 COMPARISON WITH STATE-OF-THE-ARTS

Firstly we conduct the experiments on standard benchmarks, reporting the performance on multiple different architectures. The compared methods include MTT (Cazenavette et al., 2022), SIRE (Sire et al., 2023), RDED (Sun et al., 2024), DiT (Peebles & Xie, 2023), Minimax (Gu et al., 2024a), and Img2Img (Ramesh et al., 2022). The results on ImageWoof, ImageNet-100, and the full ImageNet-1K are shown in Tab. 1 and Tab. 2, respectively.

Table 1: Performance comparison with state-of-the-art methods on ImageWoof. The superscript indicates the application of our proposed CONCORD method. Bold entries indicate best results, and underlined ones illustrate improvement over baseline.

IPC (Ratio)	Test Model	MTT	SR ² L	RDED	DiT	Minimax	Minimax ^C	unCLIP	unCLIP ^C
1 (0.08%)	ConvNet	28.6 ^{0:8}	-	18.5 ^{0:9}	20.5 ^{0:8}	16.7 ^{0:2}	<u>17.8^{0:8}</u>	20.5 ^{0:4}	19.9 ^{0:7}
	ResNet-18	-	13.3 ^{0:5}	20.8 ^{1:2}	18.3 ^{0:7}	15.3 ^{1:1}	<u>16.9^{1:0}</u>	16.7 ^{0:7}	17.4 ^{1:1}
	ResNet-101	-	13.4 ^{0:1}	19.6 ^{1:8}	17.1 ^{1:3}	14.2 ^{1:1}	<u>14.9^{1:3}</u>	14.9 ^{0:2}	<u>15.3^{1:3}</u>
10 (0.8%)	ConvNet	35.8 ^{1:8}	-	40.6 ^{2:0}	42.2 ^{1:2}	41.2 ^{0:8}	<u>43.1^{0:5}</u>	40.1 ^{0:8}	41.2 ^{0:8}
	ResNet-18	-	20.2 ^{0:2}	38.5 ^{2:1}	38.2 ^{1:1}	42.8 ^{1:1}	<u>44.4^{0:9}</u>	37.9 ^{1:1}	40.7 ^{0:4}
	ResNet-101	-	17.7 ^{0:9}	31.3 ^{1:3}	31.1 ^{0:3}	35.7 ^{0:9}	<u>36.5^{0:9}</u>	30.7 ^{0:9}	<u>31.9^{1:1}</u>
50 (3.8%)	ConvNet	-	-	61.5 ^{0:3}	59.9 ^{0:2}	61.1 ^{0:8}	<u>62.5^{0:9}</u>	59.5 ^{1:4}	60.4 ^{0:4}
	ResNet-18	-	23.3 ^{0:3}	68.5 ^{0:7}	65.9 ^{0:2}	67.8 ^{0:5}	<u>69.2^{1:0}</u>	63.6 ^{0:6}	<u>66.1^{1:1}</u>
	ResNet-101	-	21.2 ^{0:2}	59.1 ^{0:7}	60.1 ^{1:1}	62.2 ^{0:6}	<u>63.6^{0:2}</u>	60.0 ^{1:0}	<u>60.8^{0:9}</u>

Table 2: Performance comparison with state-of-the-art methods on ImageNet-100 (left) and ImageNet-1K (right). The superscript indicates the application of our proposed CONCORD method. Bold entries indicate best results, and underlined ones illustrate improvement over baseline.

Method	IPC			Method	IPC		
	1	10	50		1	10	50
SR ² L	3.0 ^{0:3}	9.5 ^{0:4}	27.0 ^{0:4}	SR ² L	0.1 ^{0:1}	21.3 ^{0:6}	46.8 ^{0:2}
RDED	8.1 ^{0:3}	36.0 ^{0:3}	61.6 ^{0:1}	RDED	6.6 ^{0:2}	42.0 ^{0:1}	56.5 ^{0:1}
DiT	8.2 ^{0:1}	29.5 ^{0:4}	59.8 ^{0:5}	DiT	6.1 ^{0:1}	41.3 ^{0:3}	56.6 ^{0:2}
Minimax	5.8 ^{0:2}	31.6 ^{0:1}	64.0 ^{0:5}	Minimax	6.0 ^{0:1}	43.4 ^{0:3}	59.1 ^{0:1}
Minimax ^C	7.1 ^{0:2}	33.3 ^{0:6}	64.9 ^{0:3}	Minimax ^C	6.4 ^{0:2}	43.8 ^{0:6}	59.4 ^{0:2}
unCLIP	7.1 ^{0:1}	26.9 ^{0:4}	64.6 ^{0:2}	unCLIP	5.9 ^{0:2}	42.0 ^{0:3}	58.1 ^{0:2}
unCLIP ^C	7.7 ^{0:2}	28.1 ^{0:7}	65.4 ^{0:4}	unCLIP ^C	6.2 ^{0:3}	42.5 ^{0:2}	58.5 ^{0:1}

Under the 1 Image-per-class (IPC) setting, previous methods MTT and RDED have demonstrated the best performance, with the vanilla DiT model also showing strong results. Minimax is re-tuned to enhance the representativeness and diversity of the generated data. Although it is less effective under small IPC settings, the performance superiority is more substantial as the IPC increases. The unCLIP Img2Img model is not specially trained or re-tuned on ImageNet, but still yields comparable performance by direct inference. When the proposed CONCORD method is applied to both baseline methods, significant performance improvements are observed across all IPC settings and architectures. These results indicate that refining instance-level conceptual completeness is essential for constructing more effective distilled datasets. However, we can also notice that the performance gain is less significant as the class number increases. A potential explanation is that the influence of instance-level quality diminishes as the overall data scale is larger. Despite this, the proposed CONCORD method achieves state-of-the-art performance on the full ImageNet-1K dataset and its subsets, especially on large IPC settings, further supporting its effectiveness in dataset distillation.

Additionally, we conduct experiments on Food-101 with unCLIP Img2Img as the baseline in Tab. 3. It simulates actual DD application scenarios for custom datasets. The results suggest that methods based on generative prior are capable and practical to perform custom DD tasks without extra training efforts. While the unCLIP baseline performs worse than random selection under the 50-IPC setting, the proposed CONCORD method still enhances the quality of distilled datasets across all IPC settings. It opens up new possibilities for resource-limited researchers to perform custom DD.

