

FlashEval: Towards Fast and Accurate Evaluation of Text-to-image Diffusion Generative Models

Lin Zhao^{2*}, Tianchen Zhao^{12*}, Zinan Lin³, Xuefei Ning^{1†}, Guohao Dai²⁴, Huazhong Yang¹, Yu Wang^{1†}
¹Tsinghua University, ²Infinigence-AI, ³Microsoft Research, ⁴Shanghai Jiao Tong University
{l11zz0309zz, suozhang1998, linzinan1995, foxdoraame}@gmail.com daiguohao@sjtu.edu.cn
{yanghz, yu-wang}@mail.tsinghua.edu.cn

Abstract

In recent years, there has been significant progress in the development of text-to-image generative models. Evaluating the quality of the generative models is one essential step in the development process. Unfortunately, the evaluation process could consume a significant amount of computational resources, making the required periodic evaluation of model performance (e.g., monitoring training progress) impractical. Therefore, we seek to improve the evaluation efficiency by **selecting the representative subset of the text-image dataset**. We systematically investigate the design choices, including the selection criteria (textural features or image-based metrics) and the selection granularity (prompt-level or set-level). We find that the insights from prior work on subset selection for **training** data do not generalize to this problem, and we propose FlashEval, an iterative search algorithm tailored to **evaluation** data selection. We demonstrate the effectiveness of FlashEval on ranking diffusion models with various configurations, including architectures, quantization levels, and sampler schedules on COCO and DiffusionDB datasets. Our searched 50-item subset could achieve comparable evaluation quality to the randomly sampled 500-item subset for COCO annotations on unseen models, achieving a 10x evaluation speedup. We release the condensed subset of these commonly used datasets to help facilitate diffusion algorithm design and evaluation, and open-source FlashEval as a tool for condensing future datasets, accessible at <https://github.com/thu-nics/FlashEval>.

1. Introduction

Diffusion models could generate high-fidelity images based on textual descriptions (referred to as prompts) and achieve substantial research interest [23, 27, 28, 39]. Researchers have designed a range of techniques aimed at en-

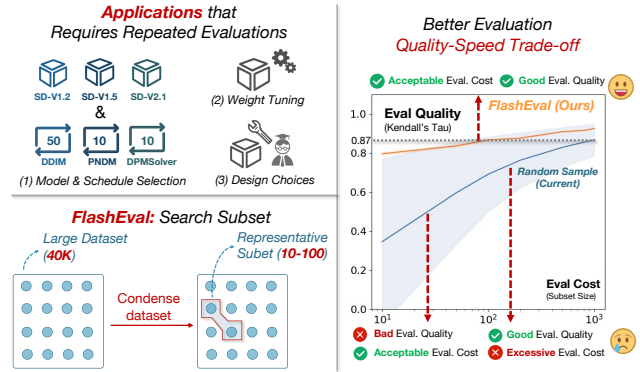


Figure 1. **The Motivation and Effectiveness of FlashEval.** Left: Applications and the search method. Right: The excessive evaluation cost and how FlashEval improves the evaluation speed-quality trade-off (evaluation quality is represented by the ranking correlation of a variety of diffusion models w.r.t subset size on COCO).

hancing these models from various angles. Some algorithm design applications often necessitate a large number of iterative evaluations, such as: **(1) Model & Schedule Selection:** The quality of image generation exhibits varying trends over different timesteps for different models and solvers. Consequently, evaluations and comparisons are needed to strike a balance between generation quality and computational cost. **(2) Weight Tuning:** In the process of training and fine-tuning the model, evaluating the performance of current parameters is required for refining the training strategy. **(3) Design Choices:** The iterative testing of multiple design options, such as layer-wise pruning rates, layer-wise mixed precision configurations, and other hyperparameter choices.

To evaluate the diffusion model, prior research [15, 21, 37] points out that the denoising loss does not necessarily reflect the generation quality. Therefore, existing evaluation schemes estimate the image quality (e.g., FID [14], IS [30], ImageReward [39]) over text annotations of existing text-image datasets [20, 36] or manually-designed textual dataset [16] (we refer them as “test prompts”). However, as

*Equal contribution

†Corresponding authors: Xuefei Ning and Yu Wang.

test prompts have a relatively large size (e.g., COCO [20] 47K, PicScore [16] 15K), the evaluation process has a significant computational burden.¹ Iterative evaluations required in the model design phase are thus impractical, and full set evaluation is often only adopted in the final evaluation (e.g., in DALL-E [26] and Stable-Diffusion [27]). To circumvent this challenge, the common practice is to **randomly sample a smaller subset** (e.g., 1K prompts in [7, 21]) for a proxy evaluation. However, we find that such random sub-sampling exhibits a poor accuracy-efficiency trade-off. In Fig. 1, we use this approach to evaluate a model zoo with various schedules. Even with 1000 samples, this approach only achieves a rank correlation of 0.87 (value of 1 indicates perfect correlation). More aggressive sub-sampling (for more speed-up) would result in an acute drop in ranking correlation and a large variation of results.

In this paper, **we aim to improve the accuracy-efficiency trade-off of diffusion evaluation by seeking “representative subsets” of the test prompts.** Firstly, we systematically explore the potential characteristics in representative sets. Then, we design baseline search methods to validate these characteristics and delve into the reasons for their failure when dealing with small item sizes. We conclude that the key issue leading to suboptimal search results is the inadequate exploration of the vast search space. Finally, inspired by the evolutionary algorithm, we design an improved iterative search algorithm on both the set and prompt levels. The proposed FlashEval search method could effectively mitigate variance and improve evaluation quality under small item sizes. Our acquired 50-item subset demonstrates a ranking correlation comparable to that of a 200-item baseline-acquired subset, and a randomly sampled 500-item subset. We will release the condensed subsets of different sizes for the commonly used datasets to assist researchers in selecting the appropriate prompt quantity for evaluation. Additionally, FlashEval could be used to search condensed evaluation subsets for new datasets.

In summary, our contributions can be listed as follows:

- We identify the need to improve the accuracy-efficiency trade-off of diffusion model evaluation and introduce the idea of **identifying a representative subset** to speed up evaluation.
- We systematically investigate the potential properties of representative sets and propose an improved search method to effectively acquire representative subsets even with small item size.
- We aim not to replace but to provide an accurate proxy of the existing evaluation scheme. We hope that FlashEval could assist researchers in selecting the proper

prompt quantity, and help accelerate and facilitate the broader diffusion algorithm design.

2. Related Work

Metrics for Text-to-Image Generation. The currently widely used metrics are often on one or more of the following aspects: image fidelity, image-text alignment, and image aesthetic. To assess fidelity, FID [14] and IS [30] are commonly used metrics that measure the feature distance between generated images and real images. CLIPScore [13] is widely used to measure image-text alignment. It calculates the similarity of features extracted from the image and text domains by dedicated encoders, respectively. As for aesthetic, Schuhmann *et al.* [31] train a score predictor using a dataset dedicated to aesthetics to facilitate evaluations.

To better align metric scores with human preference, recent methods perform manual scoring for generated images and train a predictor capable of evaluating models. PickScore [16] and HPS [38] generate two or more images for a given prompt, enabling users to choose the best one. ImageReward [39] involves users comprehensively considering both alignment and fidelity to assign explicit scores to each image. We apply metrics for all the aspects mentioned above, validating that FlashEval can be efficiently utilized to assess various aspects of models.

Diffusion Model Evaluation Datasets. To comprehensively evaluate Text-to-Image (T2I) models, there are currently various datasets available. Among them, the COCO validation set [20] is one of the most widely used datasets [7, 11, 12, 23]. In addition, to enhance benchmark precision, Saharia *et al.* [29] conduct the first attempt to consider multiple evaluation aspects and propose DrawBench, which contains 200 prompts specifically designed for evaluating T2I tasks. Cho *et al.* [8] introduce PAINTSKILLS, an evaluative dataset crafted to assess three essential visual reasoning skills: object recognition, object counting, and spatial relationships. Kirstain *et al.* [16] create a web app to build a large, open dataset of T2I prompts and real users’ preferences over generated images. Chen *et al.* [7] utilize GPT to provide scores and more detailed evaluative descriptions.

Recently, Bakr *et al.* [3] categorize different metrics and design distinct datasets for each category, aiding users in evaluating specific metrics more effectively. Similarly, Lee *et al.* [18] craft different scenarios and propose various datasets for each scenario. For any aforementioned datasets that are too large to evaluate efficiently, FlashEval can be applied to enhance the efficiency of evaluating models.

Training/Test Set Selection. There is a category of methods called "coreset selection" that accelerates the model training process by choosing subset from the training data, which aim to retain challenging examples in the subset [4, 22, 40]. Depending on the implementation, they can be categorized into two types: supervised learning [10, 24, 34] and unsuper-

¹For example, evaluating stable diffusion V1.5 on the whole COCO dataset requires 60 GPU hours (RTX 3090).

vised learning [6, 32, 33]. For supervised learning methods, it is common to assess each example during the training process to decide whether it should be included in the subset. Toneva *et al.* [34] track how often each example is forgotten during training. They assume that those forgotten less often have a smaller impact on the final results. Paul *et al.* [24] believe examples with high training loss should be retained, which indicate difficulty in learning. They train for only a few epochs for filtering, aiming to reduce the cost of subset selection. Unsupervised methods often select subsets based on clustering results. Ozan *et al.* [33] argue that examples farther from the centroid present greater learning challenges and should be retained. Sorscher *et al.* [32] use the K-center to find a diverse cover for active learning.

There are also some approaches that intend to reduce testing cost and do unbiased evaluation with unlabeled test dataset. Kossen *et al.* [17] actively select examples to be labeled by learning a surrogate model that directly predict losses at all data. Deng [9] obtain the transformed images and labels from training data to train a regression model for predicting scores of unlabeled data. Unlike them, our method wants to find the subset that are representative of the entire validation dataset. Recent work by [35] attempts to obtain the subset for evaluating language models. To best of our knowledge, we are the first to identify representative subset to accelerate evaluation for T2I models.

3. Problem Formulation

In this section, we provide a brief overview of the problem formulation for “**identifying a representative subset**” that has a good accuracy-efficiency trade-off. We choose the application scenario of ranking a variety of models (which is typically the case in the examples in Sec. 1), and use the **ranking correlation** for measuring the evaluation quality. More specifically, we consider a set of models N_m models: $\mathcal{M} = \{M_i | i \in \{1, \dots, N_m\}\}$ (with diverse model architectures, parameters, solvers, and/or step sizes). Given a dataset $\mathcal{P} = \{p_i | i \in \{1, \dots, N\}\}$ of size N , we aim to find a representative subset $\hat{\mathcal{P}} = \{p_i | i \in \mathcal{I}\}$ of item size N' with index list $\mathcal{I} = \{i_1, \dots, i_{N'}\}$ that maximizes the evaluation quality. For single-image evaluation metrics (e.g., CLIPScore, HPS, ImageReward), the metric score of prompt p_i for model M_j is described by \mathcal{S}_{ij} . We use the averaged metric scores on the entire textual dataset as the “ground-truth performance” of certain model $\hat{\mathcal{S}}_j = \sum_{i=1}^N \frac{1}{N} \mathcal{S}_{ij}$. For multiple-images metrics (e.g., FID), we use the full set evaluation as ground-truth. We adopt Kendall’s Tau (KD) [2] to measure the ranking correlation, which is computed as follows. Let’s denote the ground-truth metric and the approximated metric (computed using $\hat{\mathcal{P}}$) as x_1, \dots, x_{N_m} and y_1, \dots, y_{N_m} . For each pair of models M_i and M_j , we denote them as a *concordant* pair if $x_i < x_j$ and $y_i < y_j$, or $x_i > x_j$ and $y_i > y_j$ (i.e., they have a consistent ranking). Similarly, we denote them as

a *discordant* pair if $x_i < x_j$ and $y_i > y_j$, or $x_i > x_j$ and $y_i < y_j$. Among all pairs of models, we count the number of concordant and discordant pairs as N_c and N_d , respectively. Furthermore, we define n_1 as the number of pairs that have a tie only in the ground-truth metric (i.e., $x_i = x_j$ and $y_i \neq y_j$), and n_2 as the number of pairs that have a tie only in the approximated metric. Here, metric equivalence is deemed if the difference between values is smaller than a certain threshold, which is determined by the 3 times standard deviation across the scores of 50 randomly generated images from the same model. Given the values of N_c, N_d, n_1 , and n_2 , KD is computed as

$$\frac{N_c - N_d}{\sqrt{(N_c + N_d + n_1)(N_c + N_d + n_2)}}, \quad (1)$$

4. Methods

In T2I tasks, we can identify subsets from both the textual feature space and image metric space. In this section, we systematically investigate the potential properties of the representative subsets in the two spaces. Then, we design baseline methods to use these properties and analyze the reasons for their failure under smaller item sizes. Finally, we present an improved search algorithm that addresses such issues in Sec. 4.2.

4.1. Baseline Approaches for Subset Selection

A straightforward way of identifying the representative subset is by seeking samples that cover the entire dataset in terms of categories (**Baseline 1**) or distribution (**Baseline 2**) in the **textual feature space** following previous methods [32, 41]. However, abundant experiments in the appendix show that relying solely on textual features is not adequate for identifying the representative subset.

Therefore, we turn to examine the characteristics of the direct metric scores of generated images (**image-based metrics**). We propose to directly optimize the ranking correlation, and design search-based algorithms on the set and prompt level and validate their effectiveness on relatively medium-sized subsets. However, for smaller subset sizes like 10-50, their identified subsets exhibit significant variance and may yield poor performance. We further analyze the reasons for their failure.

Baseline 3: KD-based random search algorithm. To identify a representative set, we introduce a straightforward KD-based random search algorithm. There are two units for searching—set-wise search and prompt-wise search. For **set-wise search**, the fundamental search unit is “set”. The process involves sampling a large number of (e.g., 1 million) **subsets** randomly from the entire dataset \mathcal{P} , and then identifying the subset $\hat{\mathcal{P}}$ that yields best KD. For **prompt-wise search**, the fundamental search unit is “prompt”. We find the top K **prompts** with the highest KD to identify the representative subset with K prompts. We randomly split the

Methods	Type	$N' = 10$		$N' = 100$	
		Train	Test	Train	Test
RS	-	0.355	0.346	0.706	0.692
B3-set	Image Metrics	0.827	0.603	0.925	0.800
B3-prompt		0.822	0.610	0.914	0.764
Ours		0.934	0.796	0.966	0.867

Table 1. Comparisons of the Kendall’s Tau of CLIP-Score for subsets acquired by FlashEval and baseline methods. The “RS” represents sampling N' prompts randomly from dataset.

