# Decentralized Diffusion Models

David McAllister[1*]    Matthew Tancik[2]    Jiaming Song[2]    Angjoo Kanazawa[1†]

[1]University of California, Berkeley        [2]Luma AI

## Abstract

*Large-scale AI model training divides work across thousands of GPUs, then synchronizes gradients across them at each step. This incurs a significant network burden that only centralized, monolithic clusters can support, driving up infrastructure costs and straining power systems. We propose Decentralized Diffusion Models, a scalable framework for distributing diffusion model training across independent clusters or datacenters by eliminating the dependence on a centralized, high-bandwidth networking fabric. Our method trains a set of expert diffusion models over partitions of the dataset, each in full isolation from one another. At inference time, the experts ensemble through a lightweight router. We show that the ensemble collectively optimizes the same objective as a single model trained over the whole dataset. This means we can divide the training burden among a number of "compute islands," lowering infrastructure costs and improving resilience to localized GPU failures. Decentralized diffusion models empower researchers to take advantage of smaller, more cost-effective and more readily available compute like on-demand GPU nodes rather than central integrated systems. We conduct extensive experiments on ImageNet and LAION Aesthetics, showing that decentralized diffusion models FLOP-for-FLOP outperform standard diffusion models. We finally scale our approach to 24 billion parameters, demonstrating that high-quality diffusion models can now be trained with just eight individual GPU nodes in less than a week.*

## 1. Introduction

Diffusion models achieve breakthrough results in image generation [38, 39], video modeling [4], and robotic control through diffusion policies [10]. However, these models continue to demand greater and greater training compute. Even early successes in image diffusion underscored the immense computational requirements, with Stable Diffusion 1.5's training consuming over 6,000 A100 GPU days [7, 33].
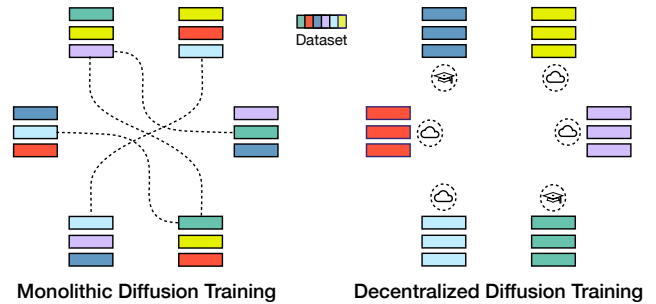
Figure 1. **Decentralized Diffusion Models (DDM).** Left: Existing diffusion models (monolithic) require synchronized, centralized training across thousands of GPUs, making high-quality training systems expensive and inaccessible. Right: DDM divides a diffusion model into an ensemble of expert models, each trained on its own data cluster in complete isolation. This ensemble collectively optimizes the same diffusion objective as a single model trained on all the data. This enables flexible training across diverse cloud or academic compute facilities. At inference, the ensemble delivers improved performance at the same FLOP-cost, making high-quality diffusion model training more efficient and accessible. See Figure 2 for large-scale DDM model samples.

Video diffusion models demand a step-function increase in these resources. For instance, Meta's Movie Gen trains on up to 6,114 H100 GPUs [34]—nearly GPT-3 scale [6]. As these models continue to grow in size and capability, they encounter significant systems-level challenges that mirror those faced in LLM training. High-bandwidth interconnect becomes a critical bottleneck, storage systems must handle massive streaming datasets and persistent hardware failures can crash or silently slow the training process [1, 14]. The result is a fragile integrated system where performance and reliability depend on complex interactions between arrays of accelerators and networking hardware. This is challenging and costly for large industry labs—untenable for academics.

We present Decentralized Diffusion Models, a scalable method for distributing the modeling burdening of diffusion across independent expert models, each trained on its own "compute island" with no cross-communication. This enables training with scattered resources, leveraging cost-effective cloud compute or combining resources across mul-

Figure 2. **Decentralized diffusion models train on readily-available hardware and generate high quality, diverse images.** We present selected samples from our 8x3B parameter model.

tiple clusters. Decentralization is also helpful for training "foundation"-scale models [3], where it is increasingly difficult to build datacenters with enough power capacity for the scale of modern training runs [30]. Independent training provides additional benefits, such as the ability to execute across heterogeneous hardware, making it possible to reuse existing training clusters alongside newer infrastructure.

Decentralized diffusion models employ a new training objective, Decentralized Flow Matching (DFM), which partitions the training data into K groups and trains a dedicated diffusion model over each. We refer to these specialized models as experts, inspired by the Mixture of Experts (MoE) literature [16]. An independently trained router model orchestrates these experts, determining which are most appropriate at test-time. We show that this ensemble of expert models collectively optimizes the same global objective as a single model trained on the entire dataset.