4.3 ABLATION STUDY AND DISCUSSION

In this section we conduct component analysis and experimental results on extended settings. By default, the experiments are conducted on ImageWoof, with unCLIP Img2Img as the baseline.

Prompt Design We employ LLMs to retrieve essential visual descriptions as the informing target. The description quality is crucial for achieving optimal informing effects. Therefore, a quantita-

Table 3: Performance comparison with unCLIP on CLIP Img2Img on Food-101 dataset.

Method	IPC		
	1	10	50
Random	5.0 _{0:1}	30.1 _{0:1}	64.0 _{0:2}
unCLIP	6.4 _{0:1}	30.7 _{0:2}	61.3 _{0:3}
unCLIP ^C	6.9 _{0:1}	32.0 _{0:2}	62.5 _{0:2}

Table 4: Comparison with different prompts for concept retrieval on ImageWoof.

Method	IPC		
	1	10	50
Classification	17.6 _{2:0}	38.2 _{1:3}	64.1 _{0:7}
Ours-3.5	16.8 _{0:5}	38.8 _{0:2}	65.4 _{1:1}
Ours-4	17.4 _{1:1}	40.7 _{0:4}	66.1 _{1:1}

Table 5: Comparison with different negative description selection on ImageWoof.

Method	IPC		
	1	10	50
Random	15.5 _{1:6}	39.5 _{1:2}	64.9 _{0:2}
Similar-10	16.1 _{0:6}	38.3 _{0:9}	65.3 _{1:2}
Similar-25	15.9 _{1:0}	37.9 _{0:4}	64.5 _{0:5}
Similar-50	15.9 _{1:4}	38.3 _{1:2}	64.6 _{0:4}
Weighted	17.4 _{1:1}	40.7 _{0:4}	66.1 _{1:1}

Table 6: Ablation study on the optimization baseline and objectives on ImageWoof.

Base	Objective	IPC		
		1	10	50
DiT	None	18.3 _{0:7}	38.2 _{1:1}	65.9 _{0:2}
	Contrastive	20.3 _{0:7}	40.5 _{1:2}	67.6 _{0:4}
unCLIP	None	16.7 _{0:7}	37.9 _{1:1}	63.6 _{0:6}
	Classifier	16.9 _{0:7}	38.5 _{1:0}	65.2 _{0:8}
	Cosine	18.2 _{1:6}	39.7 _{1:1}	63.9 _{0:2}
	Contrastive	17.4 _{1:1}	40.7 _{0:4}	66.1 _{1:1}

Comparative comparison between different prompt design and LLM models is provided in Tab. 4. Menon & Vondrick (2023) design prompts to retrieve descriptions for zero-shot classification, denoted as “Classification” in the table. While the retrieved concepts improve performance when IPC=1, the impact is less significant for larger IPCs. Accordingly, we design a new prompt (shown in Fig. 2) that emphasizes distinguishable appearance features. The descriptions are retrieved from GPT-3.5 and GPT-4, and the GPT-4 version achieves overall the best performance improvement. Detailed examples of the retrieved descriptions are shown in Fig. 7 for further investigation.

Negative Description Selection In the contrastive objective of Eq. 12, negative concepts are introduced for more accurate informing. While the extra constraint potentially brings more information, the selection of negative concepts is critical for stable optimization. Therefore, we evaluate the influence of different selection strategies in Tab. 5. Firstly, random selection from all categories considerably enhances the quality of the distilled dataset. Given that concepts from similar categories can serve as more challenging negative guidance, we narrow the random selection range to include only the top-similar categories, denoted as “Similar-#” with the number indicating the range. However, the unstable performance improvement suggests that limiting the diversity of negative concepts can harm the informing effect. Eventually, we propose to adopt a weighted sampling strategy based on category similarity. By simultaneously emphasizing similar categories and maintaining diversity, the strategy achieves the most significant and stable performance improvement.

Optimization We adopt a contrastive design of the objective to incorporate negative descriptions and stabilize the informing process. Accordingly, different objective forms are compared in Tab. 6 to evaluate the effectiveness of this design. After tuning, the informing weight for classifier guidance is set as 0.05 for best performance, which provides consistent improvement over the baseline. However, the supervision from a class-level is too coarse to refine the necessary details for sample generation. This limitation is evident in the superior performance achieved by the contrastive objective, which offers more detailed guidance. Additionally, the reliance on extra pre-trained classifiers also reduces the practicality of classifier guidance. Comparatively, the cosine objective in Eq. 11 yields even larger improvement when only 1 image is used for training. As the IPC grows, the performance improvement decreases, potentially due to limited diversity from only positive concepts. Since the proposed CONCORD method is designed to work without the need for pre-trained classifiers, we focus exclusively on concept informing in the main experiments.

We also conduct experiments on the vanilla DiT model without Minimax fine-tuning, where the top-1 accuracy improves by 2% across different IPCs. It further validates our hypothesis that instance-level conceptual completeness is essential for dataset distillation methods based on generative prior.

432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485

Figure 3: Comparison between images informed by fine-grained descriptive concepts (the first row) and the class name alone (the second row). From left to right the informing weight is gradually increased. Descriptions and corresponding image details are highlighted to illustrate better control with distinguishable concepts.

And the proposed CONCORD method can be broadly applied to existing diffusion pipelines to enhance the quality of the distilled datasets, which proves its practicality.

Sample Visualization We present example generated images with and without our proposed CONCORD method in Fig. 1. The baseline Minimax (Gu et al., 2024a) generates images with realistic texture and diverse variations. However, it overlooks the instance-level conceptual completeness, where essential concepts are often incorrect or missing (e.g., the unnatural shape of the coffee mug and the absence of beacon in the beacon image). By applying the CONCORD method, the generated images demonstrate substantial improvement in representing essential object details. In dataset distillation, where the number of samples is limited, the instance-level defects can severely affect the quality of the distilled dataset. In contrast, by emphasizing conceptual completeness at the instance level, our proposed CONCORD method enhances the overall quality of generated samples, also providing interpretability for the superior performance.