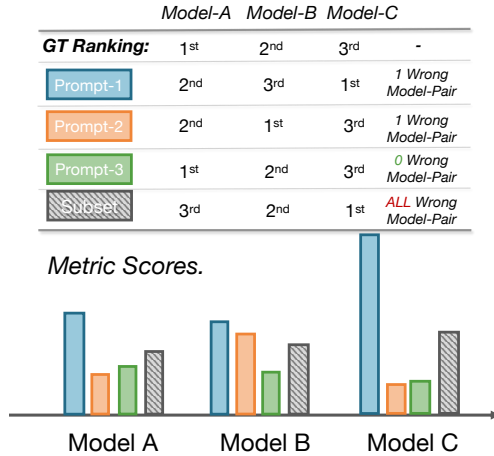


Figure 2. Illustration of the reasons for baseline-3’s failing under small item sizes N' . When combining multiple prompts with standalone high KD, the set-wise KD is lower.

model zoo to be ranked as training/testing, and adopt the above mentioned search on the training models.

Results for KD-based search. As could be seen from Tab. 1, the results of prompt-wise search (brief for “B3-prompt”) and set-wise search (brief for “B3-set”) achieve better results than RS. When $N' = 100$, the searched subset achieves a KD of 0.75-0.8 on test set, whereas when $N' = 10$, the KD is worse (≈ 0.6).

Analysis of why B3-prompt failed. For “B3-prompt”, we conduct detailed analysis for searched prompt and composed set. When combining multiple prompts with standalone high KD, the set-wise ranking KD notably decreases. As shown from the toy examples in Fig. 2: the prompt-1,2,3 have high standalone KD (rank one or none of the models wrong), but the combined subset have low KD (rank all model wrong). It reveals that searching the top prompts and combining them may not produce the optimal subset.

Analysis of why B3-set failed. Based on the analysis above, the set-wise search might be a better approach because it takes the interference between prompts into consideration. However, the improvement of set-wise search over prompt-wise search is marginal even on the training set as shown in Tab. 1. We attribute such a phenomenon to the insufficient exploration of the search space. Take $N' = 10$ as example, the complete search space of constructing 10

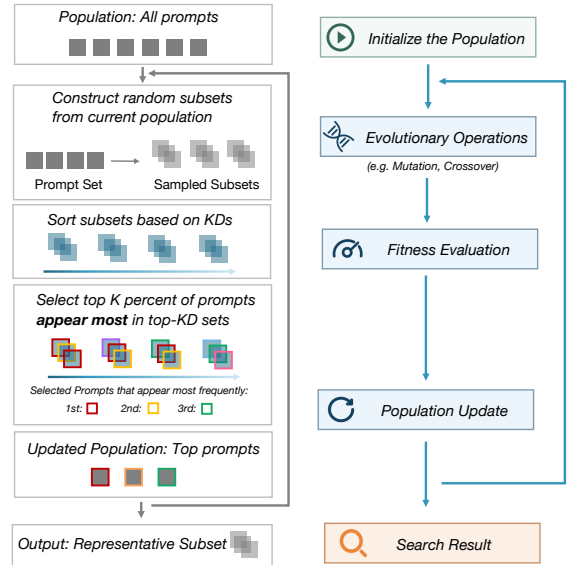


Figure 3. Illustration of FlashEval search method. Inspired by the evolutionary algorithm, we design an iterative search algorithm on both the set- and prompt-level.

prompts from 40,504 text annotations (COCO dataset), the search space size is C_{40504}^{10} , which is impractical to traverse through. Even we construct numerous subsets (i.e. 1M), it is still restricted compared with the enormous search space size. Therefore, an improved search algorithm with enhanced sample efficiency is required to generate representative subsets.

4.2. FlashEval Search Method

As discussed above, the prompt-wise search is efficient but results in suboptimal subsets because it doesn’t consider interactions between prompts. The set-wise search considers the interactions but suffers from insufficient exploration of enormous search space. The key problem is to improve the sample efficiency of set-wise searching. In light of the prior that well-performing prompts are more likely to construct well-performing sets, we propose to utilize the prompt-level information to provide guidance for the set-level search. Inspired by the evolutionary algorithms [1, 5], we design an iterative search method on both the prompt and set level. We provide an overall introduction to the workflow of our method in Sec. 4.2.1. Then, we elaborate on our frequency-based filtering method in Sec. 4.2.2. The detailed description of FlashEval search algorithm flow is listed in Algo. 1.

4.2.1 Set- and Prompt-level Iterative Filtering

The flowchart of the typical evolutionary algorithm and our FlashEval search method is presented in Fig. 3. Firstly, an initial *population* containing candidate solutions is generated. Then, the “recombination” process employs *evolutionary operators* (e.g., mutation, crossover) to generate new individuals (offspring) from existing ones, which will be the new population. After that, the *fitness* of each individual

Algorithm 1: FlashEval Search Algorithm

Input: Whole textual dataset $\mathcal{P} = \{p_i | i = 1, \dots, N\}$, The train split of the model zoo \mathcal{M}_{train} , the number of candidate subsets to rank N_{set} , the number of iterations N_{iter} , keep ratio of prompts K_p , keep ratio of sets K_s ;
Output: Representative subsets $\hat{\mathcal{P}}$ of size N' ;
Get the model ranking R_{gt} on \mathcal{P} ;
// Initialize the population
Initialize $\mathcal{P}' \leftarrow \mathcal{P}$;
for $i \leftarrow 1$ to N_{iter} **do**
 for $j \leftarrow 1$ to N_{set} **do**
 // Evolutionary Operations
 Construct random subset $\mathcal{P}_j \sim \mathcal{P}'$ of size N' , get its model ranking R_j on \mathcal{M}_{train} ;
 Calculate KD_j values between R_j and R_{gt} ;
 // Fitness Evaluation
 $S_{set} \leftarrow$ subsets among $\{\mathcal{P}_j\}_{j=1}^{N_{set}}$ with top K_s percent of KD ;
 // Population update
 $S_{prompt} \leftarrow$ top K_p percent of prompts that appear most frequently in S_{set} ;
 Update candidate prompt set $\mathcal{P}' \leftarrow S_{prompt}$;
Construct subsets $\{\mathcal{P}_j\}_{j=1}^{N_{set}}$ of size N' from \mathcal{P}' ;
return $\hat{\mathcal{P}}$ with best KD in $\{\mathcal{P}_j\}_{j=1}^{N_{set}}$

in the population is evaluated based on a pre-defined *objective*. Furthermore, based on the fitness values of each individual, a *selection* is applied to filter subpar individuals to acquire the updated population. Such process is repeated until convergence and the last population is acquired.

Inspired by the process of the evolutionary method, we design an iterative search on both the set- and prompt-level as presented in Alg. 1. We treat the well-performing prompts as the “population”, and use the entire text annotations as initialization. For each iteration, we construct a large number of subsets with the prompts within the population \mathcal{P}' . It could be viewed as the “recombination” process. Then, the generated subsets (offspring) are evaluated using the KD value as the “fitness”. A novel frequency-based prompt “selection” scheme is designed to effectively distinguish promising prompts. The population is updated after the selection process and used in the next iteration. The search algorithm progressively filters the redundant subpar prompts and significantly enhances the search space exploration efficiency, and acquires top-performing subsets with low variance.

4.2.2 Frequency-based Prompt Selection

An effective “population update” (or prompt selection) is crucial for the success of the search method. As shown in Fig. 2, the selected prompts should avoid such corner cases whose standalone KD is high, but interfere with each other, leading to poorly performing subsets. We aim to seek prompts that

satisfy the property of “**when combined with any other prompts within the population, it has a higher probability of achieving high KD** ”. With the intention of that, we design a frequency-based prompt selection method that selects prompts that appear most frequently in top-performing subsets. In the top-performing sets from a large variety of constructed subsets, the most frequently appearing prompts are likely to satisfy the desired property. The prompts selected this way exhibit high standalone KD , and retain high KD when randomly combined.

5. Experiments

5.1. Experimental Details

Dataset. We conducted our experiments on two widely used large-scale datasets: COCO validation dataset (i.e., *val2014*) [20] and DiffuisonDB dataset [36]. In COCO, each image corresponds to several prompts. We randomly select one for each image to result in 40,504 prompts. DiffuisonDB contains 2 million prompts, divided into 2,000 parts. In each part, we uniformly select samples, resulting in a total of 5,000 prompts for our experiments.

Metric setting. We adopt five popular evaluation metrics. Specifically, FID [14], CLIP [25], and Aesthetic [31] are employed to assess three distinct aspects including image fidelity, alignment, and aesthetic. ImageReward [39] and HPS [38]) are utilized for evaluating human preference.

Model zoo setting. We construct 78 **model configurations** from the StableDiffusion Models [27] to form the model zoo by considering multiple perspectives. The model zoo comprises 12 **variant models** (v1.2, v1.2-6bit quant, v1.2-8bit quant, v1.4, v1.5, small-v1.5, v2.1, dreamlike-photoreal, v1.4-6bit quant, v1.4-8bit quant, v1.5-6bit quant, v1.5-8bit quant) and 8 **schedules** (DDIM-10step, DDIM-20step, DDIM-50step, PNDM-10step, PNDM-20step, PNDM-50step, DPM-10step, DPM-20step). Quantization models are based on the q-diffusion [19] framework, which encompasses both the DDIM and DPM solvers. For each model, we maintain consistency by fixing seed=0 and guidance scale=7.5 to ensure fairness in ranking across models. This model zoo covers various scenarios discussed in Sec. 1 where users want to evaluate and compare diffusion models, including (1) parameter tuning, such as tuning training iterations/datasets (v1.2, v1.4, v1.5, and dreamlike-photoreal are from different training iterations/datasets of the same model), (2) schedule selection, and (3) validating design choices (e.g., quantization).

5.2. Comparison to Baseline Methods

We design three sub-tasks based on the composition of different models to validate the effectiveness of FlashEval across various settings. Furthermore, to ensure applicability in diverse scenarios, the representative subset needs to

models	item size methods \sub-tasks	$N'=50$			$N'=500$		
		random	model variants	schedulers	random	model variants	schedulers
Train	RS	0.607±0.000	0.594±0.000	0.632±0.000	0.857±0.000	0.858±0.000	0.857±0.000
	B3-prompt	0.900	0.909	0.862	0.895	0.917	0.872
	B3-set	0.895±0.002	0.912±0.002	0.894±0.002	0.971±0.002	0.970±0.001	0.966±0.002
	Ours	0.956±0.004	0.969±0.004	0.960±0.004	0.984±0.003	0.986±0.003	0.978±0.003
Test	RS	0.597±0.000	0.588±0.000	0.560±0.000	0.829±0.000	0.826±0.000	0.827±0.000
	B3-prompt	0.729	0.784	0.810	0.805	0.822	0.851
	B3-set	0.750±0.014	0.680±0.021	0.721±0.013	0.875±0.007	0.836±0.008	0.863±0.008
	Ours	0.851±0.004	0.800±0.008	0.850±0.005	0.906±0.003	0.899±0.004	0.909±0.003

Table 2. Comparisons of the Kendall’s Tau of CLIP-Score on COCO dataset (the error bars are standard errors). B3-prompt only requires finding prompts once, hence there’s no need to calculate standard errors. “Random”, “model variants”, “schedulers” are the three sub-tasks mentioned in 5.2. The “RS” represents sampling N' prompts randomly from dataset.

models	item size methods \sub-tasks	$N'=5$			$N'=100$		
		random	model variants	schedulers	random	model variants	schedulers
Train	RS	0.395±0.000	0.307±0.000	0.413±0.000	0.730±0.000	0.721±0.000	0.755±0.000
	B3-prompt	0.819	0.725	0.820	0.879	0.843	0.920
	B3-set	0.839±0.003	0.827±0.003	0.833±0.002	0.935±0.002	0.924±0.001	0.933±0.002
	Ours	0.908±0.004	0.891±0.004	0.897±0.006	0.966±0.002	0.968±0.002	0.967±0.003
Test	RS	0.350±0.000	0.423±0.000	0.323±0.000	0.763±0.000	0.681±0.000	0.738±0.000
	B3-prompt	0.520	0.266	0.657	0.752	0.683	0.806
	B3-set	0.643±0.031	0.516±0.033	0.601±0.018	0.826±0.010	0.792±0.018	0.836±0.014
	Ours	0.789±0.009	0.691±0.011	0.748±0.014	0.851±0.007	0.861±0.007	0.873±0.007

Table 3. Comparisons of the Kendall’s Tau of CLIP-Score on DiffusionDB dataset (the error bars are standard errors).

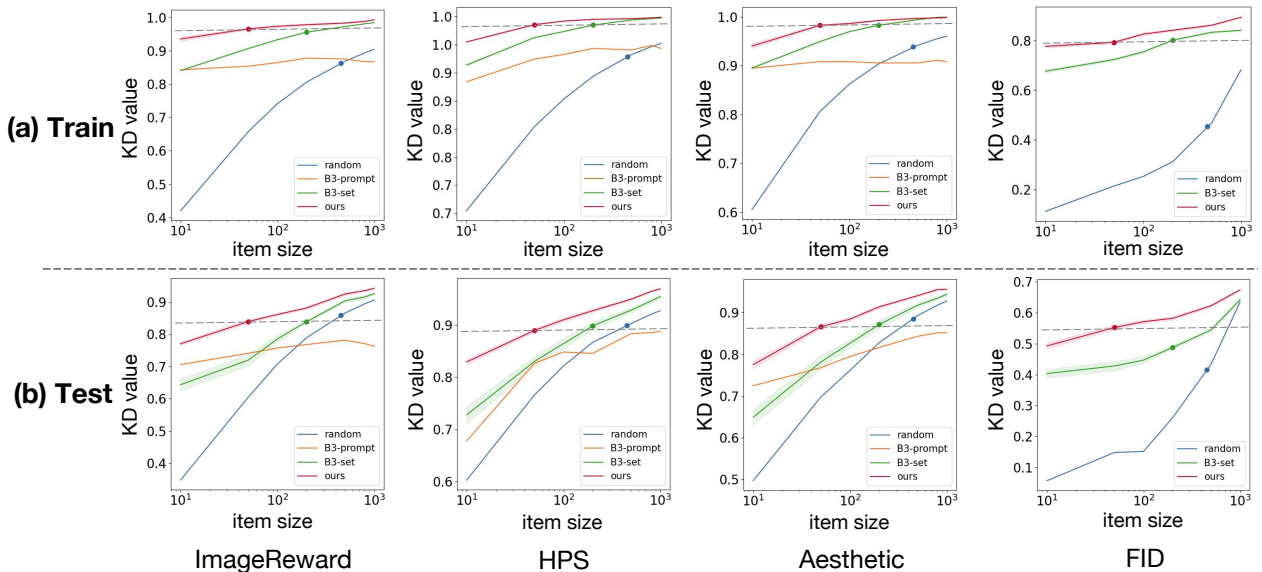


Figure 4. Comparisons with Baseline3 of the Kendall’s Tau of different metrics on COCO dataset for “Random” sub-task. The red, green, and blue dots respectively represent ours with 100 items, B3-set with 200 items, and random sample with 500 items. The shaded area are standard errors. The complete comparison results on the two datasets are in the appendix.

possess the ability to generalize across different models. In this regard, we split the model zoo with different partitioning strategies to obtain training models for searching and testing models for assessing generalization capabilities for

each sub-task. The three partitioning ways are as follows: **random split**, **across model variances** (different architectures and parameters for training and testing models), **across schedules** (different solvers and step sizes for training and

testing models). The detailed strategies are in the appendix.