During each inference step, the router predicts each expert's relevance to the input noise and condition. The expert predictions are then linearly combined, weighed by their associated relevance scores. Experts learn to specialize, so many are irrelevant to a given input, and it is more efficient to use a subset of them. If not all experts are selected, it serves as a sparse model, leveraging selective computation similar to MoE. Sparse inference is compute-efficient but remains memory-intensive. We demonstrate that we can distill experts into a single dense model, achieving convenient inference while preserving sample quality.

We evaluate our approach on ImageNet and a filtered subset of LAION-Aesthetics, where we compare DDMs against existing monolithic diffusion model training. We systematically vary the number of experts to analyze its impact on downstream generation tasks, where we find that DDMs with eight experts consistently achieve the best performance, even outperforming a single model trained with the same compute. This is made possible by the additional parametrization provided by multiple specialized experts to outperform a generalist model. Finally, we demonstrate the scaling ability of DDMs by training a large decentralized model with eight 3-billion parameter experts, resulting in high definition image generation, as shown in Figure 2.

Decentralized Diffusion Models also integrate easily into existing diffusion training environments. In practice, DDMs involve clustering a dataset then training a standard diffusion model on each cluster. This means nearly everything from existing diffusion infrastructure can be reused—training code, dataloading, systems optimizations, noise schedules and architectures. See our blog post for a visual walkthrough and pseudocode.

## 2. Related Works

**Accelerating Diffusion Models** Flow matching [26, 27] generalizes diffusion models [20, 43–45] and makes

tractable large-scale training of continuous normalizing flows (CNFs) [8]. Diffusion and flow matching models produce state of the art generations in high dimensional continuous spaces but are typically expensive to train to high quality. Recent works propose token masking [42] and aligning internal states with pretrained image representation models [51] to accelerate the model learning process. PixArt-$\alpha$ [7] shows that curated datasets with synthetic captions reduce training costs. We explore a complementary direction that decentralizes the training of the model across multiple GPU clusters without cross-communication, thereby increasing the robustness and efficiency of training.

Existing approaches all use data parallel training, which distributes batches across GPUs and synchronizes gradients to produce larger effective batch sizes. Our work introduces an orthogonal form of parallelism that partitions the model's training across isolated compute centers with no gradient synchronization. Our approach is complementary with existing techniques—in fact, we employ Fully Sharded Data Parallel (FSDP) [52] training within each cluster while maintaining isolation between them in our experiments.

**Mixture of Experts** Mixture of Experts (MoE) is a popular and powerful approach to increase model capacity without a proportional increase in computational cost. This is achieved through sparse parameter activation: MoE replaces each dense feed forward network (FFN) layer of a transformer with a lightweight router that selects $k$ of $N$ learned FFN layers per token. The approach gained prominence with Switch Transformers [16], followed by works that refined MoE via various routing strategies [53] and systems improvements [17].

These advances have led to notable successes in language modeling, exemplified by models like Mixtral [23] and DeepSeek-V3 [11], which achieve strong performance while maintaining computational efficiency. While our work focuses on decentralized training, we draw inspiration from MoE techniques at test-time, adapting the principle of increasing parametrization through sparsely-activated experts to the diffusion domain. Both MoE and DDMs increase total model capacity without increasing computational costs, but we achieve this by routing inputs to different data experts rather than routing tokens to different FFNs.

MoE architectures introduce significant systems challenges, particularly the problems of expert load balancing, managing capacity factors, and the challenge of holding many times more parameters in memory during training and inference. Conversely, DDMs alleviate systems constraints by training experts and routers individually, extending diffusion training to readily available hardware configurations.

**Low Communication Learning** Federated learning [29] establishes foundational techniques for communication-

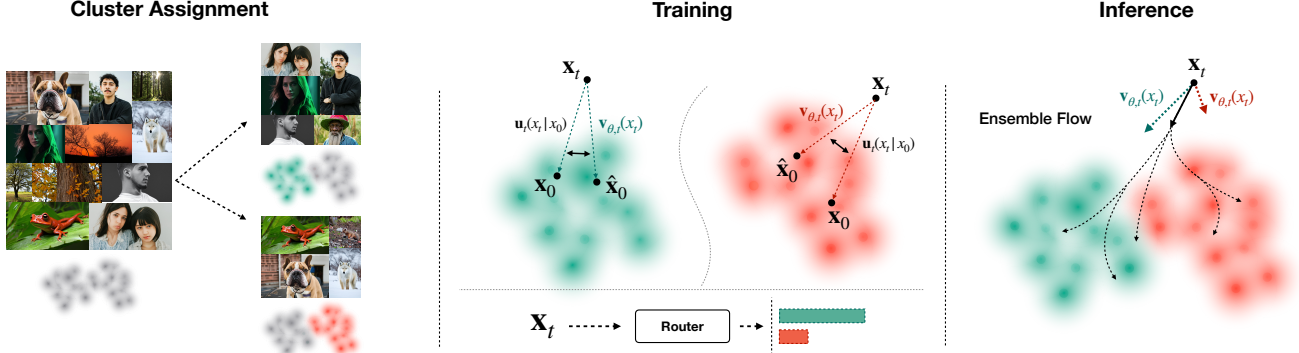**Cluster Assignment**  **Training**  **Inference**