Effectiveness of Descriptions We employ descriptive attributes generated by language models as concepts to inform the denoising process. In Fig. 3 we compare the informing effects using detailed descriptions versus class names. For avoiding the influence of objective forms, we perform the experiments using cosine similarity as in Eq. 11, and only match positive concepts. Several conclusions can be drawn from the comparison of results. Firstly, while class names provide certain level of concept understanding, fine-grained descriptions offer more precise control over the diffusion process. For instance, when informed by a description like “legs covered in thick fur”, the length of the leg fur visibly increases as the informing weight grows, whereas images constrained by only the class name do not show a similar trend. Secondly, as the informing weight increases, images constrained by class names tend to collapse more quickly. It indicates that fine-grained concept informing provides better stability during the diffusion process compared with relying solely on class names. Thirdly, crucial descriptions such as “otter-shaped head with expressive eyes” effectively constrain the diffusion process. Even as images start to collapse, the head shape remains similar to the original generation result. In contrast, without explicit constraints from fine-grained descriptions, images informed by the class name show concept shift in these discriminative details.

IPC Scale-up An advantage offered by distillation methods based on generative prior is the flexibility to create surrogate datasets of varying sizes. Beyond the standard small-size benchmarks, we further extend the dataset size to 200 IPC in Fig. 4a. Across all IPC settings, the proposed CONCORD method provides consistent improvement upon both Minimax and unCLIP baselines. Notably, with 200 images per class, Minimax with CONCORD achieves the top-1 accuracy attained with the entire original ImageWoof dataset, following the same validation protocol.

486
487
488
489
490
491
492
493
494
495
496
497
498

(a) (b) (c)

Figure 4: (a) Applying the proposed concept informing brings consistent improvement across all IPC settings. With 200 Images per class, our method achieves the performance attained with the full original set. (b/c) Parameter analysis on informing weight and negative sample number.

502
503
504

4.4 PARAMETER ANALYSIS

505
506
507

There are multiple hyper-parameters involved in the proposed method. In this section, we perform analysis by adjusting the parameter to observe the influence on the performance.

508
509
510
511
512
513
514

Informing Weight The informing weight controls the degree of influence applied to the denoising diffusion process. As shown in Fig. 4b, setting $\omega = 0$ results in standard inference without concept informing. As ω increases within a reasonable range, the performance is also improved, indicating that the injected concept information enhances the quality of distilled datasets. However, if ω is too high, it disrupts the standard denoising process, leading to performance drop. Through comparison, we set the value of ω as 2.0 for balance between sufficient control and stable denoising.

515
516
517
518
519
520
521
522

Negative Sample Number The number of negative concepts is critical for constructing an effective contrastive loss. Therefore, we investigate the influence of negative sample number in Fig. 4c. When zero negative samples are used, cosine objective is applied for informing as in Eq. 11. Both too few or too many negative samples lead to unstable optimization and sub-optimal performance. Unlike standard contrastive learning, where the encoder separates different instances, the goal in DD is to focus on emphasizing essential object concepts. Therefore, enhancing positive concepts is more important. Based on our analysis, we adopt 10 negative samples in the contrastive objective to provide an appropriate constraint while maintaining stable optimization.

523
524

5 CONCLUSION

525
526
527
528
529
530
531
532
533
534

In this work, we propose to incorporate the conceptual understanding of large language models (LLMs) to enhance instance-level image quality for dataset distillation. Specifically, distinguishable concepts are retrieved based on category labels, and are subsequently utilized to inform the diffusion-based sample generation process. The conceptual completeness obtained by the proposed CONCEPT-inFORMed Diffusion (CONCORD) process mitigates the information loss caused by image defects, leading to higher overall quality of distilled datasets. CONCORD is evaluated on multiple baselines, and achieves state-of-the-art performance on the full ImageNet-1K dataset. The generated real-looking images with necessary details provide explicit interpretability for their effectiveness, and also prompt new possibilities of down-stream applications of dataset distillation.

535
536
537
538
539

Limitations and Future Works The proposed concept informing method significantly improves the instance-level concept completeness, and thereby enhances the performance of the distilled data. But simultaneously, it also involves extra computational cost. Since the informing is conducted throughout the diffusion denoising process, the method might not be applicable to few-step diffusion techniques, which aim to reduce computational overhead. In future works, we will explore efficient diffusion inference techniques for more practical dataset distillation.

540 Reproducibility Statement We have provided implementation details regarding the baseline
541 preparation, the proposed CONCORD method as well as the evaluation process in the Appendix
542 Sec. B. We use the publicly available ImageNet dataset as well as its subsets for conducting experi-
543 ments. Additionally, the utilized source code is attached in the supplementary material, and will be
544 made public upon acceptance.

546 REFERENCES

- 547
548 Yutong Bai, Xinyang Geng, Karttikeya Mangalam, Amir Bar, Alan Yuille, Trevor Darrell, Jitendra
549 Malik, and Alexei A Efros. Sequential modeling enables scalable learning for large vision models.
550 arXiv preprint arXiv:2312.00785, 2023.
- 551 Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101—mining discriminative compo-
552 nents with random forests. *ECCV*, pp. 446–461. Springer, 2014.
- 553
554 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
555 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
556 few-shot learners. *NeurIPS* volume 33, pp. 1877–1901, 2020.
- 557 George Cazenavette, Tongzhou Wang, Antonio Torralba, Alexei A Efros, and Jun-Yan Zhu. Dataset
558 distillation by matching training trajectories. *CVPR*, pp. 4750–4759, 2022.
- 559
560 George Cazenavette, Tongzhou Wang, Antonio Torralba, Alexei A Efros, and Jun-Yan Zhu. Gener-
561 alizing dataset distillation via deep generative prior. *CVPR*, pp. 3739–3748, 2023.
- 562 Jin Chen, Zheng Liu, Xu Huang, Chenwang Wu, Qi Liu, Gangwei Jiang, Yuanhao Pu, Yuxuan
563 Lei, Xiaolong Chen, Xingmei Wang, et al. When large language models meet personalization:
564 Perspectives of challenges and opportunities. *World Wide Web* 27(4):42, 2024.
- 565
566 Guillaume Couairon, Jakob Verbeek, Holger Schwenk, and Matthieu Cord. Diffedit: Diffusion-
567 based semantic image editing with mask guidance. *ICLR*, 2023.
- 568 Justin Cui, Ruochen Wang, Si Si, and Cho-Jui Hsieh. Scaling up dataset distillation to imagenet-1k
569 with constant memory. *ICML*, pp. 6565–6590. PMLR, 2023.
- 570
571 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale
572 hierarchical image database. *CVPR*, pp. 248–255, 2009.
- 573 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep
574 bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805, 2018.
- 575
576 Prafulla Dhariwal and Alexander Nichol. Diffusion Models Beat GANs on Image Synthesis. In
577 *NeurIPS* volume 34, pp. 8780–8794, 2021.
- 578 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
579 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszko-
580 reit, and Neil Houlsby. An Image is Worth 16x16 Words: Transformers for Image Recognition at
581 Scale. *ICLR*, 2022.
- 582
583 Upol Ehsan, Elizabeth A Watkins, Philipp Wintersberger, Carina Manger, Sunnie SY Kim, Niels
584 Van Berkel, Andreas Riener, and Mark O Riedl. Human-centered explainable ai (hcxai): Reload-
585 ing explainability in the era of large language models (llms). *Extended Abstracts of the CHI*
586 *Conference on Human Factors in Computing Systems*, pp. 1–6, 2024.
- 587 Fastai. Fastai/imagenette: A smaller subset of 10 easily classif ed classes from imagenet, and a little
588 more french. URL <https://github.com/fastai/imagenette>
- 589
590 Jianyang Gu, Saeed Vahidian, Vyacheslav Kungurtsev, Haonan Wang, Wei Jiang, Yang You, and
591 Yiran Chen. Efficient dataset distillation via minimax diffusion. *CVPR*, pp. 15793–15803,
592 2024a.
- 593 Jianyang Gu, Kai Wang, Wei Jiang, and Yang You. Summarizing stream data for memory-restricted
online continual learning. *AAAI*, pp. 12217–12225, 2024b.