FastEval finds subsets with better accuracy-efficiency trade-off. We conduct comparisons with the previously mentioned baseline approaches on each metric for every sub-task. For the sake of fair comparison, we ensure that the search space for each method, i.e., the total number of candidate subsets, remains consistent ($N = 10^7$). In addition, each set-based search algorithm is executed 10 times to compute standard errors. Due to FID is only applicable to the evaluation of two distributions and cannot evaluate a single image, it does not support prompt-based searches. Therefore, we just compare FlashEval with set-based methods for FID.

From Tab. 2 and Tab. 3, we can observe that FlashEval can easily achieve KD values around 0.9 on training models when the number of prompts is relatively low. FlashEval surpasses all methods with equal search space, particularly evident when the item size is small. FlashEval at $N' = 100$ performs comparably to B3-set at $N' = 200$ (Red and green dots in Fig. 4-(a)) across all metrics. For prompt-based method, as the number of items in the subset increases, the KD values do not improve, which validates our hypothesis in Sec. 4.1 that prompt-based method is not suitable for searching representative subset. In addition, Fig. 4-(a) shows that the commonly used random sampling strategy’s performance drops notably below 500 items (especially on FID).

FlashEval generalizes across diverse model setting. We aim for users to utilize the representative set directly on the new models. Thus, we test the generalization performance of FlashEval on **test models** across three sub-tasks. We find that FlashEval’s set with $N' = 5$ surpasses the average KD values of randomly selected sets with $N' = 100$ on DiffusionDB (Tab. 3). Moreover, the $N' = 50$ set identified by FlashEval is on par with the randomly selected $N' = 500$ set on COCO (and Tab. 2). As depicted in Fig. 4-(b), FlashEval exhibits superior generalization across all metrics compared to the baseline methods. Additionally, it notably excels when the representative set has a limited number of prompts. The dots in Fig. 4 also signify that our method achieves comparable performance with $N' = 50$ compared to B3-set’s $N' = 200$ and random sampling’s $N' = 500$ even for unseen models. Furthermore, the standard error of FlashEval is relatively smaller compared to B3-set, indicating its higher robustness and stability.

5.3. Use Cases of FlashEval

To validate the practicality of FlashEval, we present two examples of downstream use cases here.

Application A: Model selection from different versions. **Experiment settings:** We adopt a 10-item RS and FlashEval to rank 44 testing models, and the full set (40K) averaged ImageReward is used as ground truth. We illustrate the superiority of FlashEval using 3 variants of SD1.5 (with DDIM solver & 20 steps) in testing models: *dreamlike* (fine-tuned),

set	small-sd1.5	weight-8bit	dreamlike
RS-10 item (KD=0.286)	20	28	25
Ours-10 item (KD=0.761)	34	23	3
Ground Truth	39	25	4

Table 4. The ranking of three models among 44 testing models.

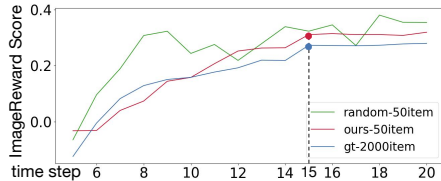


Figure 5. The evolution of model scores (sd1.5, DPM solver) with increasing steps.

small-sd1.5 (reduce params), and *weight-8bit* (quantized). **Result:** Tab. 4 provides specific rankings for 3 models. FlashEval with significantly higher KD (0.761) not only reasonably predicts that the dreamlike is better than *small-sd1.5* but also accurately ranks the two efficient models (*small-sd1.5*, *weight-8bit*) that are not intuitively comparable. However, RS failed in both scenarios.

Application B: Given a model & solver, the performance increase saturates as steps increase. Identifying the “saturate point” is vital for striking good performance-efficiency trade-off. **Experiment settings:** We search the “saturate point” by observing that the score remains relatively unchanged in several steps. **Result:** Fig. 5 shows that FlashEval identifies an accurate saturation point (15) matching ground truth (15). In contrast, RS fails to find a saturation point in 20 steps due to continuous performance fluctuations.

5.4. Analysis

5.4.1 Ablation Studies

Effectiveness of frequency-based selection. To validate the impact of frequency-based selection, the algorithm is modified by selecting prompts with higher KD values, instead of choosing prompts based on their frequency.

As shown in Tab. 5, the results are worse than FlashEval, especially on testing models. In addition, the gap is larger when $N' = 10$, further illustrating the effectiveness of the frequency-based selection module when N' is small.

Effectiveness of iterative filtering. As shown in Tab. 5, our algorithm demonstrates significantly superior performance on both the training and testing models even without iterative filtering compared to the average performance achieved through random selection. In addition, from Fig. 6, it is evident that with an increase in the number of iterations, the KD values of the representative set steadily increase. In each iteration, the KD values of the filtered set are consistently lower than those of selected set. The results validate our algorithm’s iterative process of gradually filtering out

Methods	N'=10		N'=100	
	Train	Test	Train	Test
RS	0.355	0.142	0.706	0.692
Ours (w/o FS)	0.876	0.640	0.953	0.760
Ours (w/o IF)	0.821	0.667	0.839	0.812
Ours	0.934	0.796	0.966	0.867

Table 5. **Comparison of ablation experiments of CLIP-score on COCO dataset.** ‘‘Ours (w/o FS)’’ is our algorithm without frequency-based selection, ‘‘Ours (w/o IF)’’ is our algorithm without iteration filtering.

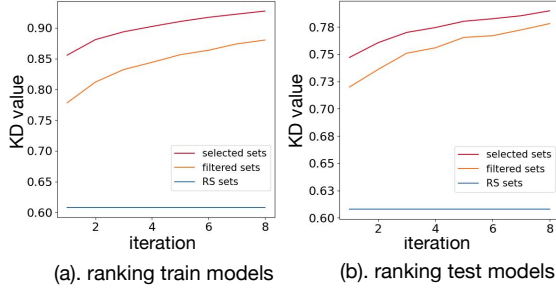


Figure 6. **The change in KD values as the algorithm progresses through iterations** ($N' = 50$, CLIP-Score, COCO, random split).

relatively inferior prompts in each iteration.

5.4.2 Consistent Superiority of FlashEval

A high KD value may arise from the correct ranking of only some specific positions (e.g. only the ranking of bottom models are correct). In such a scenario, a high KD value is irrelevant for users aiming to select the top model. We add two criteria to verify that FlashEval consistently outperforms RS, rather than focusing on correct ranking at specific positions: (1) **Top- K Ranking**: KD of the true top- K models. (2) **Top- K Proportion**: the proportion of true top- K models in predicted top- K models. Tab. 6 shows that FlashEval consistently outperforms RS across K s.

5.4.3 Search Cost of FlashEval

The search of FlashEval is a one-time cost; the subset found by FlashEval can be used in all evaluations of the same dataset. Here, we aim to further reduce this one-time search cost for identifying the representative set for a new dataset.

Assume that user has a new dataset of size N_{all} , and we aim to identify a representative set N' of N_{all} . The main computational cost of FlashEval is to evaluate the scores of N_{all} samples against all models in the model zoo. As shown in Sec. 5.2, random sampling shows poor accuracy with small item size. However, when the number of items is large enough, it is still a cheap approach to achieve reasonable accuracy. Inspired by this observation, we propose the following two-step approach: To find N' , we first randomly select a N_{eval} subset ($N' < N_{eval} < N_{all}$) and treat it as

Top- K	Top- K Ranking		Top- K Proportion	
	RS	Ours	RS	Ours
$K = 5$	0.400	0.800	40%	80%
$K = 10$	0.520	0.810	50%	80%
$K = 20$	0.458	0.793	55%	95%

Table 6. **Results of 10-item subsets obtained by RS and Ours for various K values.** The total number of testing models is 39.

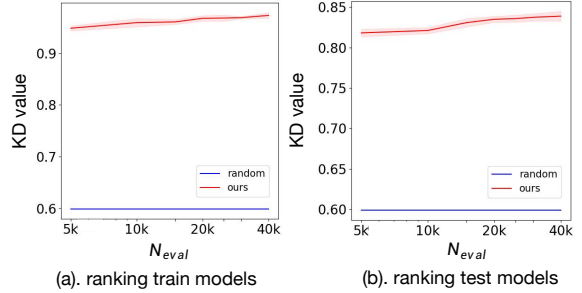


Figure 7. **The variation of KD values with the increase in N_{eval}** ($N' = 50$, CLIP-Score, COCO, random split).

the ground-truth when running FlashEval. This approach reduces the one-time cost from $\mathcal{O}(N_{all})$ to $\mathcal{O}(N_{eval})$. N_{eval} is the critical parameter that controls the trade-off between search efficiency and accuracy. On one hand, N_{eval} should not be too small, as otherwise, it is not representative enough for the true ranking of models (Sec. 5.2). On the other hand, when N_{eval} is larger, this approach provides less speed-up compared to the vanilla FlashEval.

To validate the searched set’s effectiveness, we compare the ranking results of N' with the rankings of N_{all} . As depicted in Fig. 7, when reducing N_{eval} to 5000, there is no significant decrease in the results. In addition, the trade-off in Fig. 7 illustrates the relationship between search cost and performance, aiding users in making their own choices.

6. Conclusion

In this paper, we identify the excessive cost and suboptimal quality-speed trade-off of the current diffusion model evaluation. We propose to address the issues by identifying representative subsets to condense the evaluation datasets and designing a search method to solve the challenging task. Extensive experiments demonstrate the effectiveness and generalization capability of the representative sets identified by FlashEval. We hope that FlashEval can assist researchers in selecting an appropriate prompt quantity in diffusion model evaluation, and help accelerate and facilitate the broader diffusion algorithm design.

Acknowledgements. This work was supported by National Natural Science Foundation of China (No. 62325405, 62104128, U19B2019, U21B2031, 61832007, 62204164), Tsinghua EE Xilinx AI Research Fund, and Beijing National Research Center for Information Science and Technology (BNRist). We thank for all the support from Infinigence-AI.

References

- [1] Evolutionary algorithm. https://en.wikipedia.org/wiki/Evolutionary_algorithm. 4
- [2] Hervé Abdi. The kendall rank correlation coefficient. *Encyclopedia of Measurement and Statistics*. Sage, Thousand Oaks, CA, pages 508–510, 2007. 3
- [3] Eslam Mohamed Bakr, Pengzhan Sun, Xiaogian Shen, Faizan Farooq Khan, Li Erran Li, and Mohamed Elhoseiny. Hrs-bench: Holistic, reliable and scalable benchmark for text-to-image models. In *ICCV*, pages 20041–20053, 2023. 2
- [4] Zalán Borsos, Mojmír Mutný, and Andreas Krause. Coresets via bilevel optimization for continual learning and streaming. *NeurIPS*, 33:14879–14890, 2020. 2
- [5] Thomas Bäck. *Evolutionary algorithms in theory and practice - evolution strategies, evolutionary programming, genetic algorithms*. Oxford University Press, 1996. 4
- [6] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. *NeurIPS*, 33:9912–9924, 2020. 3
- [7] Yixiong Chen. X-iqe: explainable image quality evaluation for text-to-image generation with visual large language models. *arXiv preprint arXiv:2305.10843*, 2023. 2
- [8] Jaemin Cho, Abhay Zala, and Mohit Bansal. Dall-eval: Probing the reasoning skills and social biases of text-to-image generation models. In *ICCV*, pages 3043–3054, 2023. 2
- [9] Weijian Deng and Liang Zheng. Are labels always necessary for classifier accuracy evaluation? In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 15069–15078, 2021. 3
- [10] Vitaly Feldman and Chiyuan Zhang. What neural networks memorize and why: Discovering the long tail via influence estimation. *NeurIPS*, 33:2881–2891, 2020. 2
- [11] Jiayi Guo, Chaoqun Du, Jiangshan Wang, Huijuan Huang, Pengfei Wan, and Gao Huang. Assessing a single image in reference-guided image synthesis. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 753–761, 2022. 2
- [12] Jiayi Guo, Chaofei Wang, You Wu, Eric Zhang, Kai Wang, Xingqian Xu, Humphrey Shi, Gao Huang, and Shiji Song. Zero-shot generative model adaptation via image-specific prompt learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11494–11503, 2023. 2
- [13] Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A reference-free evaluation metric for image captioning. *arXiv preprint arXiv:2104.08718*, 2021. 2
- [14] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *NeurIPS*, 30, 2017. 1, 2, 5
- [15] Dongjun Kim, Seung-Jae Shin, Kyungwoo Song, Wanmo Kang, and Il-Chul Moon. Soft truncation: A universal training technique of score-based diffusion model for high precision score estimation. In *ICML*, 2021. 1
- [16] Yuval Kirstain, Adam Polyak, Uriel Singer, Shahbuland Matiana, Joe Penna, and Omer Levy. Pick-a-pic: An open dataset of user preferences for text-to-image generation. *arXiv preprint arXiv:2305.01569*, 2023. 1, 2
- [17] Jannik Kossen, Sebastian Farquhar, Yarin Gal, and Thomas Rainforth. Active surrogate estimators: An active learning approach to label-efficient model evaluation. *Advances in Neural Information Processing Systems*, 35:24557–24570, 2022. 3
- [18] Tony Lee, Michihiro Yasunaga, Chenlin Meng, Yifan Mai, Joon Sung Park, Agrim Gupta, Yunzhi Zhang, Deepak Narayanan, Hannah Benita Teufel, Marco Bellagente, et al. Holistic evaluation of text-to-image models. *arXiv preprint arXiv:2311.04287*, 2023. 2
- [19] Xiuyu Li, Long Lian, Yijiang Liu, Huanrui Yang, Zhen Dong, Daniel Kang, Shanghang Zhang, and Kurt Keutzer. Q-diffusion: Quantizing diffusion models. *ICCV*, 2023. 5
- [20] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*, 2014. 1, 2, 5
- [21] Enshu Liu, Xuefei Ning, Zinan Lin, Huazhong Yang, and Yu Wang. Oms-dpm: Optimizing the model schedule for diffusion probabilistic models. *ICML*, 2023. 1, 2
- [22] Baharan Mirzasoleiman, Jeff Bilmes, and Jure Leskovec. Coresets for data-efficient training of machine learning models. In *ICML*, pages 6950–6960. PMLR, 2020. 2
- [23] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021. 1, 2
- [24] Mansheej Paul, Surya Ganguli, and Gintare Karolina Dziugaite. Deep learning on a data diet: Finding important examples early in training. *NeurIPS*, 34:20596–20607, 2021. 2, 3
- [25] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 5
- [26] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *ICML*, pages 8821–8831. PMLR, 2021. 2
- [27] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, pages 10684–10695, 2022. 1, 2, 5
- [28] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *NeurIPS*, 35:36479–36494, 2022. 1
- [29] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael

- Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *NeurIPS*, 35:36479–36494, 2022. 2
- [30] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training gans. *NeurIPS*, 29, 2016. 1, 2
- [31] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *NeurIPS*, 35:25278–25294, 2022. 2, 5
- [32] Ozan Sener and Silvio Savarese. Active learning for convolutional neural networks: A core-set approach. *ICLR*, 2018. 3
- [33] Ben Sorscher, Robert Geirhos, Shashank Shekhar, Surya Ganguli, and Ari Morcos. Beyond neural scaling laws: beating power law scaling via data pruning. *NeurIPS*, 35:19523–19536, 2022. 3
- [34] Mariya Toneva, Alessandro Sordani, Remi Tachet des Combes, Adam Trischler, Yoshua Bengio, and Geoffrey J Gordon. An empirical study of example forgetting during deep neural network learning. *ICLR*, 2018. 2, 3
- [35] Rajan Vivek, Kawin Ethayarajh, Diyi Yang, and Douwe Kiela. Anchor points: Benchmarking models with much fewer examples. *arXiv preprint arXiv:2309.08638*, 2023. 3
- [36] Zijie J. Wang, Evan Montoya, David Munechika, Haoyang Yang, Benjamin Hoover, and Duen Horng Chau. Diffusiondb: A large-scale prompt gallery dataset for text-to-image generative models. *ArXiv*, abs/2210.14896, 2022. 1, 5
- [37] Daniel Watson, Jonathan Ho, Mohammad Norouzi, and William Chan. Learning to efficiently sample from diffusion probabilistic models. *ArXiv*, abs/2106.03802, 2021. 1
- [38] Xiaoshi Wu, Keqiang Sun, Feng Zhu, Rui Zhao, and Hongsheng Li. Better aligning text-to-image models with human preference. *arXiv preprint arXiv:2303.14420*, 2023. 2, 5
- [39] Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao Dong. Imagereward: Learning and evaluating human preferences for text-to-image generation. *arXiv preprint arXiv:2304.05977*, 2023. 1, 2, 5
- [40] Zeyuan Yin, Eric Xing, and Zhiqiang Shen. Squeeze, recover and relabel: Dataset condensation at imagenet scale from a new perspective. *arXiv preprint arXiv:2306.13092*, 2023. 2
- [41] Daquan Zhou, Kai Wang, Jianyang Gu, Xiangyu Peng, Dongze Lian, Yifan Zhang, Yang You, and Jiashi Feng. Dataset quantization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 17205–17216, 2023. 3

Supplementary Materials

1. Detailed Description of Textual Feature-Based Baseline Methods

In this section, we supplement the implementation details of the textual feature space baseline methods (B1/B2) in Sec. 4.1 .

1.1. POS (Part of Speech) Tagging Analysis

The POS (Part of Speech) tagging is a process to mark up the words in text format for a particular part of a speech based on its definition and context. The “part of speech” is a typical low-level textual feature, which could potentially reflect the characteristic of the representative textual subset. We conduct the POS tagging analysis on COCO annotations. We use the NLTK’s (Nature Language Toolkit) tokenizer and POS tagging tool to analyze the POS distribution. We present the POS that occupies more than 1% in the legend. The “NN” represents “singular noun”, “DT” represents “determiner” (e.g., “my”), “IN” represents “Preposition and conjunction” (e.g., “at/in”), “JJ” represents “Adjective”, “.” represents the punctuation, “VBG/VBZ” represents different forms of verbs, “NNS” represents “Proper and plural noun”, and “CC” represents “conjunction of coordinating” (e.g., “that/which”). As demonstrated in Fig. 1, the POS distributions of the entire set and subsets of 100 items with high/low standalone KD are **quite similar**. Therefore, **it is not feasible to select representative subsets according to this characteristic**.

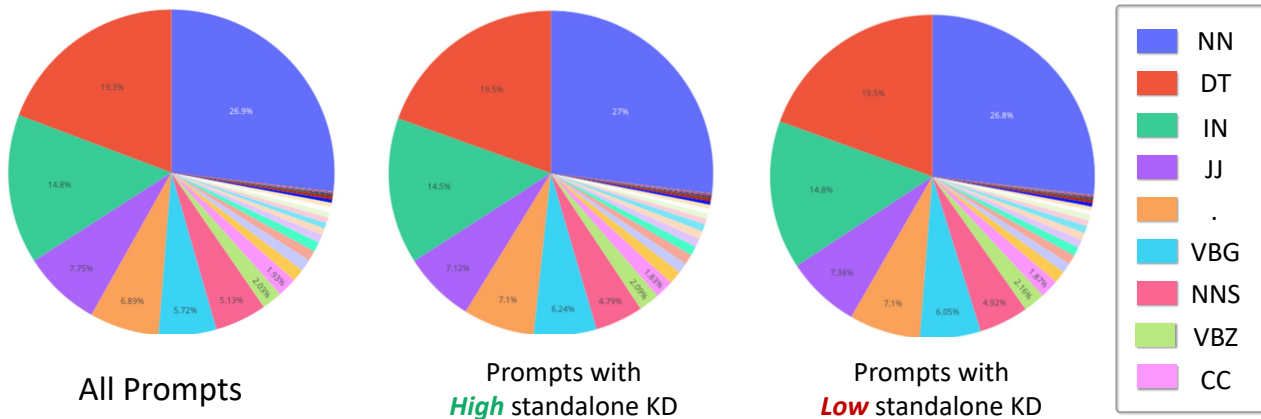


Figure 1. Comparison of pos tagging analysis for the entire COCO annotations, 100 item subsets with high/low standalone KD.

1.2. Semantic Class Diversity

To effectively illustrate the inability to select a representative subset through semantic class diversity, we conduct comprehensive experiments on this matter. We partition CLIP/BERT features into 10 clusters using the K-means algorithm, ensuring that prompts in each cluster encompass similar semantic information categories (e.g., humans, cats, etc). Subsequently, we selected $N/10$ prompts from each cluster based on the Euclidean distance to the cluster centroid (near or far), or through random selection. As depicted in Tab. 1 and Tab. 2, all the methods yields results similar to random sampling. Therefore, **ensuring semantic class diversity does not lead to an improvement in performance**.

Methods	$N' = 10$		$N' = 100$	
	Train	Test	Train	Test
RS	0.355	0.346	0.706	0.692
Near Centroid	0.455	0.451	0.474	0.570
Far Centroid	0.551	0.539	0.617	0.713
RS-each cluster	0.341	0.313	0.617	0.682

Table 1. Comparisons of the Kendall’s Tau of CLIP-Score for subsets acquired by clustering CLIP features.

Methods	$N' = 10$		$N' = 100$	
	Train	Test	Train	Test
RS	0.355	0.346	0.706	0.692
Near Centroid	0.057	0.335	0.490	0.540
Far Centroid	0.285	-0.116	0.659	0.545
RS-each cluster	0.228	0.251	0.686	0.651

Table 2. Comparisons of the Kendall’s Tau of CLIP-Score for subsets acquired by clustering Bert features.

Algorithm 1: Greedy Search Algorithm

Input: Whole textual dataset $\mathcal{P} = \{p_i | i = 1, \dots, N\}$, The corresponding features of dataset $\mathcal{F} = \{f_i | i = 1, \dots, N\}$, The mean μ_{all} and covariance Σ_{all} of multivariate Gaussian distribution of \mathcal{F} ;

Output: Subsets $\hat{\mathcal{P}}$ of size N' ;

Random Initialize $\hat{\mathcal{P}} \leftarrow p_k, \hat{\mathcal{F}} \leftarrow f_k (1 \leq k \leq N)$;

for $i \leftarrow 2$ **to** N' **do**

$index = INF, D_{min} = INF$;

for f_j **in** $\mathcal{F} - \hat{\mathcal{F}}$ **do**

$\mathcal{F}_j = \hat{\mathcal{F}} \cup f_j$;

 Calculate μ_j and Σ_j of \mathcal{F}_j ;

 // KL distance between \mathcal{F}_j and \mathcal{F}

$D_j = KL(\mu_j, \Sigma_j, \mu_{all}, \Sigma_{all})$

if $D_{min} > D_j$ **then**

$index = j, D_{min} = D_j$;

$\hat{\mathcal{P}} = \hat{\mathcal{P}} \cup p_{index}, \hat{\mathcal{F}} = \hat{\mathcal{F}} \cup f_{index}$;

return $\hat{\mathcal{P}}$

1.3. Statistical Property Similarity

As shown in Algo. 1, we design a greedy search algorithm that selects the prompt that minimizes the KL divergence between the CLIP feature distribution of the subset and the entire set in each iteration. Since the algorithm requires computing the KL divergence, the feature distribution is assumed to be a multivariate Gaussian distribution. To better fit the feature distribution, we expand the Gaussian distribution into the Gaussian mixture distribution. It can be depicted in Fig. 2, when the Gaussian Mixture distribution contains five Gaussian functions, it effectively fits the features. Then, we run the algorithm in each Gaussian function, identifying $N'/5$ prompts in each function (a total of N' prompts). We show the results of different items in Tab. 3. **It was evident that matching the statistical property similarity does not improve the ranking outcome.**

Methods	$N' = 100$		$N' = 200$		$N' = 500$		$N' = 1000$	
	Train	Test	Train	Test	Train	Test	Train	Test
RS	0.706	0.692	0.784	0.763	0.857	0.829	0.896	0.867
Gaussian	0.714	0.696	0.764	0.791	0.873	0.847	0.913	0.872
Gaussian Mixture	0.565	0.420	0.809	0.770	0.799	0.802	0.881	0.880

Table 3. Comparisons of the Kendall’s Tau of CLIP-Score for subsets acquired by statistical property similarity.

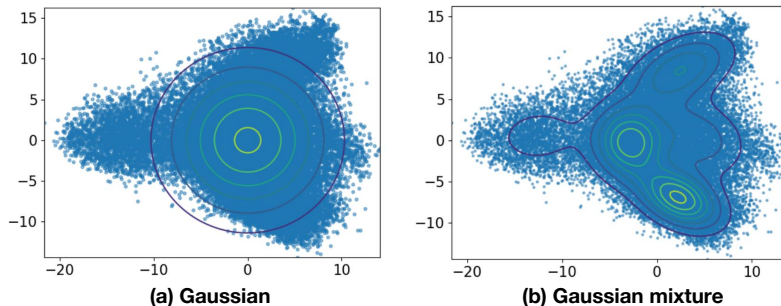


Figure 2. T-SNE Feature Visualization of Gaussian and Gaussian Mixture (#component = 5) distributions fitted to all features.

2. FlahEval’s Capability of Estimating Model Scores

In some scenarios, users may require the exact metric scores of models rather than the relative model rankings. FlashEval also supports fitting model scores using a small subset of N' by **searching for the mean square error**. As shown in Fig. 3 and Fig. 4, we demonstrate a comparison of our method against random sampling for model scores prediction. It is evident that our method can accurately estimate model scores even with smaller N' and it generalizes well to testing models across three sub-tasks.

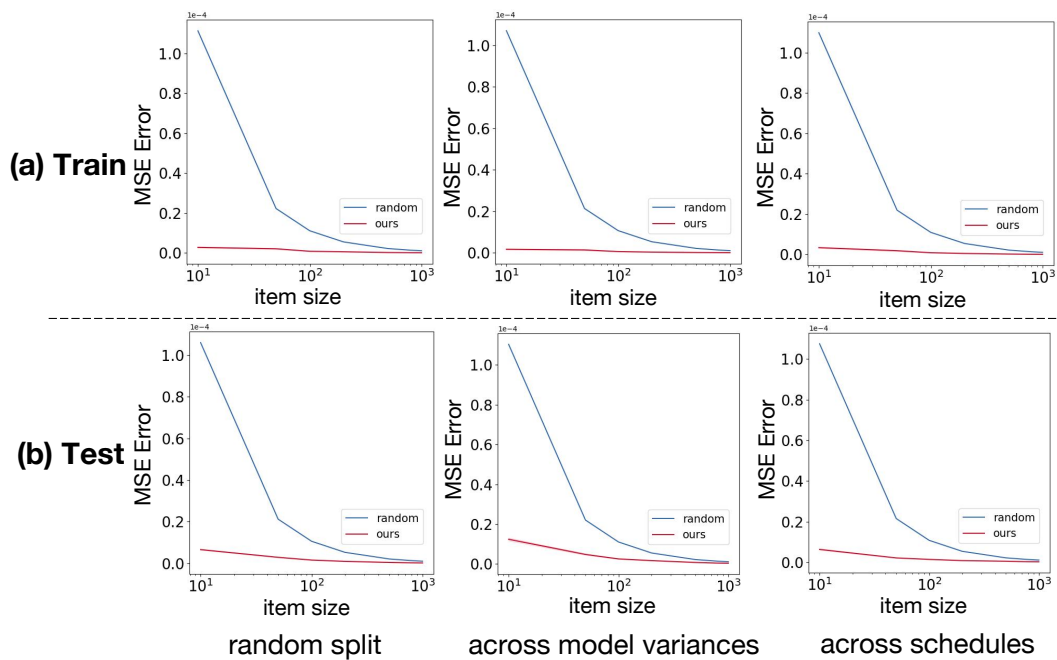


Figure 3. Comparisons with random sampling of the estimation quality for CLIP-Score on COCO dataset. Small MSE Error is better.