Figure 3. **Decentralized Diffusion Model (DDM) Training Overview.** DDMs follow a three-step training process. We first cluster the dataset using off-the-shelf representation extraction models. We train a diffusion model over each of these clusters and a router that associates any input $x_t$ with its most likely clusters. At test-time, given a noisy sample, each expert (in red and green) predict their own flows, which combine linearly via the weights predicted by the router. The combined flow samples the entire distribution and is illustrated on the right.

constrained training, originally motivated by data privacy concerns. Follow-up works including Federated Averaging [47], Adaptive Federated Optimization [37] and DiLoCo [13], adapt these ideas to address large-scale training challenges. DiLoCo, uses an inner and outer optimization loop to balance local training with periodic global synchronization. This parallels developments in large model training, where methods like Branch-Train-Merge [25] and Diffusion Soup [2, 49] have explored training and combining data-specialist models. While our approach shares technical insights with these methods, we focus on leveraging distributed computation for robustness and flexibility rather than data sovereignty or training on mobile devices. In fact, our method should complement many of these approaches for further decentralization, which we hope to see in future work.

## 3. Decentralized Flow Matching

We introduce Decentralized Flow Matching (DFM), the training objective for DDMs, to distribute the diffusion modeling burden across an ensemble of expert denoisers, each trained individually without cross-communication. We show that the flow matching objective divides naturally into this distributed arrangement. Since diffusion models and rectified flows can be seen as special cases of flow matching, DFM applies directly across these popular algorithms. Specifically, we train each expert over a pre-determined and sharded subset of the data distribution. Each expert's predicted flow is then combined at test-time using a router that trains with a simple classification objective. This yields an ensemble that samples the entire data distribution.

### 3.1. Preliminary: Flow Matching Objective

Flow matching [26] defines a fixed forward corruption process and a learned reverse denoising process. We index this

process with timesteps $t$ interpolating the data distribution at $t = 0$ and the per-pixel Gaussian distribution at $t = 1$. Concretely, $x_t = \alpha_t x_0 + \sigma_t \epsilon$, where the scaling coefficients $\alpha_t$ and $\sigma_t$ are chosen ahead of time, typically to preserve variance across timesteps [45]. The learned model, $v_{\theta,t}(x_t)$, reverses this corruption process to transport from the noise distribution to the data distribution. These transport paths collectively represent the marginal flow, $u_t(x_t)$, which is a vector field at each timestep that can be regressed in a model to interpolate the data distribution or computed analytically over a dataset, which will reproduce training samples.

Flow transports a latent variable, $x_t$, to the data distribution through a multi-step sampling process. At each step, the model predicts the flow from $x_t$ toward a weighted average of the data distribution. In its continuous definition, we express this as an integral:

$$u_t(x_t) = \int_{x_0} u_t(x_t|x_0) \frac{p_t(x_t|x_0)q(x_0)}{p_t(x_t)} dx_0 \qquad (1)$$

Where $u_t(x_t)$ is the marginal flow, $x_t$ is our noisy latent, $x_0$ is a data sample, $p_t(x_t|x_0)$ is Gaussian and $u_t(x_t|x_0)$ is the conditional flow from $x_0$ to $x_t$. Flow matching regresses this analytical form of $u_t$ through a parametric model. Since we train over a discrete dataset, the integral becomes a summation,

$$u_t(x_t) = \frac{1}{p_t(x_t)} \sum_{x_0} u_t(x_t|x_0)p_t(x_t|x_0)q(x_0). \qquad (2)$$

### 3.2. Decentralized Flow Matching Objective

DFM decomposes the marginal flow into a series of expert flows. We partition the data into $K$ disjoint clusters $\{S_1, S_2, \ldots, S_K\}$, and each expert trains on an assigned subset ($x_0 \in S_i$). This is a natural choice, since empirical

4

results suggest that image data lies on a disjoint union of manifolds [5, 24] and this may be one reason diffusion models are effective in the image domain [48]. Splitting across these subsets yields marginal flow of the form:

$$u_t(x_t) = \sum_{k=1}^{K} \frac{1}{p_t(x_t)} \sum_{x_0 \in S_k} u_t(x_t|x_0) p_t(x_t|x_0) q(x_0). \quad (3)$$