- 594 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-
595 nition. In CVPR, pp. 770–778, 2016.
- 596 Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollar, and Ross Girshick. Masked
597 autoencoders are scalable vision learners. In CVPR, pp. 16000–16009, 2022.
- 599 Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-or.
600 Prompt-to-prompt image editing with cross-attention control. In ICLR, 2023.
- 601 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising Diffusion Probabilistic Models. In NeurIPS
602 volume 33, pp. 6840–6851, 2020.
- 604 Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron
605 Maschinot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. In NeurIPS volume 33,
606 pp. 18661–18673, 2020.
- 607 Gwanghyun Kim, Taesung Kwon, and Jong Chul Ye. Diffusionclip: Text-guided diffusion models
608 for robust image manipulation. In CVPR, pp. 2426–2435, 2022a.
- 610 Jang-Hyun Kim, Jinuk Kim, Seong Joon Oh, Sangdoon Yun, Hwanjun Song, Joonhyun Jeong, Jung-
611 Woo Ha, and Hyun Oh Song. Dataset condensation via efficient synthetic-data parameterization.
612 In ICML, pp. 11102–11118, 2022b.
- 613 Diederik Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational diffusion models. In
614 NeurIPS volume 34, pp. 21696–21707, 2021.
- 616 Nicholas Kroeger, Dan Ley, Satyapriya Krishna, Chirag Agarwal, and Himabindu Lakkaraju. Are
617 large language models post hoc explainable? arXiv preprint arXiv:2310.05797, 2023.
- 618 Vyacheslav Kungurtsev, Yuanfang Peng, Jianyang Gu, Saeed Vahidian, Anthony Quinn, Fadwa
619 Idlahcen, and Yiran Chen. Dataset distillation from first principles: Integrating core information
620 extraction and purposeful learning. arXiv preprint arXiv:2409.01411, 2024.
- 622 Gihyun Kwon and Jong Chul Ye. Diffusion-based image translation using disentangled style and
623 content representation. In ICLR, 2023.
- 624 Songhua Liu, Kai Wang, Xingyi Yang, Jingwen Ye, and Xinchao Wang. Dataset distillation via
625 factorization. In NeurIPS 35:1100–1113, 2022.
- 627 Yanqing Liu, Jianyang Gu, Kai Wang, Zheng Zhu, Wei Jiang, and Yang You. Dream: Efficient
628 dataset distillation by representative matching. In ICCV, pp. 17314–17324, 2023.
- 629 Noel Loo, Ramin Hasani, Alexander Amini, and Daniela Rus. Efficient dataset distillation using
630 random feature approximation. In NeurIPS volume 35, pp. 13877–13891, 2022.
- 631 Noel Loo, Ramin Hasani, Mathias Lechner, and Daniela Rus. Dataset distillation with convexified
632 implicit gradients. arXiv preprint arXiv:2302.06755, 2023.
- 634 Aru Maekawa, Satoshi Kosugi, Kotaro Funakoshi, and Manabu Okumura. Dilm: Distilling dataset
635 into language model for text-level dataset distillation. arXiv preprint arXiv:2404.00264, 2024.
- 636 Sachit Menon and Carl Vondrick. Visual classification via description from large language models.
637 In ICLR, 2023.
- 639 Brian B Moser, Federico Raue, Sebastian Palacio, Stanislav Frolov, and Andreas Dengel. Latent
640 dataset distillation with diffusion models. arXiv preprint arXiv:2403.03881, 2024.
- 641 Timothy Nguyen, Zhouong Chen, and Jaehoon Lee. Dataset meta-learning from kernel ridge-
642 regression. In ICLR, 2021a.
- 643 Timothy Nguyen, Roman Novak, Lechao Xiao, and Jaehoon Lee. Dataset distillation with infinitely
644 wide convolutional networks. In NeurIPS 34:5186–5198, 2021b.
- 645 Alexander Quinn Nichol and Prafulla Dhariwal. Improved Denoising Diffusion Probabilistic Mod-
646 els. In ICML, pp. 8162–8171, 2021.