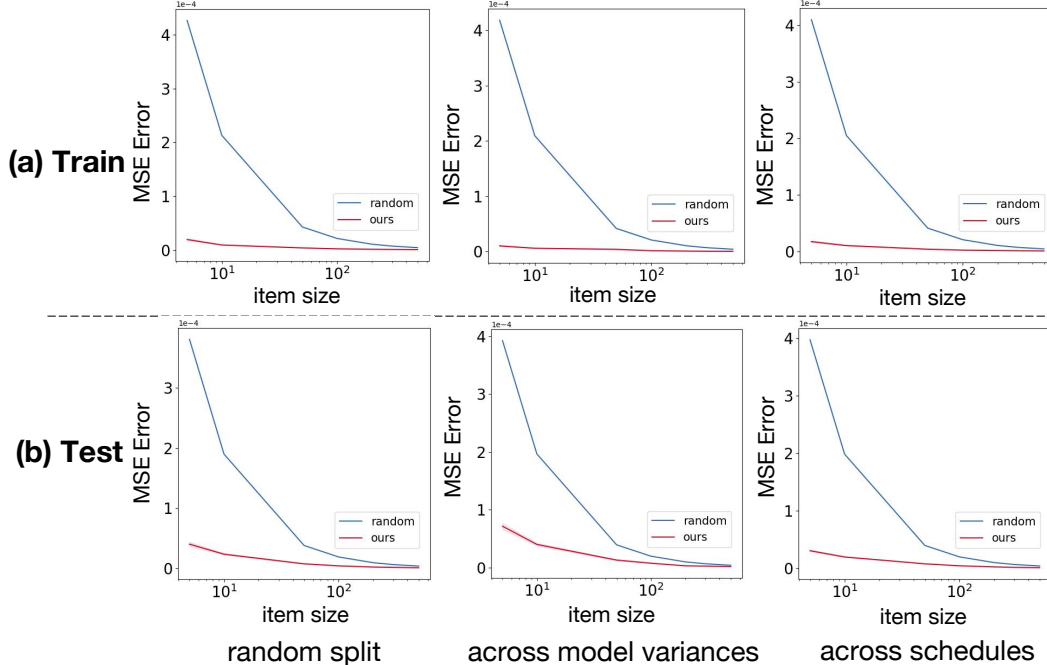


Figure 4. Comparisons with random sampling of the estimation quality for CLIP-Score on DiffusionDB dataset. Small MSE Error is better.

3. Analysis of Alternative Approaches to Acquire Representative Subsets

In the main paper, we recognize that “some subsets with high KD has large estimation error”, which reflects such sets lack precise estimation of the performance. Therefore, an intuitive question arises “could we use the estimation error as the metric for recognizing representative subsets?”. We could measure the “estimation error” by the MSE (Mean Squared Error) loss between the ground-truth performance (averaged on the entire set) and the subset’s performance $S'_j = \sum_{i \in \mathcal{I}} \frac{1}{N'} \mathcal{S}_{ij}$ for each model setting. It could be described as follows (N_m is the number of model settings):

$$\Omega_{EA} = \frac{1}{N_m} \sum_{i=1}^{N_m} (\hat{\mathcal{S}}_j - S'_j)^2. \quad (1)$$

In this section, we investigate this assumption. We discuss some intuitive alternative approaches to acquire representative subsets based on estimation error (“prompts with least estimation error”, “sets with least estimation error”) and demonstrate their failure. Furthermore, We analyze the underlying reason for their failure and why the FlashEval search method is necessary.

As shown in Fig. 5, we visualize the scores of all prompts in COCO dataset across two randomly selected models. It is noticeable that scores obtained by prompts for a certain model generally follows a Gaussian distribution. The mean of the gaussian distribution is the “ground-truth” model performance evaluation across the entire set. Therefore, if there exists prompts/subsets that could correctly approximate the mean of the gaussian distribution for all models, using the prompt/subsets for evaluation will have both high KD values and low estimation error.

Choose prompts with lowest estimation errors. To validate this existence of such prompts, we separately select 50 prompts with high KD values and 50 prompts with low estimation errors, examining their performance in the other criterion. From Tab. 4, we observe that prompts with high KD values have large estimation errors, while prompts with low estimation errors also exhibit small KD values. We conclude that **single prompt could not simultaneously satisfy high KD and estimation error, which means that “prompt that approximates the mean for all models” does not exist.** Therefore, simply choosing prompts with least estimation error to construct subsets is not feasible.

Choose subsets with lowest estimation errors. Similarly, we design experiments to validate the existence of the subsets that could precisely estimate the performance of each model. We generate 10,000,000 random samplings at various N' values, and select a set with the highest KD value and a set with the lowest estimation error. From Fig. 6, it is evident that the set we

discovered with the lowest error reaches only a KD value of roughly around 0.7 in training models when $N' = 10$. However, the subset with the highest KD value which b3-set discovered can achieve KD values above 0.8. Therefore, simply choosing subsets with least estimation error to construct subsets is also not feasible.

The underlying reason for their failure: As discussed above, in Fig. 6, we can see that when N' is small, the subset chosen based on estimation error performs poorly, indicating a lack of consistency between KD values and estimation errors. However, as N increases, the subsets acquired by “error-set” achieves relatively high KD. To further validate the above findings, we present the KD values and estimation errors with respect to N' in Fig. 7. We discover that the MSE error decreases sharply with $N' < 500$. Such phenomenon reveals the potential reason for “high KD and low estimation error could not be simultaneously satisfied for a smaller N' ”. When the N' is small, the estimation error is still relatively large (larger than 0.002), in which case, relatively lower estimation error does not guarantee higher KD. However, as N' increases, the estimation error is sufficiently small to ensure high correlation between lower estimation error and higher KD.

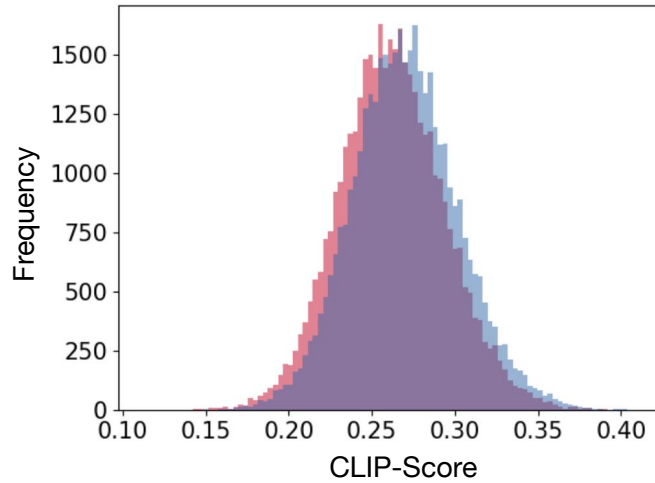


Figure 5. **Visualization of CLIP-Score of all prompts for two models.** The x-axis representing CLIP-Score values and the y-axis indicating the count of prompts achieving that value.

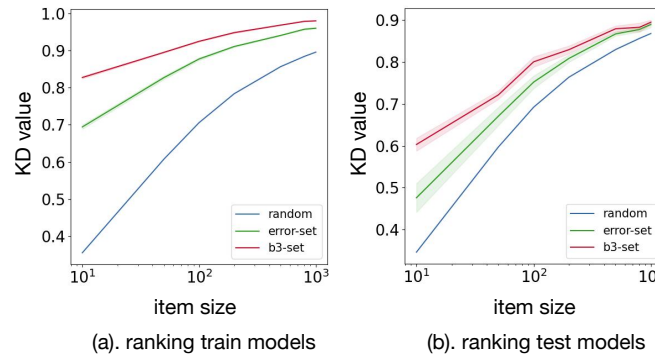


Figure 6. **Comparisons of set-based baselines using Kendall’s Tau for CLIP-Score on COCO dataset.** The shaded area are standard errors.

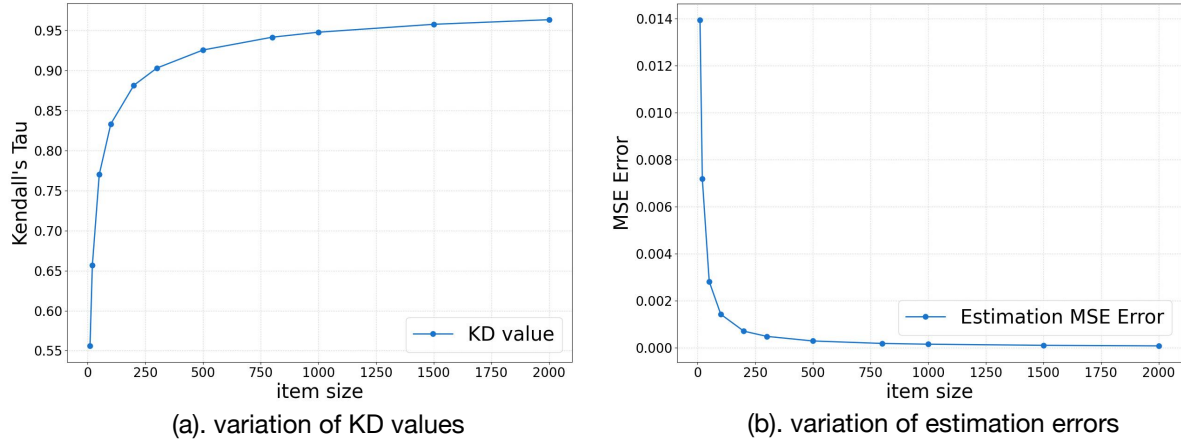


Figure 7. **The variation of average KD values and estimation errors with increasing N' .** The blue dots represent the average values obtained from 100,000 random samplings at the current N' .

Methods	Kendall's Tau \uparrow	Estimation Error \downarrow
All prompts mean	0.269/0.201	0.142/0.138
prompts with low errors	0.391/0.245	0.052/0.062
prompts with high KD values	0.678/0.463	0.117/0.117

Table 4. **The average performance of the top 50 prompts in one criterion on the other criterion for training/testing models.** The “All prompts mean” represents the averaged value for all independent prompts.

4. Qualitative Results of FlashEval Acquired Subsets

In this section, we present and analyze the FlashEval searched textual subset, to demonstrate the searched subset effectively distinguishes different model settings. We choose the $N' = 10$ item subset for the COCO annotations. Tab. 5 and Tab. 6 respectively outline the prompts included in the representative set for CLIP and ImageReward. Additionally, we show the images generated by prompts across different models in Fig. 8 and Fig. 9. The upper row displays the generated images under different schedules of the model “dreamlike”, sorted by the ground-truth metric values (the right has higher metric values). As could be seen, the image visual quality for this prompt (“A plate with rice topped with scallops and a side of broccoli.”) aligns with the ground-truth performance. It verifies that the searched representative subsets contain prompts correctly reflecting the model performance. Similarly, the lower row demonstrates that when using the same schedule, the prompt could also correctly rank different models. Aside from the CLIP-Score, we also present the generated images picked by human-preference-based metric ImageReward in Fig. 9, similar alignment could be witnessed.

item size	prompts in the representative set
$N' = 10$	"A plate with rice topped with scallops and a side of broccoli." "A decorative Asian lantern sculpture in the garden with flower ornaments." "A man is riding his skateboard down the road." "A close-up of a plate of food that has been eaten." "A small pizza in the middle of a table." "A white plate topped with a hot dog, french fries and condiments." "A black and white cat sitting on top of a pile of clothes." A giraffe and a zebra are on the grass near trees & cars. "A busy bus station with ramp going downstairs." "Heavily loaded green truck with passengers on back in roadway in urban area."

Table 5. Prompts in the representative set of CLIP-Score when $N' = 10$ on COCO dataset.

item size	prompts in the representative set
$N' = 10$	"Two giraffes standing around on the grassy plains." "there is someone holding a remote in there hand " "A man adjusts numbers at a tennis match." "A skateboarder catches major air during this stunt." "A little boy in a yellow shirt feeding a giraffe. " "An old building with two rusty, very old pick up trucks parked in front." "A little girl sitting on top of a bed next to a lamp." "An old rusty fire hydrant standing on a cracked sidewalk." "A man driving a car across an airport runway." "The german shepherd is guarding the front of the barred door,"

Table 6. Prompts in the representative set of ImageReward when $N' = 10$ on COCO dataset.

“A plate with rice topped with scallops and a side of broccoli.”

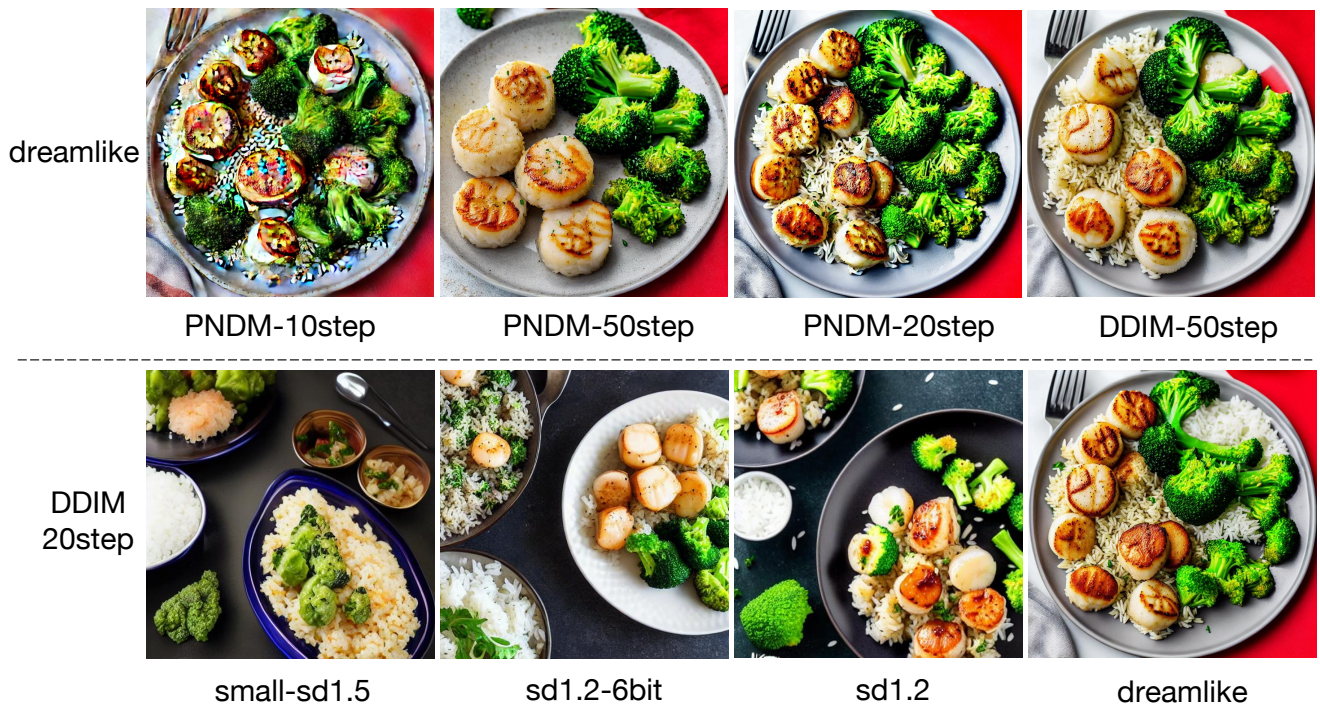


Figure 8. Visualization of images generated by the prompt from the representative set of CLIP-Score across various models. The first row represents variations among different schedules in the same model, while the second row depicts variances among different model architectures or parameters under the same schedule.