Finally, we find that marginal flow can be written as a linear combination of expert flows,

$$u_t(x_t) = \sum_{k=1}^{K} \underbrace{\frac{p_{t,S_k}(x_t)}{p_t(x_t)}}_{\text{Router}} \underbrace{\sum_{x_0 \in S_k} \frac{u_t(x_t|x_0) p_t(x_t|x_0) q(x_0)}{p_{t,S_k}(x_t)}}_{\text{Data-Expert Flow}}, \quad (4)$$

where the outer sum results in a categorical probability distribution over all experts that we refer to as a 'router', and the inner sum is the marginal flow over each subset that we refer to as a data-expert flow. Please see the supplement for an alternative derivation based on score matching.

This division is highly convenient for large-scale training. Instead of one monolith training run, we can train a number of independent expert denoisers. We train each with the standard flow-matching objective and zero cross-model communication. Separately, we train a router that predicts the probability that $x_t$ is drawn from each data subset, which can be formulated as the classification task we describe below. Similar to MoE, this learned router delegates computation to different parameters during inference. A key difference is that our router can be explicitly trained in isolation, instead of being trained end-to-end with gradients flowing through the entire model. At inference time, the experts combine in an ensemble that, as we show above, collectively optimize the same global objective as a monolithic training.

### 3.3. Router Training

The router predicts the probability of a given $x_t$ being drawn from each of the data partitions.

$$p(k|x_t, t) = \frac{p_{t,S_k}(x_t)}{p_t(x_t)}, \quad k \in [K] \quad (5)$$

To learn this, we sample input $x_t$ then supervise with a cross-entropy loss over the cluster labels associated to each data sample. This loss is minimized when the model, $r_\theta(x_t, t)$, exactly predicts $p(k|x_t, t)$. See Algorithm 1 for details. Note that router trains independently of the denoisers.

We parameterize $r_\theta(x_t, t)$ with a small diffusion transformer (DiT) [32] and use the same conditioning mechanisms as a DiT denoiser. We append a learned classification token [12] that we decode to cluster logits through a linear head following usual conventions.

---

**Algorithm 1** Flow Router Training

**Require:** Training data $\{x_0, k\}$, schedule $\{\alpha_t, \sigma_t\}_{t=1}^{T}$
1: **while** not converged **do**
2:     $x_0, k \sim \mathcal{D}$
3:     $t \sim \{1, \dots, T\}$
4:     $\epsilon \sim \mathcal{N}(0, \mathbf{I})$
5:     $x_t = \alpha_t x_0 + \sigma_t \epsilon$
6:     $\mathbf{z} = r_\theta(x_t, t) \in \mathbb{R}^{|K|}$
7:     Update $\theta$ using $\nabla_\theta \mathcal{L}_{\text{CE}}(\mathbf{z}, \text{OneHot}(k))$
8: **end while**

---

### 3.4. Expert Training

The DFM objective divides the training task into learning a series of expert denoisers. Conveniently, these experts employ the same objective as standard flow matching:

$$\mathcal{L}_{\text{CFM}}(\theta) = \mathbb{E}_{t,q(x_0),p_t(x_t|x_0)} \|\mathbf{v}_{\theta,t}(x_t) - \mathbf{u}_t(x_t|x_0)\|^2. \quad (6)$$

We parameterize each expert as its own DiT, though it is important to note that our method is architecture-agnostic. We can even assign each denoiser a different architecture as long as it optimizes the above objective.

### 3.5. Inference Strategy

At test-time, we merge the predictions of our experts following the weights predicted by the router:

$$u_t(x_t) = \sum_{k=1}^{K} \underbrace{r_\theta(x_t, t)}_{\text{Router}} \underbrace{v_{\theta,t}(x_t)}_{\text{Expert}}. \quad (7)$$

While the full ensemble is necessary to exactly match the global flow matching objective, it is much less expensive to use a subset of experts or even only one expert for each test-time prediction. The model's FLOP cost scales linearly with the number of active experts in each forward pass. This introduces a trade-off between theoretical correctness and computational efficiency. We can interpret this selection as instantiating model sparsity, where the model activates only a subset of its parameters in each forward pass. In MoE, this is demonstrated to improve performance with the same computational budget [16]. We compare different inference strategies empirically to outline this efficiency-correctness tradeoff, where we find selecting only the top expert is the most efficient approach and does not sacrifice quality.