- 648 Maitreya Patel, Changhoon Kim, Sheng Cheng, Chitta Baral, and Yezhou Yang. Eclipse: A
649 resource-efficient text-to-image prior for image generation. *arXiv preprint arXiv:2312.04655*
650 2023.
- 651 William Peebles and Saining Xie. Scalable diffusion models with transformers. *ICCV*, pp. 4195–
652 4205, 2023.
- 653 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
654 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
655 models from natural language supervision. *ICML*, pp. 8748–8763. PMLR, 2021.
- 656 Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-
657 conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.
- 658 Herbert E Robbins. An empirical bayes approach to statistics. *Breakthroughs in Statistics:
659 Foundations and basic theory*, pp. 388–394. Springer, 1992.
- 660 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
661 Resolution Image Synthesis With Latent Diffusion Models. *ICMPR*, pp. 10684–10695, 2022.
- 662 Naveen Sachdeva and Julian McAuley. Data distillation: A survey. *Transactions on Machine
663 Learning Research*, 2023.
- 664 Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi
665 Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An
666 open large-scale dataset for training next generation image-text models. *NeurIPS* volume 35,
667 pp. 25278–25294, 2022.
- 668 Ozan Sener and Silvio Savarese. Active learning for convolutional neural networks: A core-set
669 approach. *International Conference on Learning Representation*, 2018.
- 670 Zhiqiang Shen and Eric Xing. A fast knowledge distillation framework for visual recognition. In
671 *ECCV*, pp. 673–690. Springer, 2022.
- 672 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv
673 preprint arXiv:2010.02502*, 2020.
- 674 Duo Su, Junjie Hou, Weizhi Gao, Yingjie Tian, and Bowen Tang. D^4 : Dataset distillation via
675 disentangled diffusion model. *ICVPR*, pp. 5809–5818, 2024.
- 676 Peng Sun, Bei Shi, Daiwei Yu, and Tao Lin. On the diversity and realism of distilled dataset: An
677 efficient dataset distillation paradigm. *ICVPR*, 2024.
- 678 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
679 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and
680 efficient foundation language model. *arXiv preprint arXiv:2302.13971*, 2023.
- 681 Saeed Vahidian, Mingyu Wang, Jianyang Gu, Vyacheslav Kungurtsev, Wei Jiang, and Yiran
682 Chen. Group distributionally robust dataset distillation with risk minimization. *arXiv preprint
683 arXiv:2402.04676*, 2024.
- 684 Kai Wang, Jianyang Gu, Daquan Zhou, Zheng Zhu, Wei Jiang, and Yang You. Dim: Distilling
685 dataset into generative models. *arXiv preprint arXiv:2303.04707*, 2023a.
- 686 Tongzhou Wang, Jun-Yan Zhu, Antonio Torralba, and Alexei A Efros. Dataset distillation. *arXiv
687 preprint arXiv:1811.10959*, 2018.
- 688 Ziyu Wang, Yue Xu, Cewu Lu, and Yong-Lu Li. Dancing with images: Video distillation via static-
689 dynamic disentanglement. *arXiv preprint arXiv:2312.00362*, 2023b.
- 690 Xing Wei, Anjia Cao, Funing Yang, and Zhiheng Ma. Sparse parameterization for epitomic dataset
691 distillation. In *NeurIPS* volume 36, 2024.
- 692 Max Welling. Herding dynamical weights to learn. *ICML*, pp. 1121–1128, 2009.

702 Sebastien C Wong, Adam Gatt, Victor Stamatescu, and Mark D McDonnell. Understanding data
703 augmentation for classification: when to warp? *DMCTA*, pp. 1–6. IEEE, 2016.
704

705 Yuanhao Xiong, Ruochen Wang, Minhao Cheng, Felix Yu, and Cho-Jui Hsieh. FedDM: Iterative
706 Distribution Matching for Communication-Efficient Federated Learning. *CVPR*, pp. 16323–
707 16332, 2023.

708 Teresa Yeo, Andrei Atanov, Harold Benoit, Aleksandr Alekseev, Ruchira Ray, Pooya Esmail
709 Akhoondi, and Amir Zamir. Controlled training data generation with diffusion models. *arXiv
710 preprint arXiv:2403.15309*, 2024.
711

712 Zeyuan Yin, Eric Xing, and Zhiqiang Shen. Squeeze, recover and relabel: Dataset condensation at
713 imagenet scale from a new perspective. *NeurIPS*, pp. 73582–73603, 2023.

714 Ruonan Yu, Songhua Liu, and Xinchao Wang. Dataset distillation: A comprehensive review.
715 *Transactions on Pattern Analysis and Machine Intelligence*, 2023.
716

717 Sangdoon Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo.
718 Cutmix: Regularization strategy to train strong classifiers with localizable features. *CVPR*, pp.
719 6023–6032, 2019.

720 Bo Zhao and Hakan Bilen. Dataset condensation with distribution matching. *WACV*, pp. 6514–
721 6523, 2023.

722 Bo Zhao, Konda Reddy Mopuri, and Hakan Bilen. Dataset condensation with gradient matching. In
723 *ICLR*, 2021.
724

725 Ganlong Zhao, Guanbin Li, Yipeng Qin, and Yizhou Yu. Improved Distribution Matching for
726 Dataset Condensation. *CVPR*, pp. 7856–7865, 2023.
727

728 Daquan Zhou, Kai Wang, Jianyang Gu, Xiangyu Peng, Dongze Lian, Yifan Zhang, Yang You, and
729 Jiashi Feng. Dataset quantization. *Proceedings of the IEEE/CVF International Conference on
730 Computer Vision*, pp. 17205–17216, 2023.

731 Yongchao Zhou, Ehsan Nezhadarya, and Jimmy Ba. Dataset distillation using neural feature regres-
732 sion. In *NeurIPS* volume 35, pp. 9813–9827, 2022.
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

APPENDIX

The appendix is organized into the following sections: In Sec. A we provide additional justification from the literature as far as the utility of using LLM-based concept informed learning. In Sec. B we introduce the implementation details of our method. In Sec. C we present more experiment results and analysis on the proposed CONCORD method. In Sec. D we show example generated images to better illustrate the effect of the proposed concept informing method.

A FURTHER ANALYTICAL GROUNDING

Perspective from XAI Explainable Artificial Intelligence (XAI) is an emerging area within machine learning that aims to provide users with greater insight into the functioning and mechanisms of black-box models, such as neural networks. Standard practice often involves knowledge extraction techniques, where broader, less precise, but simpler and thus more intuitive models are presented to explain the behavior of complex machine learning models. However, it has been noticed that this paradigm can be amended with the success and proliferation of LLMs (Ehsan et al., 2024). In particular, LLMs enable a more iterative and active learning procedure, where user prompts can directly inform the learning process to accommodate user needs. Simultaneously, LLMs can periodically generate language-based explanations, offering updates on model progress and adjustments.

One of the key advantages of LLMs is their potential for personalization (Chen et al., 2024). Given the rich variety of concepts derived from the extensive training data and the depth of developed models, LLMs are capable to foster a detailed and more human-centric understanding. This allows the models to tune the learning process towards specific application-driven concerns. The effectiveness of LLMs as explainers (Kroeger et al., 2023) provides the clear potential to address important use case concerns that are difficult to represent through standard analytical loss functions. This adaptability allows LLMs to bridge gaps between the learning objectives and real-world applications.