“An old rusty fire hydrant standing on a cracked sidewalk.”



Figure 9. Visualization of images generated by the prompt from the representative set of ImageReward across various models. The first row represents variations among different schedules in the same model, while the second row depicts variances among different model architectures or parameters under the same schedule.

5. Verification of the Generalization Ability of FlashEval

As mentioned in Table.2 and Table.3 in the main paper, we design 3 distinct model settings train-test split to verify the generalization ability of FlashEval. They act as 3 diverse subtask to test the generalization ability across random settings, across models, and across schedules. In this section, we firstly provide a detailed description of the 3 train-test split. Then, we present the results on these subtasks on COCO annotations (Fig. 10, Fig. 11, Fig. 12, Fig. 13, and Fig. 14) and DiffusionDB datasets (Fig. 15, Fig. 16, Fig. 17, Fig. 18, and Fig. 19).

5.1. Detailed List of Models

In Tabs. 7 to 9, we present detailed lists of models encompassed by both the training and testing models for each subtask. We conduct the search on the training models and validate the ranking correlation on the testing models to verify the generalization ability of FlashEval. The “random split” randomly splits the model setting zoo into two halves, the models of the training and testing splits have diverse model variants, solver, and timesteps. The “across model variants” uses only 4 models (SD-1.4, SD-1.5, SD-1.5-6bit,SD-2.1) as the training models, and the rest as the test models. The “across schedules” uses the “PNM-10”, “DDIM-10”, “DDIM-20”, “DDIM-50” as the train models.

5.2. Results under Different Model Settings Split on COCO annotations

We present the comparative results of all methods across various metrics in Figs. 10 to 14. It is evident that FlashEval consistently outperforms others in all scenarios, which demonstrates the generalization ability of FlashEval across diverse model or schedules.

5.3. Results under Different Model Settings Split on DiffusionDB

As mentioned in Sec. 5.1 of the main paper, due to the substantial size of the complete diffusionDB dataset, it is not feasible to iterate through all prompts (2,000,000) to get the ground-truth performance. Considering the cost, we select a relatively large amount of samples (5000) and use the averaged performance on them as the proxy “ground-truth” of DiffusionDB (treating the 5,000 prompts as the “whole set” to further condense it into smaller representative subset).

As could be witnessed from the Figs. 15 to 19, our acquired subsets achieves superior evaluation compared with the random sample baseline (adopted by recent literature Lee *et al.* [41]) quality for different metrics across different model splits.

training models			testing models		
model	solver	step	model	solver	step
stablediffusion1.2	PNDM	10	stablediffusion1.2	PNDM	20
stablediffusion1.2	DDIM	10	stablediffusion1.2	PNDM	50
stablediffusion1.2	DDIM	20	stablediffusion1.2-6bit	DDIM	10
stablediffusion1.2	DDIM	50	stablediffusion1.2-8bit	DDIM	10
stablediffusion1.2	DPM	10	stablediffusion1.2-8bit	DPM	10
stablediffusion1.2	DPM	20	stablediffusion1.2-8bit	DPM	20
stablediffusion1.2-6bit	DDIM	20	stablediffusion1.4	PNDM	20
stablediffusion1.2-6bit	DDIM	50	stablediffusion1.4	PNDM	50
stablediffusion1.2-6bit	DPM	10	stablediffusion1.4	DDIM	20
stablediffusion1.2-6bit	DPM	20	stablediffusion1.4	DPM	20
stablediffusion1.2-8bit	DDIM	20	stablediffusion1.4-6bit	DDIM	20
stablediffusion1.2-8bit	DDIM	50	stablediffusion1.4-6bit	DPM	20
stablediffusion1.4	PNDM	10	stablediffusion1.4-8bit	DDIM	10
stablediffusion1.4	DDIM	10	stablediffusion1.4-8bit	DDIM	20
stablediffusion1.4	DDIM	50	stablediffusion1.4-8bit	DDIM	50
stablediffusion1.4	DPM	10	stablediffusion1.4-8bit	DPM	20
stablediffusion1.4-6bit	DDIM	10	stablediffusion1.5	PNDM	20
stablediffusion1.4-6bit	DDIM	50	stablediffusion1.5	PNDM	50
stablediffusion1.4-6bit	DPM	10	stablediffusion1.5	DDIM	50
stablediffusion1.4-8bit	DPM	10	stablediffusion1.5	DPM	10
stablediffusion1.5	PNDM	10	stablediffusion1.5-6bit	DDIM	10
stablediffusion1.5	DDIM	10	stablediffusion1.5-6bit	DDIM	50
stablediffusion1.5	DDIM	20	stablediffusion1.5-6bit	DPM	20
stablediffusion1.5	DPM	20	stablediffusion1.5-8bit	DDIM	20
stablediffusion1.5-6bit	DDIM	20	stablediffusion1.5-8bit	DPM	10
stablediffusion1.5-6bit	DPM	10	stablediffusion1.5-8bit	DPM	20
stablediffusion1.5-8bit	DDIM	10	small-stablediffusion1.5	PNDM	20
stablediffusion1.5-8bit	DDIM	50	small-stablediffusion1.5	PNDM	50
small-stablediffusion1.5	PNDM	10	small-stablediffusion1.5	DPM	20
small-stablediffusion1.5	DDIM	10	stablediffusion2.1	PNDM	10
small-stablediffusion1.5	DDIM	20	stablediffusion2.1	PNDM	20
small-stablediffusion1.5	DDIM	50	stablediffusion2.1	DDIM	10
small-stablediffusion1.5	DPM	10	stablediffusion2.1	DDIM	50
stablediffusion2.1	PNDM	50	stablediffusion2.1	DPM	10
stablediffusion2.1	DDIM	20	stablediffusion2.1	DPM	20
dreamlike-photoreal	PNDM	20	dreamlike-photoreal	PNDM	10
dreamlike-photoreal	DDIM	10	dreamlike-photoreal	PNDM	50
dreamlike-photoreal	DDIM	20	dreamlike-photoreal	DPM	10
dreamlike-photoreal	DDIM	50	dreamlike-photoreal	DPM	20

Table 7. The detailed list of models in sub-task with "random split".

training models			testing models		
model	solver	step	model	solver	step
stablediffusion1.4	PNDM	10	stablediffusion1.2	PNDM	10
stablediffusion1.4	PNDM	20	stablediffusion1.2	PNDM	20
stablediffusion1.4	PNDM	50	stablediffusion1.2	PNDM	50
stablediffusion1.4	DDIM	10	stablediffusion1.2	DDIM	10
stablediffusion1.4	DDIM	20	stablediffusion1.2	DDIM	20
stablediffusion1.4	DDIM	50	stablediffusion1.2	DDIM	50
stablediffusion1.4	DPM	10	stablediffusion1.2	DPM	10
stablediffusion1.4	DPM	20	stablediffusion1.2	DPM	20
stablediffusion1.5	PNDM	10	stablediffusion1.2-6bit	DDIM	10
stablediffusion1.5	PNDM	20	stablediffusion1.2-6bit	DDIM	20
stablediffusion1.5	PNDM	50	stablediffusion1.2-6bit	DDIM	50
stablediffusion1.5	DDIM	10	stablediffusion1.2-6bit	DPM	10
stablediffusion1.5	DDIM	20	stablediffusion1.2-6bit	DPM	20
stablediffusion1.5	DDIM	50	stablediffusion1.2-8bit	DDIM	10
stablediffusion1.5	DPM	10	stablediffusion1.2-8bit	DDIM	20
stablediffusion1.5	DPM	20	stablediffusion1.2-8bit	DDIM	50
stablediffusion1.5-6bit	DDIM	10	stablediffusion1.2-8bit	DPM	10
stablediffusion1.5-6bit	DDIM	20	stablediffusion1.2-8bit	DPM	20
stablediffusion1.5-6bit	DDIM	50	stablediffusion1.4-6bit	DDIM	10
stablediffusion1.5-6bit	DPM	10	stablediffusion1.4-6bit	DDIM	20
stablediffusion1.5-6bit	DPM	20	stablediffusion1.4-6bit	DDIM	50
stablediffusion1.5-8bit	DDIM	10	stablediffusion1.4-6bit	DPM	10
stablediffusion1.5-8bit	DDIM	20	stablediffusion1.4-6bit	DPM	20
stablediffusion1.5-8bit	DDIM	50	stablediffusion1.4-8bit	DDIM	10
stablediffusion1.5-8bit	DPM	10	stablediffusion1.4-8bit	DDIM	20
stablediffusion1.5-8bit	DPM	20	stablediffusion1.4-8bit	DDIM	50
stablediffusion2.1	PNDM	10	stablediffusion1.4-8bit	DPM	10
stablediffusion2.1	PNDM	20	stablediffusion1.4-8bit	DPM	20
stablediffusion2.1	PNDM	50	small-stablediffusion1.5	PNDM	10
stablediffusion2.1	DDIM	10	small-stablediffusion1.5	PNDM	20
stablediffusion2.1	DDIM	20	small-stablediffusion1.5	PNDM	50
stablediffusion2.1	DDIM	50	small-stablediffusion1.5	DDIM	10
stablediffusion2.1	DPM	10	small-stablediffusion1.5	DDIM	20
stablediffusion2.1	DPM	20	small-stablediffusion1.5	DDIM	50
			small-stablediffusion1.5	DPM	10
			small-stablediffusion1.5	DPM	20
			dreamlike-photoreal	PNDM	10
			dreamlike-photoreal	PNDM	20
			dreamlike-photoreal	PNDM	50
			dreamlike-photoreal	DDIM	10
			dreamlike-photoreal	DDIM	20
			dreamlike-photoreal	DDIM	50
			dreamlike-photoreal	DPM	10
			dreamlike-photoreal	DPM	20

Table 8. The detailed list of models in sub-task with "across model variances".

training models			testing models		
model	solver	step	model	solver	step
stablediffusion1.2	PNDM	10	stablediffusion1.2	PNDM	20
stablediffusion1.2	DDIM	10	stablediffusion1.2	PNDM	50
stablediffusion1.2	DDIM	20	stablediffusion1.2	DPM	10
stablediffusion1.2	DDIM	50	stablediffusion1.2	DPM	20
stablediffusion1.2-6bit	DDIM	10	stablediffusion1.2-6bit	DDIM	50
stablediffusion1.2-6bit	DDIM	20	stablediffusion1.2-6bit	DPM	20
stablediffusion1.2-6bit	DPM	10	stablediffusion1.2-8bit	DDIM	50
stablediffusion1.2-8bit	DDIM	10	stablediffusion1.2-8bit	DPM	20
stablediffusion1.2-8bit	DDIM	20	stablediffusion1.4	PNDM	20
stablediffusion1.2-8bit	DPM	10	stablediffusion1.4	PNDM	50
stablediffusion1.4	PNDM	10	stablediffusion1.4	DPM	10
stablediffusion1.4	DDIM	10	stablediffusion1.4	DPM	20
stablediffusion1.4	DDIM	20	stablediffusion1.4-6bit	DDIM	50
stablediffusion1.4	DDIM	50	stablediffusion1.4-6bit	DPM	20
stablediffusion1.4-6bit	DDIM	10	stablediffusion1.4-8bit	DDIM	50
stablediffusion1.4-6bit	DDIM	20	stablediffusion1.4-8bit	DPM	20
stablediffusion1.4-6bit	DPM	10	stablediffusion1.5	PNDM	20
stablediffusion1.4-8bit	DDIM	10	stablediffusion1.5	PNDM	50
stablediffusion1.4-8bit	DDIM	20	stablediffusion1.5	DPM	10
stablediffusion1.4-8bit	DPM	10	stablediffusion1.5	DPM	20
stablediffusion1.5	PNDM	10	stablediffusion1.5-6bit	DDIM	50
stablediffusion1.5	DDIM	10	stablediffusion1.5-6bit	DPM	20
stablediffusion1.5	DDIM	20	stablediffusion1.5-8bit	DDIM	50
stablediffusion1.5	DDIM	50	stablediffusion1.5-8bit	DPM	20
stablediffusion1.5-6bit	DDIM	10	small-stablediffusion1.5	PNDM	20
stablediffusion1.5-6bit	DDIM	20	small-stablediffusion1.5	PNDM	50
stablediffusion1.5-6bit	DPM	10	small-stablediffusion1.5	DPM	10
stablediffusion1.5-8bit	DDIM	10	small-stablediffusion1.5	DPM	20
stablediffusion1.5-8bit	DDIM	20	stablediffusion2.1	PNDM	20
stablediffusion1.5-8bit	DPM	10	stablediffusion2.1	PNDM	50
small-stablediffusion1.5	PNDM	10	stablediffusion2.1	DPM	10
small-stablediffusion1.5	DDIM	10	stablediffusion2.1	DPM	20
small-stablediffusion1.5	DDIM	20	dreamlike-photoreal	PNDM	20
small-stablediffusion1.5	DDIM	50	dreamlike-photoreal	PNDM	50
stablediffusion2.1	PNDM	10	dreamlike-photoreal	DPM	10
stablediffusion2.1	DDIM	10	dreamlike-photoreal	DPM	20
stablediffusion2.1	DDIM	20			
stablediffusion2.1	DDIM	50			
dreamlike-photoreal	PNDM	10			
dreamlike-photoreal	DDIM	10			
dreamlike-photoreal	DDIM	20			
dreamlike-photoreal	DDIM	50			

Table 9. The detailed list of models in sub-task with "across schedules".