### 3.6. Distillation

In production systems, it is difficult to deploy models with extremely high parameter counts. A DDM model with single-expert selection at test-time achieves nearly the same FLOPs per forward pass as a monolithic model, even though it has

$|K|$ times more parameters. Still, it requires $|K|$ times more memory than the monolithic model to load in the first place. As model sizes scale past the memory budget of today's accelerators, this becomes a challenging distributed systems problem. While this is achievable in many cases, we present a distillation approach as a convenient alternative.

Specifically, the sparse DDM model can be distilled into a dense model. After the distillation, inference can be done just like any other diffusion model trained in a non-decentralized manner. Note that the goal of distillation here is to reduce $K$ experts into a single model, not to reduce the number of denoising steps, as done in other diffusion distillation works [40, 46, 50]. We follow a teacher-student training procedure [18]. Specifically, we replace the conditional flow $u_t(x_t|x_0)$ supervision target with a prediction from the teacher model:

$$\mathcal{L}_{\text{distill}}(\theta) = \mathbb{E}_{t,q(x_0),p_t(x_t|x_0)} \left\| \mathbf{v}_{\theta,t}(x_t) - \mathbf{v}_{\text{teacher},t}(x_t) \right\|^2 .$$

In our case, we use the cluster label assigned to each data point to select a single teacher expert per training example.

## 4. Experiments

We evaluate the effectiveness of decentralized diffusion models on the academic ImageNet dataset and a subset of the LAION [41] dataset filtered for aesthetic scores of 5 or higher, which more closely resembles real-world training scenarios. Our goal is to analyze the design space of decentralized diffusion models applied to widely-adopted settings so that practitioners can confidently adapt the method to their needs. For this reason, we use standard architectures, hyperparameters and data when possible.

### 4.1. Implementation Details

**Dataset Partitioning** We partition the dataset in an efficient manner that facilitates learning the decentralized flow matching objective. Naive k-means scales quadratically with input data size, which is a non-starter for massive Internet datasets. Ma et al. [28] propose a multi-stage algorithm that efficiently consolidates a large number of fine-grained clusters into a small number of partitions. We adopt this approach to focus our efforts on generative modeling.

Following [28], we compute image features for the images in the dataset by pooling the output from DINOv2 [31] to incorporate both pixel features and semantics, cluster these features to 1024 fine-grained centroids, and then further consolidate to $k$ coarse centroids. We assign each data point to the nearest of the coarse centroids to produce the final set of partitions. This partitioning process is computationally negligible compared to training.

**Evaluations** We control for known sensitivities in FID calculations that arise from different implementations and evaluation splits. We designate a fixed evaluation set of 50,000 samples from each training dataset for computing features

| Inference Strategy | GFLOPs ↓ | FID ↓ | CLIP FID ↓ |
|---|---|---|---|
| Monolith | **308** | 12.81 | 5.58 |
| Oracle | **308** | 10.46 | 5.83 |
| Full | 2490 | 10.52 | 5.83 |
| Top-1 | 334 | **9.84** | **5.48** |
| Top-2 | 642 | 10.31 | 5.74 |
| Top-3 | 950 | 10.37 | 5.77 |
| Sample-1 | 334 | 157.05 | 51.17 |
| Sample-2 | 642 | 10.27 | 5.73 |
| Sample-3 | 950 | 10.44 | 5.78 |
| Threshold-0.01 | - | 10.46 | 5.81 |
| Threshold-0.05 | - | 10.37 | 5.75 |
| Threshold-0.1 | - | 10.15 | 5.72 |
| Nucleus ($T = 0.5$) | 334 | 188.66 | 60.09 |
| Nucleus ($T = 1.0$) | 334 | 152.16 | 48.37 |
| Nucleus ($T = 2.0$) | 334 | 33.9 | 14 |

Table 1. **Test-Time Combination Strategies.** We ablate strategies to sample from the ensemble at test-time and find that simply selecting the top expert outperforms more sophisticated alternatives.

and evaluating all models. Our FID figures align well with published results and our consistent implementation provides precise relative comparisons between different approaches. This standardization eliminates confounding variables that often complicate comparisons of generative models.

**Training Details** For ImageNet experiments, we adopt the hyperparameters and architecture from diffusion transformers [32], using a batch size of 256, EMA decay rate of 0.9999, and learning rate of 1e-4 without warmup or decay. We aim to replicate plausible real-world training runs with our LAION experiments, so we scale the batch size to 1024.

We reimplement the DiT XL/2 architecture for our denoising models, with each containing 895M parameters. In decentralized diffusion models, the total parameter count scales linearly with the number of experts. However, the computational cost remains constant during inference when using single expert selection.

For LAION experiments, we implement text conditioning following the Pix-Art Alpha [7] architecture, using SDXL's CLIP [9, 35] model to incorporate text via cross-attention. The router uses the smaller DiT B/2 architecture (158M parameters) augmented with a learned CLS token that decodes linearly to a probability distribution over the clusters.