Perspective from Instrumental DD In the recent work by Kungurtsev et al. (2024), it has been argued that an important analytical consideration often overlooked in most optimization formulations for dataset distillation is its instrumentality. Specifically, synthetic data is typically not just used to solve the same learning problem in the same setting, but rather the dataset is expected to be used in some broader applications of interest to the user. These applications may have information needs that are not inherently condensed by standard off-the-shelf DD algorithms. By including additional custom criteria into the DD optimization formulation, while still incorporating existing powerful tools, DD can be more effectively steered towards performance on desired use cases. In this work, concepts are employed to facilitate natural human taxonomy with respect to object identification and recognition, and this consideration substantially improves the process by aligning the synthetic dataset with desired test performance outcomes.

B MORE IMPLEMENTATION DETAILS

Baselines We adopt Minimax (Gu et al., 2024a) and Stable Diffusion unCLIP Img2Img (Ramesh et al., 2022) as the baselines to illustrate the efficacy of our proposed concept informing method. These two baselines represent two different application scenarios, as outlined below.

For Minimax, a η -tuning process is conducted on ImageNet-1K. While η -tuning on target datasets yields superior performance, it also demands more resource consumption. Additionally, class labels are utilized for conditioning the denoising process, which might be inconvenient when extending the model to broader datasets. We adopt the default parameter setting in the original paper. The entire ImageNet-1K is partitioned into 50 subsets, each containing data of 20 classes. For each subset, a DiT model (Peebles & Xie, 2023) is η -tuned for 8 epochs. The mini-batch size, representative weight and diversity weight are set as 8, 0.002 and 0.008, respectively. During inference, the corresponding η -tuned model is loaded to generate data for specific classes.

For unCLIP Img2Img, we utilize the pre-trained model without any η -tuning adjustment. Random real images are fed into the model simultaneously with text prompts to generate high-quality

¹<https://huggingface.co/radames/stable-diffusion-2-1-unclip-img2img>

810 samples without losing the information of original data distribution. We adopt 28 prompt templates
 811 for generating images, e.g., “a photo of a nice $\$classname$ ” (Radford et al., 2021). The utilization
 812 of text prompts for conditioning provides significant flexibility, enabling data generation for custom
 813 datasets without extra training efforts. While the absence of fine-tuning may lead to slight reduction
 814 in generation quality, it allows for direct application of the proposed CONCORD method to any custom
 815 data given relevant text descriptions. During inference, the same pre-trained model is adopted
 816 for generating images for all target categories.

817
 818 **Concept Acquisition** We use GPT-4o to retrieve descriptive concepts for different categories.
 819 The full adopted prompt is as follows:

820
 821 You are an expert in computer vision and image analysis. Here is the $\langle task \rangle$ I want
 822 to use some visual descriptions to identify different categories in ImageNet dataset. Please
 823 first consider whether there exist categories with similar appearance to $\langle classname \rangle$. Then
 824 please give 10 short descriptions describing the appearance features of $\langle classname \rangle$
 825 has and can be used to distinguish it from other classes. The phrases should only focus on
 826 visual appearance of body parts or components instead of functioning. Each phrase should
 827 be detailed but also shorter than 128 characters. Each phrase starts with non-capitalized
 828 characters. $\langle task \rangle$ Give the answer in the form of answer: [$\$classname$, [$\$phrase1$,
 829 $\$phrase2$, $\$phrase3$, $\$phrase4$, $\$phrase5$, $\$phrase6$, $\$phrase7$, $\$phrase8$, $\$phrase9$,
 830 $\$phrase10$]] /answer .

831 After retrieving the original concepts, we perform a similarity calculation between the textual con-
 832 cepts and real images of the corresponding category. The top 5 most similar concepts are selected
 833 for the subsequent informed diffusion process, as described in Sec. 3.2. This approach helps ensure
 834 that the selected concepts align closely with the real images, thereby enhancing the validity of the
 835 concepts used in the diffusion process to a certain extent.

836
 837 **Informing** The informing process involves similarity calculation between embeddings of images
 838 and textual concepts. We use a CLIP model with ViT-L as the visual encoder, pre-trained on LAION-
 839 2B data (Schuhmann et al., 2022) to encode these embeddings. The model weights can be down-
 840 loaded from Hugging Face. The generation process involves 50 denoising steps for each sample.
 841 Prior to denoising, 5 descriptive concepts from the same class as well as 10 negative concepts each
 842 from a different class are retrieved for the sample. Before extracting text embeddings, the concepts
 843 are grouped with the corresponding class name using the following format:

844
 845 $f \ \$classname$ with $f \ \$concept$.

846
 847 During each denoising step, the similarity between the generated sample and corresponding con-
 848 cepts is calculated for the informing objective in Eq. 12. The informing weights are set as 1 for
 849 optimal performance. The concept informing guides the denoising process to obtain completeness
 850 on essential details, and thereby enhances the instance-level quality of the generated images.

851
 852 **Validation** We adopt the validation protocol in RDED (Sun et al., 2024) to evaluate the perfor-
 853 mance of distilled data. We mainly employ a ResNet-18 (He et al., 2016) architecture for experi-
 854 ments, with additional ones run on ResNet-101 and ConvNets as shown in Tab. 1. Specifically, for
 855 ImageWoof, ImageNet-100, ImageNet-1K, we adopt 5-layer, 6-layer, and 4-layer ConvNets, respec-
 856 tively, consistent with the settings in RDED. For ImageWoof, the images are resized to 224×224
 857 on ConvNet-5, while for all other cases, the images are resized to 224 for evaluation.

858 For ImageNet and its subsets, we employ pre-trained models to generate soft labels and apply
 859 Fast Knowledge Distillation (Shen & Xing, 2022). The models are trained for 300 epochs using
 860 the AdamW optimizer, with an initial learning rate of 0.001 and a weight decay of 0.01. A cosine
 861 annealing scheduler is used to adjust the learning rate. The mini-batch size for evaluation is set the
 862 same as IPC, e.g., a mini-batch size of 10 is adopted for evaluating 10-IPC sets. The applied data

863 ²<https://huggingface.co/laion/CLIP-ViT-L-14-laion2B-s32B-b82K>

³<https://github.com/LINs-lab/RDED>

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

Table 7: Comparison with different optimization objectives and their combination on ImageWoof.