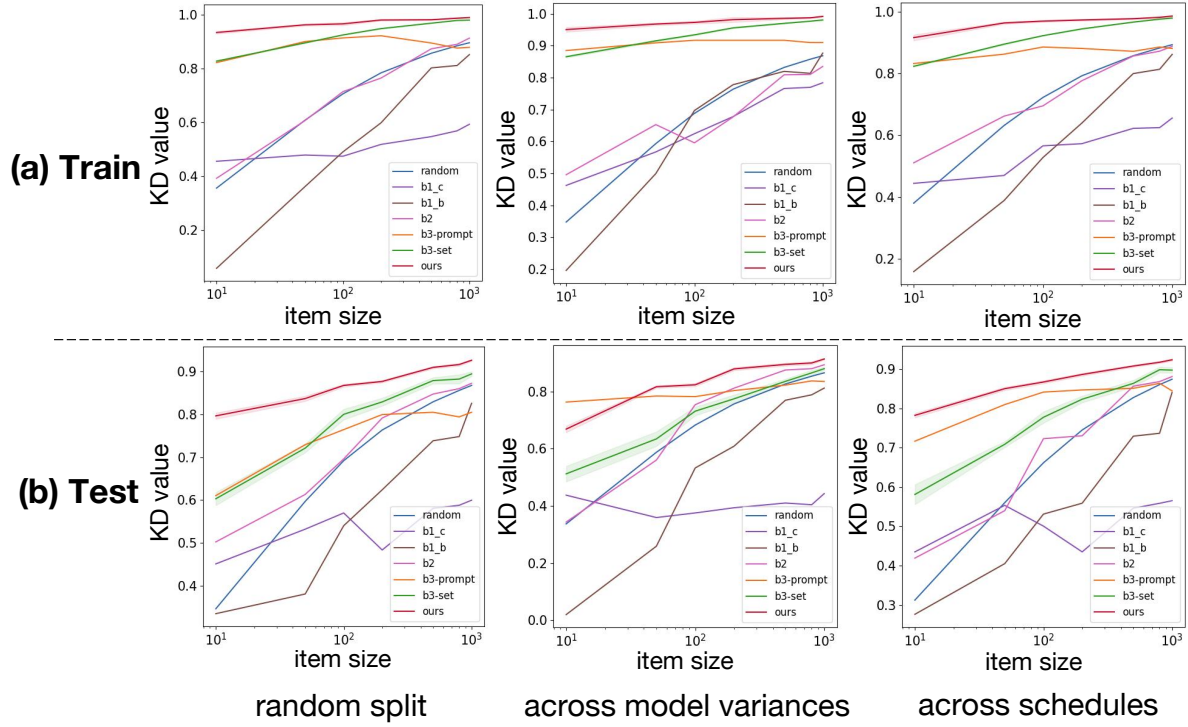


Figure 10. Comparisons with all baselines of the Kendall's Tau for CLIP-Score on COCO dataset.

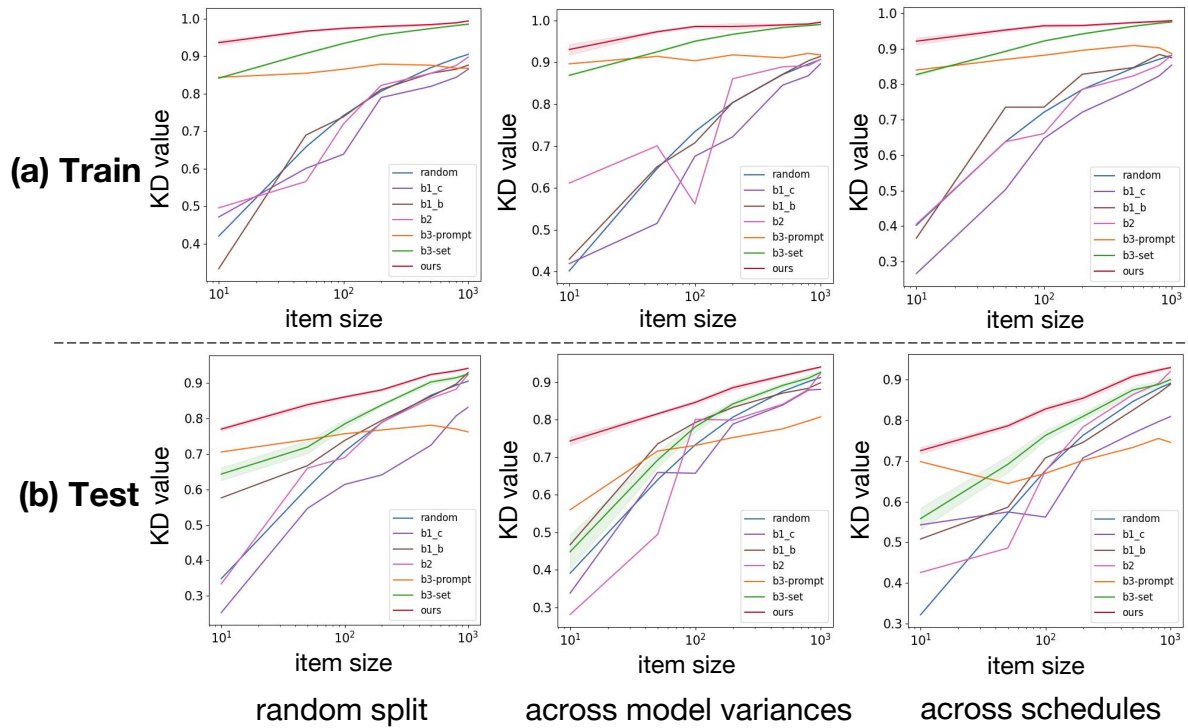


Figure 11. Comparisons with all baselines of the Kendall's Tau for ImageReward on COCO dataset.

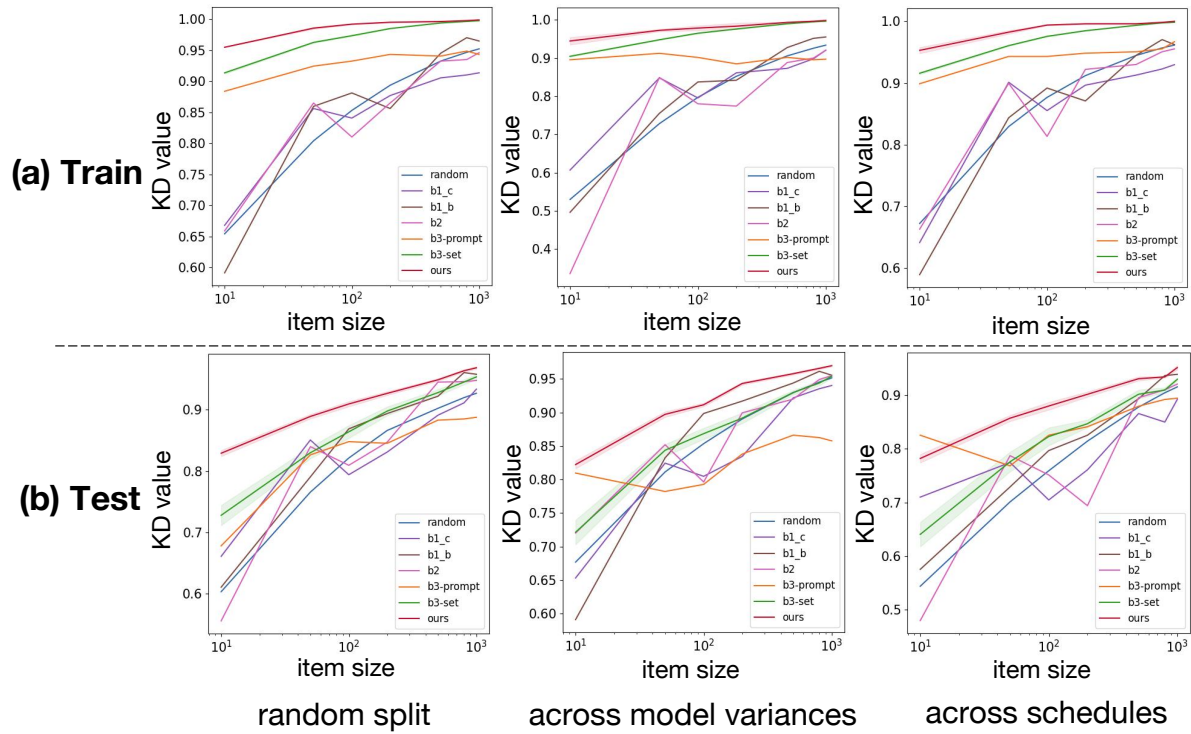


Figure 12. Comparisons with all baselines of the Kendall's Tau for HPS on COCO dataset.

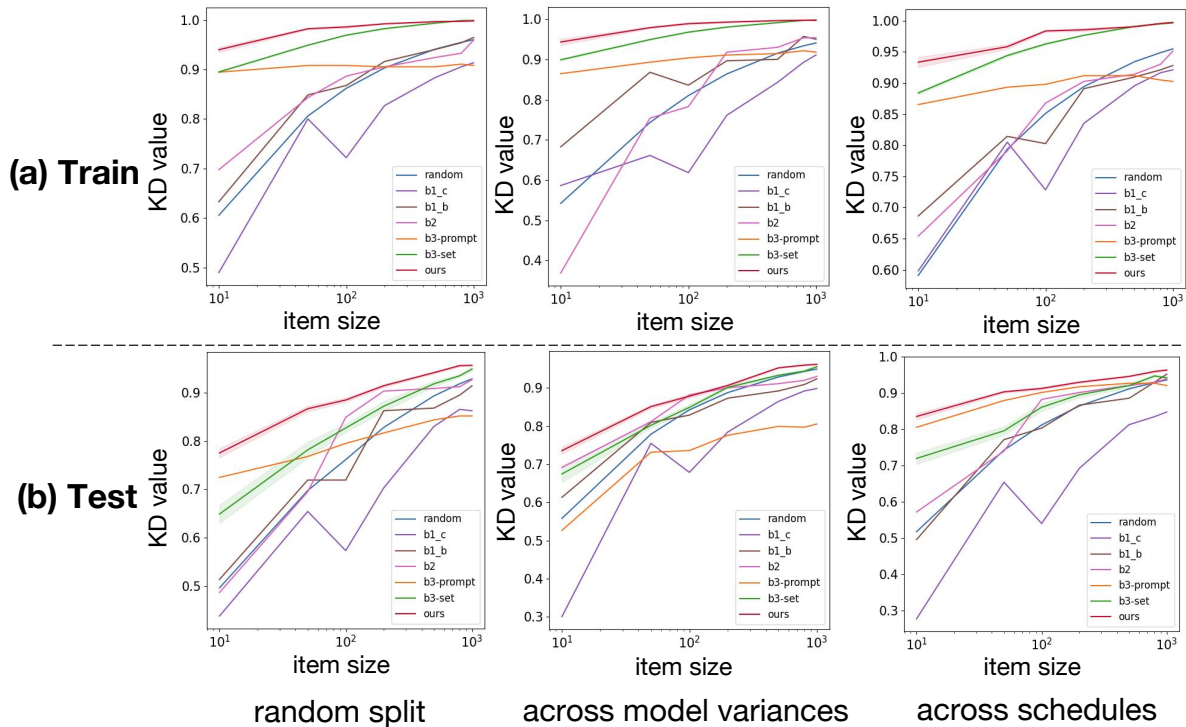


Figure 13. Comparisons with all baselines of the Kendall's Tau for Aesthetic on COCO dataset.

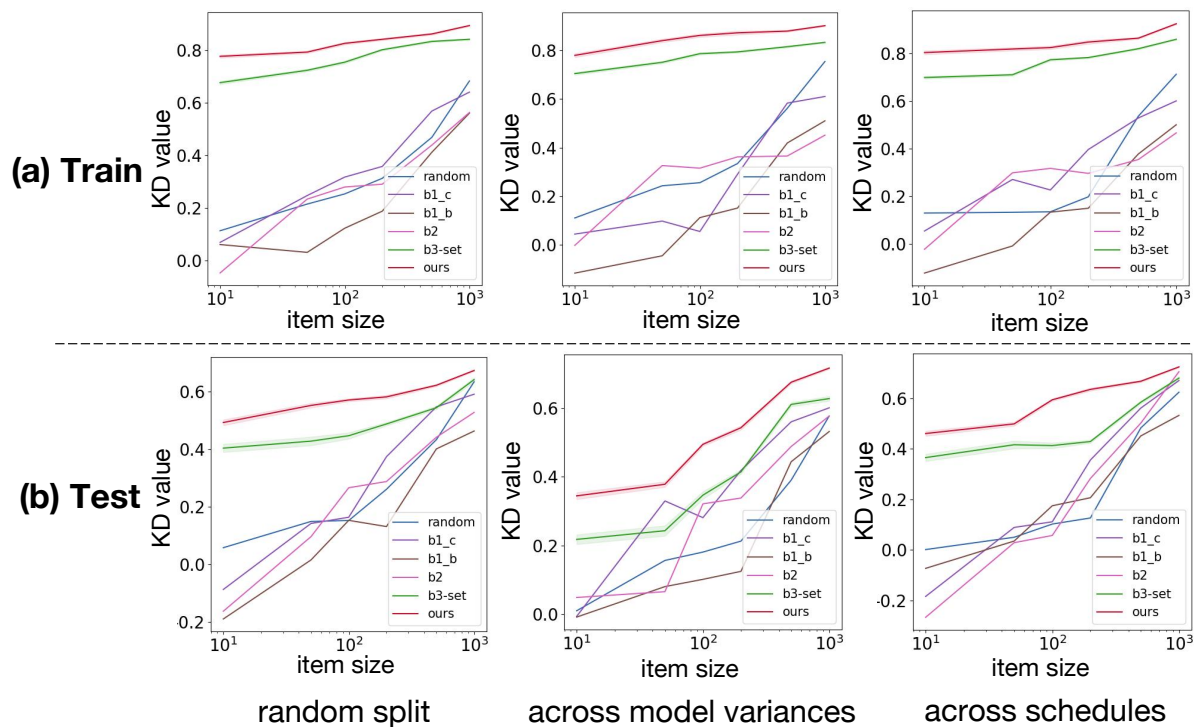


Figure 14. Comparisons with all baselines of the Kendall's Tau for FID on COCO dataset.