We ensure fair comparison by maintaining consistent total computation across decentralized diffusion models and baseline monolith models. We achieve this simply by dividing total batch size evenly among experts. For example, with eight experts, a 256 monolith batch size corresponds to eight expert batches of size 32. This equalizes the total training FLOPs between DDMs and baselines. The router introduces an additional 4% measured training FLOPs overhead, which we also incorporate in our comparisons.
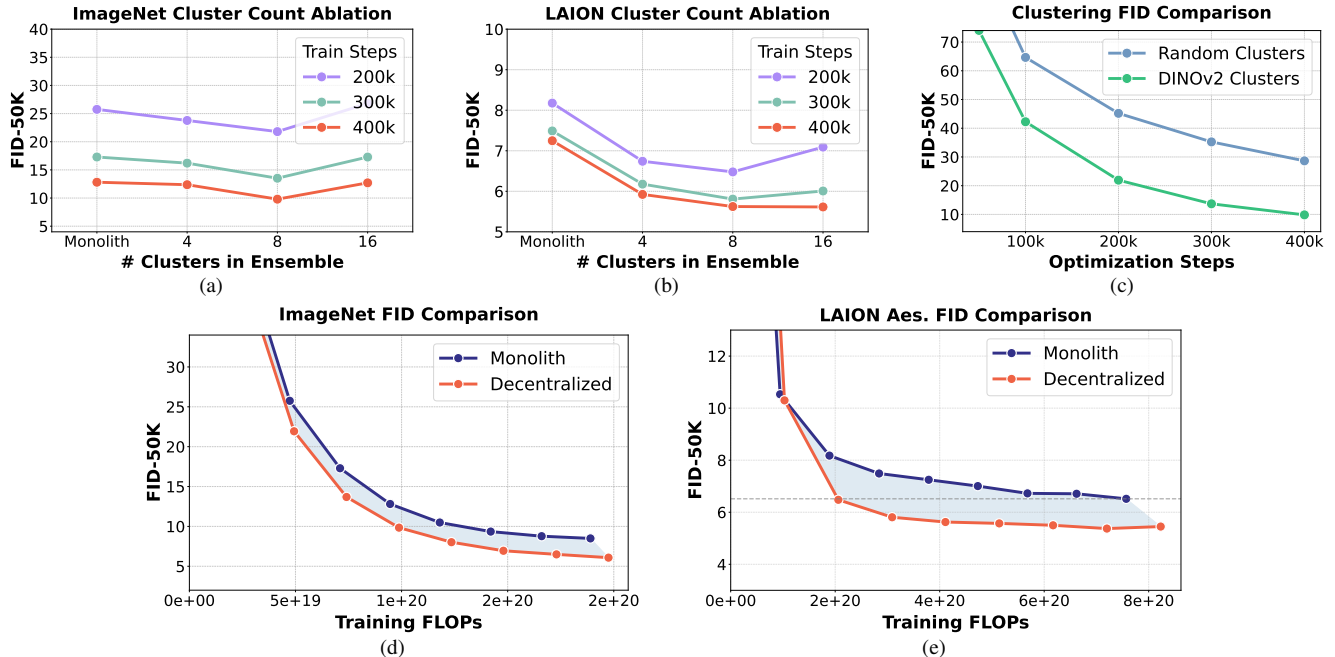
Figure 4. **Ablations at the DiT XL model scale.** Eight-expert DDMs display the best consistent performance on ImageNet (a) and LAION Aesthetics (b). We show the importance of image-based clustering on ImageNet compared to random clustering (c). Finally, FLOP-for-FLOP, decentralized diffusion models outperform monolith diffusion models on both datasets (d, e).

## 4.2. Ensembling at Test-Time

We first compare different strategies to combine expert predictions at test-time. A full estimate of the marginal flow involves linearly combining all expert predictions:

$$u_t(x_t) = \sum_{k=1}^{K} \underbrace{r_\theta(x_t, t)}_{\text{Router}} \underbrace{v_{\theta,t}(x_t)}_{\text{Expert}}. \quad (8)$$

In practice, selecting only important experts saves on computation. We evaluate the following strategies:

- **Full.** Compute the weighted combination of all expert predictions. This strategy's FLOP cost scales linearly with the number of experts.
- **Sample.** Sample from the router's predicted softmax distribution to select a single expert. This is an unbiased Monte-Carlo estimate of the marginal flow.
- **Top-k.** Simply use the $k$ experts with the maximum predicted router probability. Top-1 selection is the most efficient option at test-time.
- **Nucleus.** Sample one expert according to the nucleus (top-$p$) sampling strategy [21] commonly used in large language models. We use $p = 0.9$ and ablate softmax temperature in Table 1.
- **Oracle.** Select one expert according to the cluster label associated with an evaluation image. Used only to evaluate the effectiveness of the learned router.