Method	IPC		
	1	10	50
None	16.7 _{0:7}	37.9 _{1:1}	63.6 _{0:6}
Classifier	16.9 _{0:7}	38.5 _{1:0}	65.2 _{0:8}
Contrastive	17.4 _{1:1}	40.7 _{0:4}	66.1 _{1:1}
Combination	16.7 _{0:3}	38.8 _{1:9}	65.7 _{0:7}

Table 8: Inference time cost comparison of generating one sample under the mini-batch size of 1 between the baselines and the proposed CONCORD method.

Method	Minimax	Minimax ^C	unCLIP	unCLIP ^S
Time (s)	2.1	5.3	9.7	22.8

Figure 5: Parameter analysis on the diffusion denoising step number.

augmentation techniques include patch shuffling (Sun et al., 2024), random crop resize (Wong et al., 2016), random cropping and CutMix (Yun et al., 2019). After training on the distilled dataset, the model is then evaluated with the original validation set, and Top-1 accuracy is used as the validation performance. Each experiment is performed for three times, and the mean accuracy and standard variance are reported in the results.

For the Food-101 dataset, since no pre-trained models are provided by RDED, we train a ResNet-18 model on the original training set for 300 epochs, and use it for soft-labeling. It is important to note that the utilization of pre-trained models is independent from the sample generation process, and is only for fair comparison with state-of-the-art methods, which can be omitted in actual applications.

C EXTENDED EXPERIMENTS AND ANALYSIS

Ablation on Denoising Steps In the main experiments, we adopt 50 denoising steps for sample generation. The effect of varying the number of denoising steps is evaluated and presented in Fig. 5, with the unCLIP Img2Img model as the baseline. As the number of denoising steps increases, the accuracy of the baseline distilled data shows an upward trend under the IPC setting of 10, while is relatively consistent for the 50-IPC setting. For concept informing, fewer denoising steps result in insufficient informing, leading to performance similar to that of original images. Conversely, when too many denoising steps are used, the informing starts to disrupt the standard denoising process, leading to a drop in performance. While more fine-grained tuning of the informing weight could potentially mitigate the negative effect, more denoising steps also lead to extra computational cost. Therefore, we adopt 50 denoising steps as the standard setting.

Combining Concept Informing and Classifier Guidance The proposed CONCORD method informs the diffusion process to contain more discriminative details for enhancing instance-level image quality. While concept informing shares similarity to classifier guidance, the key difference is that CONCORD utilizes the similarity between generated samples and descriptive concepts as optimization targets, without relying on pre-trained classifiers. We also conduct the experiment to combine these two kinds of constraints together to simulate scenarios where pre-trained classifiers are available. As shown in Tab. 7, when functioning independently, our proposed contrastive concept informing outperforms classifier guidance. It supports our hypothesis that detailed descriptions provide richer information compared with the category-level labels, and are more helpful in refining the instance-level sample quality. However, when both types of guidance are combined, classifier guidance does not provide additional information, and disrupts the concept informing process. As a result, the combined approach shows less effective performance improvement on the generated images. Therefore, in the main experiments, we exclusively use the proposed concept informing, as it delivers better overall results and saves extra computational consumption.

Table 9: Performance comparison with state-of-the-art methods on ImageWoof. The superscript indicates the application of our proposed CONCORD method. Bold entries indicate best results, and underlined ones illustrate improvement over baseline.

IPC (Ratio)	Test Model	Random	K-Center	Herding	IDM	Minimax	Minimax ^C	Full
1 (0.08%)	ConvNet	16.3 ^{0:5}	15.8 ^{0:6}	16.8 ^{1:1}	17.1 ^{0:2}	16.7 ^{0:2}	<u>17.8</u> ^{0:8}	69.0 ^{0:2}
	ResNet-18	15.1 ^{0:2}	15.7 ^{0:8}	16.1 ^{0:4}	16.7 ^{0:5}	15.3 ^{1:1}	<u>16.9</u> ^{1:0}	76.9 ^{0:1}
	ResNet-101	14.0 ^{0:6}	13.7 ^{1:2}	14.1 ^{0:6}	16.3 ^{0:6}	14.2 ^{1:1}	<u>14.9</u> ^{1:3}	77.6 ^{0:2}
10 (0.8%)	ConvNet	40.5 ^{1:5}	37.1 ^{0:9}	41.2 ^{0:4}	38.5 ^{0:6}	41.2 ^{0:8}	<u>43.1</u> ^{0:5}	69.0 ^{0:2}
	ResNet-18	34.3 ^{1:6}	33.1 ^{0:5}	36.8 ^{0:6}	36.5 ^{1:2}	42.8 ^{1:1}	<u>44.4</u> ^{0:9}	76.9 ^{0:1}
	ResNet-101	32.1 ^{1:0}	31.6 ^{0:3}	33.8 ^{0:4}	30.8 ^{1:2}	35.7 ^{0:9}	<u>36.5</u> ^{0:9}	77.6 ^{0:2}
50 (3.8%)	ConvNet	60.9 ^{0:9}	57.7 ^{1:2}	60.4 ^{0:8}	61.0 ^{0:6}	61.1 ^{0:8}	<u>62.5</u> ^{0:9}	69.0 ^{0:2}
	ResNet-18	67.1 ^{1:0}	64.3 ^{0:9}	67.6 ^{0:5}	64.9 ^{0:6}	67.8 ^{0:5}	<u>69.2</u> ^{1:0}	76.9 ^{0:1}
	ResNet-101	61.4 ^{0:7}	58.8 ^{0:4}	60.8 ^{0:3}	57.2 ^{0:6}	62.2 ^{0:6}	<u>63.6</u> ^{0:2}	77.6 ^{0:2}

Extra Computational Cost We report the inference time cost for generating an image on both Minimax and unCLIP Img2Img in Tab 8. Comparatively, introducing CONCORD increases the original inference cost by approximately 1-1.5 times. unCLIP Img2Img involves Stable Diffusion v2-1 model (Rombach et al., 2022), which demands more computational resources compared with Minimax, which uses a DiT model (Peebles & Xie, 2023) as the denoising backbone. During inference, Minimax with CONCORD only requires about half the time of unCLIP baseline. Although Minimax performs better as a baseline, the advantage is based on extra re-tuning processes on the target dataset. Therefore, the model choice in real-world applications should consider multiple factors, including the balance between training and inference time consumption.