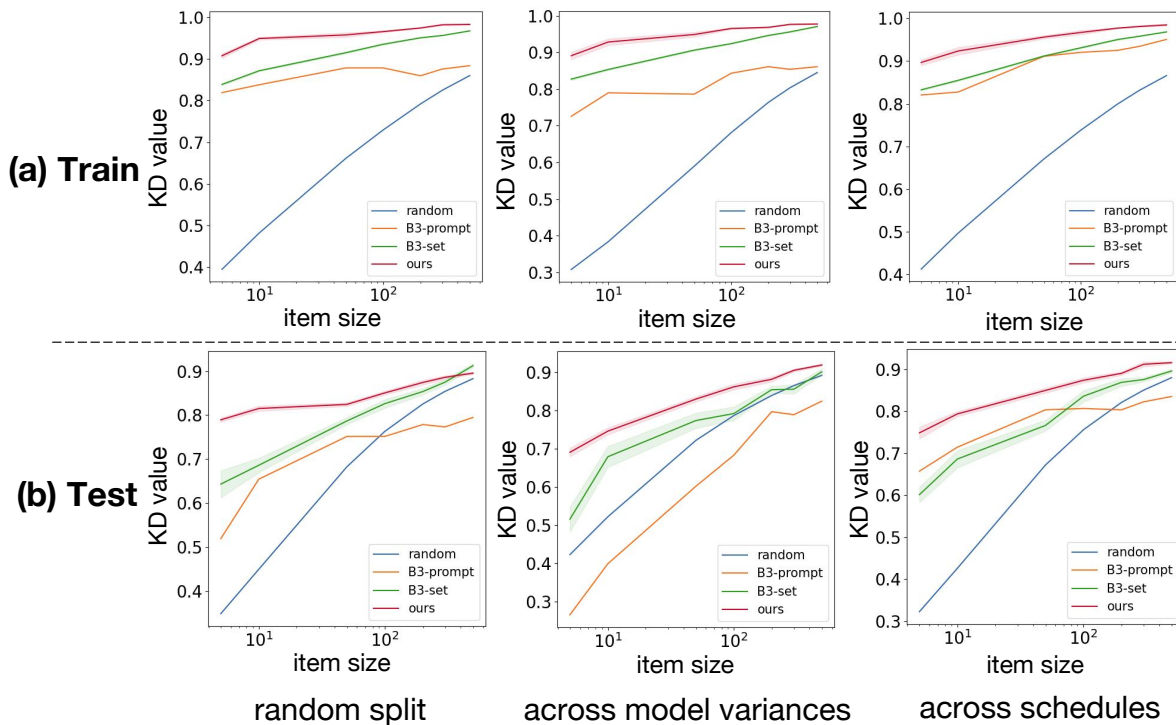


Figure 15. Comparisons with metric space baselines of the Kendall's Tau for CLIP-Score on DiffusionDB dataset.

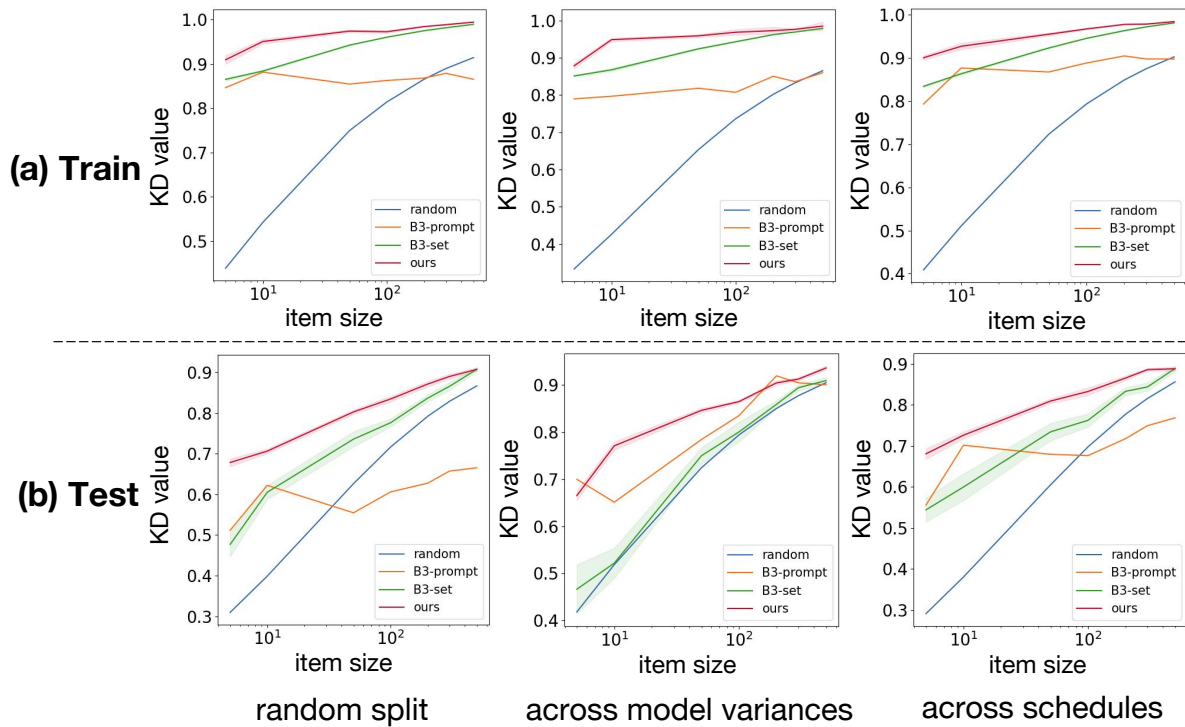


Figure 16. Comparisons with metric space baselines of the Kendall's Tau for ImageReward on DiffusionDB dataset.

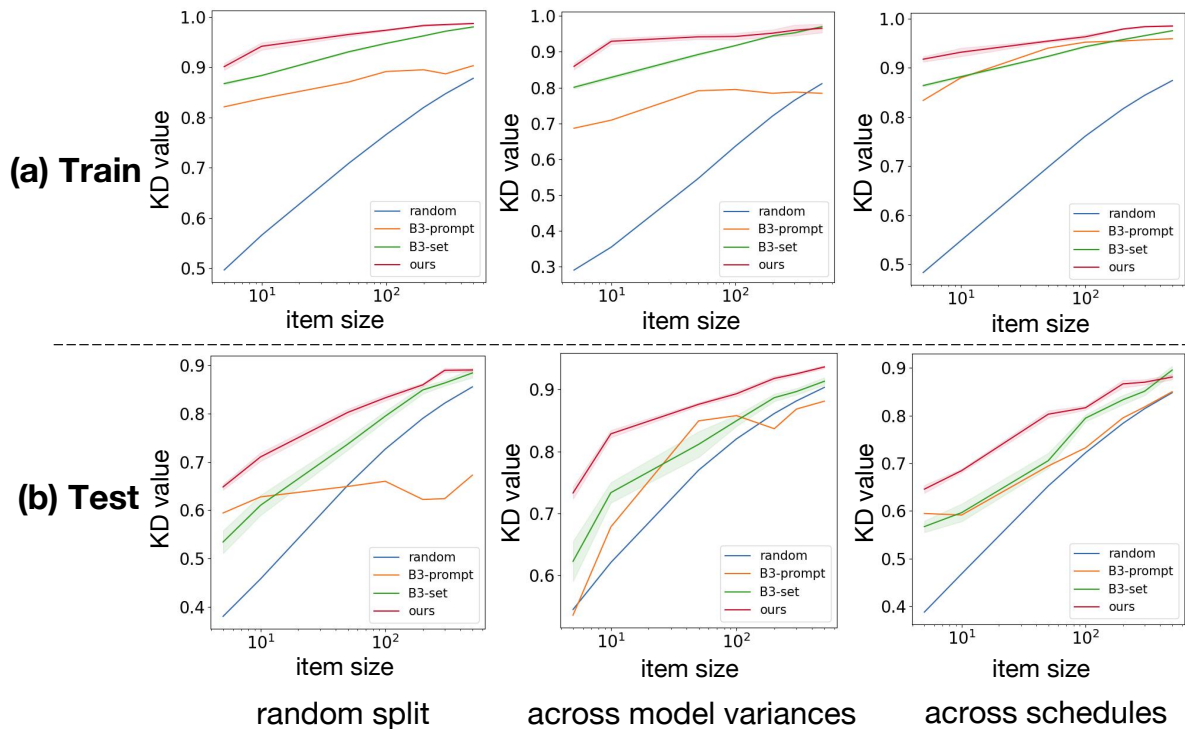


Figure 17. Comparisons with metric space baselines of the Kendall's Tau for HPS on DiffusionDB dataset.

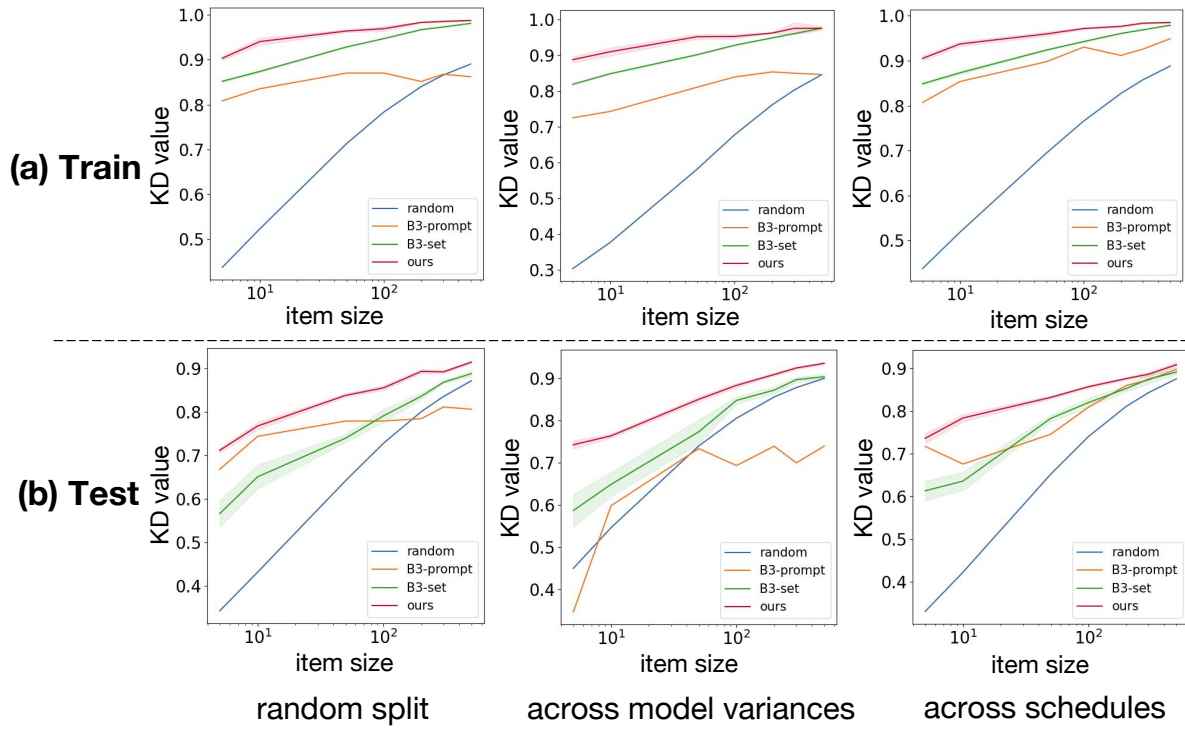


Figure 18. Comparisons with metric space baselines of the Kendall's Tau for Aesthetic on DiffusionDB dataset.

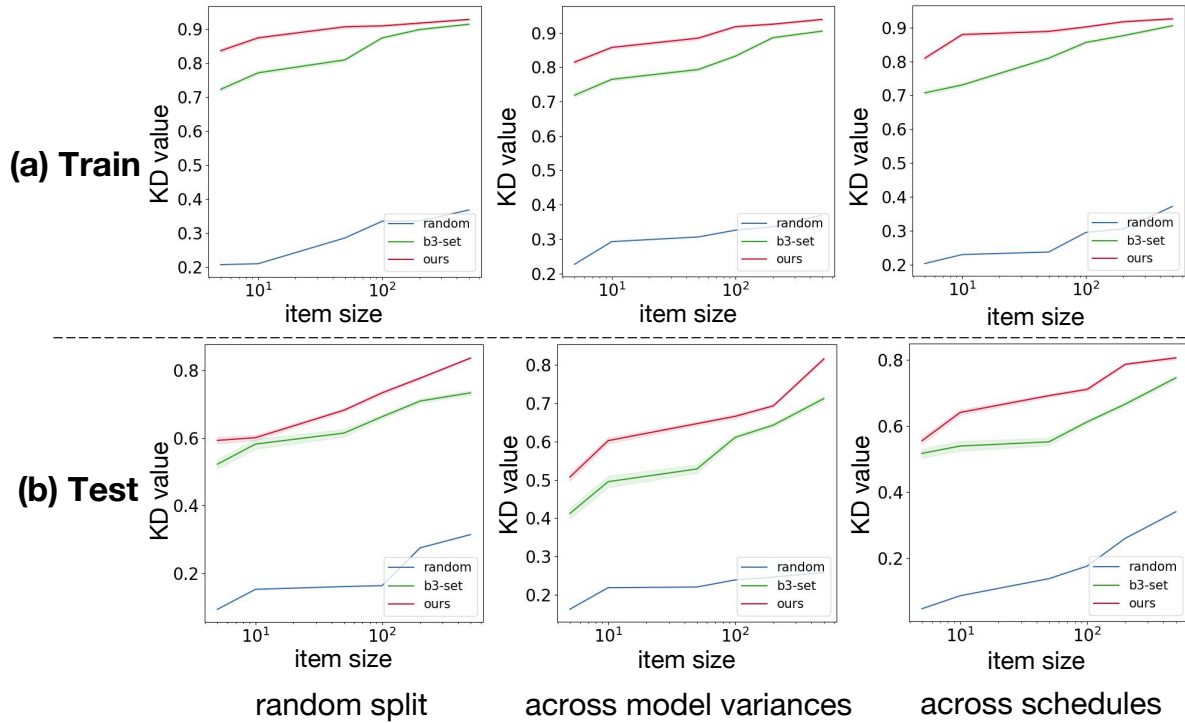


Figure 19. Comparisons with metric space baselines of the Kendall's Tau for FID on DiffusionDB dataset.