Figure 5. **DDMs optimize the global diffusion objective.** We average samples from the monolithic and DDM ImageNet models using a deterministic sampler with matching random seeds (left) and compare them to outputs generated with random noise samples (right). The left samples are highly correlated, appearing less blurry.

In Table 1, we evaluate these inference strategies on ImageNet. We find that top-1 selection outperforms all other alternatives, while also incurring the lowest FLOP cost. For all other comparisons in Figures 4 and 6, we use top-1 selection because it nearly matches the computational cost of a dense model.

## 4.3. Selecting the Right Number of Experts

The DFM objective theoretically supports any number of experts, but practically we find this is an important hyperparameter for DDMs. This choice determines both the degree of decentralization and the total parameter count of the system. In the theoretical limit, where the number of experts approaches the number of training samples, the system would reduce to a nearest neighbor lookup that can only reproduce
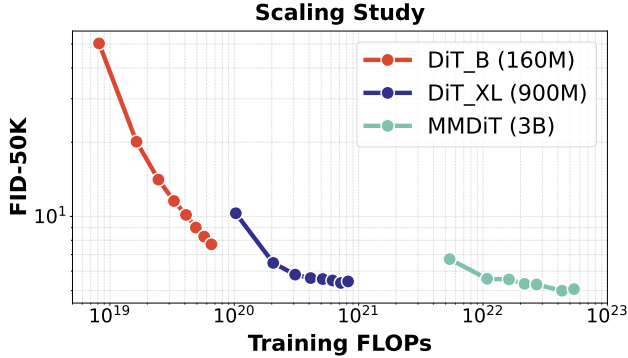
7

Figure 6. **Decentralized diffusion models scale gracefully to billions of parameters.** Throughout training, we plot the FID over LAION Aesthetics as a function of training compute. We find that increasing expert model capacity and training compute predictably improves performance.
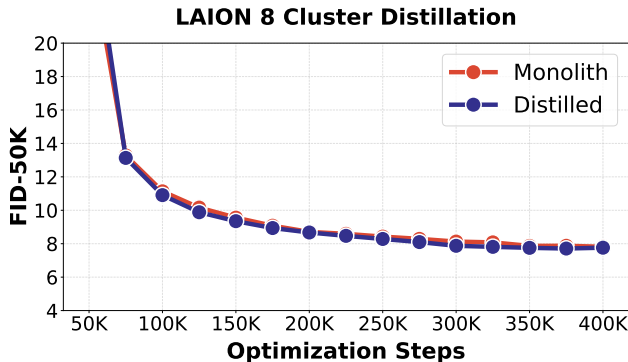


Figure 7. **Distilling a DDM into a dense model.** Dense models are often more convenient than sparse models to serve in production settings. Distilling a decentralized diffusion model into a dense model matches the performance of training a monolith from scratch at one third of the FLOP-cost (1/4 batch size).

training examples—equivalent to the analytical form of flow matching. We also find that individual experts train poorly as the global batch size is divided too aggressively, as in the 16 expert (batch size 16 per expert) experiment in Figure 4a.

We compare DDMs with 4, 8 and 16 experts in Figures 4a and 4b, and we find that eight experts achieves the best performance consistently on ImageNet and LAION. This configuration enables strong decentralization while maintaining reasonable test-time memory requirements. The eight experts appear to specialize meaningfully while preserving coherent coverage of the data distribution. Empirically, we find this is a sweet spot for the competing factors of model capacity, decentralization, and practical deployment.

### 4.4. DDMs vs. Monoliths

We compare decentralized diffusion models against a fair monolithic baseline on both datasets. We find that decentralized diffusion models with eight experts consistently outper-

form standard diffusion models on a FLOP-for-FLOP basis. In Figure 4d, we compare FID on ImageNet up to 800k training steps. The decentralized diffusion model achieves a 28% lower FID at 800k steps (6.081 vs. 8.494). Note that we plot FID as a function of training FLOPs to account for the 4% additional training cost of the router.

We find that the performance uplift is also significant on captioned internet data with LAION in Figure 4e. In fact, it achieves a lower FID of 6.48 at 200k optimization steps than the monolith's 6.52 FID at 800k steps. This represents a 4x training speedup as well as a lower convergence FID.

We also visually verify the correctness of the DFM objective. The standard diffusion objective predictably pairs noise samples and data samples, so a well-trained DDM should sample a similar image as a monolith for the same input noise. We verify this in Figure 5. Please see more samples and analysis of DDM in the supplemental.