Comparison to More Baselines In addition to the results in Tab. 1, we also conduct experimental comparison with random sampling, K-Center (Sener & Savarese, 2018), Herding (Welling, 2009) and IDM (Zhao et al., 2023) in Tab. 9. For the methods based on original samples, we first resize the images to 128128 for ConvNet and 224224 for ResNet before running validation.

K-Center and Herding are two methods selecting coresets from the original data, with unstable performance improvement compared with random sampling. IDM is a dataset distillation method based on distribution matching, which is effective under small IPC settings. However, as the required sample number increases, the generated images often perform worse than random selected original samples. The baseline Minimax comparatively provides more stable information condensation across different IPC settings. When combined with the proposed CONCORD method, the overall dataset quality is significantly enhanced, surpassing all other methods in terms of accuracy. Especially for ConvNet and ResNet-18 architectures, training with 50 images per class achieves less than 10% performance gap from training with the entire original set. As larger models (ResNet-101) require more data and training iterations to get good performance, there still remains certain performance margin between distilled data and original full-set.

Feature Distribution Visualization We provide the feature distribution comparison in Fig. 6 to illustrate the effects of our proposed CONCORD method.

Firstly, the left figure shows the t-SNE features of samples generated with and without the informing of CONCORD. CONCORD works as a training-free guidance at the inference stage, without changing the main object in the images. By refining essential details in the generated samples, CONCORD enhances instance-level conceptual completeness, and improves the overall quality of the distilled datasets. However, these detail refinements have a mild effect on the feature distribution, indicating that with an already well-structured distribution, CONCORD can further improve performance without disrupting the underlying data distribution.

Secondly, the middle figure compares the generated images with the original ones used as conditioning in unCLIP Img2Img. The generated images closely align with the original data distribution, validating their effectiveness in capturing the properties of the original dataset. It demonstrates the suitability of using these generated images for training models.

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

(a) (b) (c)

Figure 6: Feature distribution visualization of (a) samples generated by unCLIP with and without CONCORD; (b) samples generated by unCLIP with CONCORD and original samples used for conditioning; (c) samples generated by Minimax with CONCORD and original samples. Different colors indicate different categories.

Figure 7: The comparison between concepts retrieved by different prompts and LLMs. “Cls” refers to prompts used for zero-shot classification attribute retrieval. Example images of corresponding classes are also presented to show their appearance features.

Lastly, the right figure shows the distribution of samples generated by Minimax with CONCORD and the entire original set. The generated samples demonstrate comprehensive coverage over the original distribution, ensuring that they represent a wide range of instances. While with sufficient diversity brought by samples distributed near decision boundaries, the generated samples also reduces noise in the overlapping regions between categories. It makes the generated dataset stable and effective for training models, when computational resources are limited.

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

Figure 8: Example generated image comparison on Minimax with and without the proposed CONCORD method (denoted as \mathcal{C}).

Analysis on the Retrieved Concepts The effects of different prompts and LLMs have been quantitatively investigated in Tab. 4. We further conduct qualitative comparison for the retrieved descriptions to explicitly analyze the informing effects of different concepts. As shown in Fig 7, descriptions of indigo bunting and streetcar are retrieved based on three settings: prompt for zero-shot classification on GPT-4o (denoted as “Cls”), our adopted prompt on GPT-3.5, and our prompt on GPT-4o. The “Cls” prompt retrieves general descriptions about the object. However, in many cases the retrieved descriptions are still too coarse for fine-grained informing. The descriptions retrieved by GPT-3.5 are more detailed, but contain a large number of non-visual attributes, which cannot provide valid signal during concept informing. Comparatively, our adopted prompt on GPT-4 successfully emphasizes the detailed visual features of corresponding categories. These fine-grained descriptions enable the proposed CONCORD method to effectively enhance instance-level conceptual completeness and further improve the overall quality of the distilled datasets.

D SAMPLE COMPARISON

We further present more example generated images in the following sections.

Comparison with Baselines Firstly, we show comparison on Minimax and unCLIP Img2Img with and without applying the proposed CONCORD method in Fig. 8 and Fig. 9, respectively. When baseline methods fail to present essential features and often lead to image defects, CONCORD significantly enhances the conceptual completeness in samples.

Failure Cases We also present failure cases where the proposed CONCORD method fails to correct or supplement essential features in the images in Fig. 10. It can be seen that the informing tries to modify some defects in the original image, but the eventual refinement is limited. There are also some cases where the informing fails to find the missing or incorrect details. Especially the informing fails to refine the details when the number of body parts is incorrect or the body part is

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133

Figure 9: Example generated image comparison on unCLIP Img2Img with and without the proposed CONCORD method (denoted as \mathcal{C}).

Figure 10: Example cases where CONCORD fails to supplement or modify incorrect concepts in the images^Q (indicates the application of CONCORD).



1150 Figure 11: Example images generated by the proposed CONCORD method on the Food-101 dataset.
1151 The class names are annotated below the images.



1174 Figure 12: Example animal images generated by the proposed CONCORD method on the ImageNet-1K dataset. The employed diffusion pipeline is the fine-tuned Minimax model. The class names are annotated below the images.
1175
1176

1177
1178 completely missing in the original generation results. There is still much space for further improving
1179 the instance-level sample quality for dataset distillation.
1180

1181 **More Sample Visualization** Additionally, we present more example samples across various categories
1182 to demonstrate the overall high quality of the dataset generated by the proposed CONCORD
1183 method. Specifically, in Fig. 11 we present samples generated for the Food-101 dataset. In Fig. 12
1184 images of animal categories in the ImageNet-1K dataset are generated by Minimax with CONCORD
1185 applied. In Fig. 13 we show images of other categories in the ImageNet-1K dataset. The high-quality
1186 generated samples form an effective surrogate dataset, which achieves state-of-the-art performance.
1187

1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205 station wagon birdhouse bookstore broom car mirror chest cradle dome
1206
1207
1208
1209
1210 fireboat goblet greenhouse home theater jeep lipstick manhole cover mitten
1211
1212
1213
1214
1215 mortar nail oscilloscope park bench printer pitcher poncho quilt
1216
1217
1218
1219
1220 rotisserie schooner ski sock sports car stone wall sunglasses teapot
1221
1222
1223
1224
1225 totem pole trolleybus vase whistle yurt ice cream broccoli custard apple
1226

1227 Figure 13: Example other images generated by the proposed CONCORD method on the ImageNet-
1228 1K dataset. The employed diffusion pipeline is the fine-tuned Minimax model. The class names are
1229 annotated below the images.

1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241