### 4.5. Data Clustering Ablation

DFM imposes no explicit constraints on cluster size or composition, so we ablate two clustering strategies. We find that the chosen strategy significantly impacts model performance. Comparing feature-based clustering using DINO against random cluster assignments, which would maintain i.i.d. properties across partitions, reveals that feature-based clustering improves results significantly. We hypothesize that there is more mutual information within feature clusters than random clusters, meaning expert models can more efficiently compress and specialize in their assigned subdistributions. When data possesses semantic or low-level feature similarities, experts are free to learn more focused representations.

### 4.6. Distillation

While our method achieves computational efficiency through top-1 expert selection at inference time, the total memory footprint of many expert models can be substantial. We address this limitation through knowledge distillation, compressing the ensemble's capabilities into a dense model. Our approach supervises a student model with predictions from the expert ensemble, selecting the appropriate expert for each training example based on its cluster label. This can be seen as distilling the top-1 sparse model and incurs the same cost as standard teacher-student distillation.

Our distilled model matches the performance of directly training on the dataset, despite using only one quarter of the batch size and, consequently, one third the training FLOPs (assuming a backward pass costs double a forward pass). After 400k training steps, the distilled model achieves an FID of 7.76, comparable to the baseline model's FID of 7.82 (Figure 7). Many diffusion distillation works focus on reducing the number of sampling steps, whereas we just aim to replicate the ensemble in a dense model. We leave the

exploration of combining our method with sampling-focused distillation techniques as promising future work.

## 4.7. Scaling Experiment

We perform a scaling study of decentralized diffusion models. At each scale, we follow best practices gleaned from our ablations and train a system of eight expert models, each based on the FLUX MMDiT architecture [15]. We use a hidden dimension of 2560, depth of 30, and separate text and visual token streams, which total 3B parameters per expert. We encode text prompts using a single T5 XL model [36] and mix their features with image features through self attention.

Crucially, each expert can be trained independently on readily available hardware. With 16 GPUs per expert, we train at 0.28 seconds per iteration for 1M pretraining steps. Using gradient accumulation, this is equivalent to training each expert on a single on-demand cloud GPU node for six and a half days. This demonstrates that our method enables training large-scale diffusion models without specialized infrastructure or large integrated compute clusters.

We evaluate our large-scale ensemble against smaller DiT B and XL [32] ensembles in Figure 6. DDM performance improves as a function of expert parametrization and does not saturate at any scale we tried. We finally finetune our largest ensemble for 60k steps on high-resolution data and display some selected samples in Figure 2.

## 5. Discussion

Decentralized diffusion models enable high-quality generative model training across isolated compute clusters, greatly broadening the possible hardware configurations for diffusion model training. While we focus on distributing computational resources, DFM theoretically permits decentralizing data as well—a property with potential privacy implications for domains like medical imaging. Experts can train locally where sensitive data resides, and a router can train on samples from these experts rather than the raw data. These ideas allow DDM to preserve data privacy and sovereignty, as private data never leaves its original location. Furthermore, combining DDMs with low-bandwidth training methods could push the boundaries of decentralization—perhaps enabling large-scale model training on true commodity hardware. While we experiment on image modeling, the principles proposed by DDM may be applied to other domains, such as medical imaging, robotic policies and video modeling. We look forward to future works in these directions.

## References

[1] Wei An, Xiao Bi, Guanting Chen, Shanhuang Chen, Chengqi Deng, Honghui Ding, Kai Dong, Qiushi Du, Wenjun Gao, Kang Guan, et al. Fire-flyer ai-hpc: A cost-effective software-hardware co-design for deep learning. *arXiv preprint arXiv:2408.14158*, 2024. 1

[2] Benjamin Biggs, Arjun Seshadri, Yang Zou, Achin Jain, Aditya Golatkar, Yusheng Xie, Alessandro Achille, Ashwin Swaminathan, and Stefano Soatto. Diffusion soup: Model merging for text-to-image diffusion models. *arXiv preprint arXiv:2406.08431*, 2024. 4

[3] Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021. 3

[4] Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe Taylor, Troy Luhman, Eric Luhman, et al. Video generation models as world simulators. 2024. *URL https://openai. com/research/video-generation-models-as-world-simulators*, 3, 2024. 1

[5] Bradley CA Brown, Anthony L Caterini, Brendan Leigh Ross, Jesse C Cresswell, and Gabriel Loaiza-Ganem. Verifying the union of manifolds hypothesis for image data. *arXiv preprint arXiv:2207.02862*, 2022. 5

[6] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, and Sandhini Agarwal et al. Language models are few-shot learners, 2020. 1

[7] Junsong Chen, Jincheng Yu, Chongjian Ge, Lewei Yao, Enze Xie, Yue Wu, Zhongdao Wang, James Kwok, Ping Luo, Huchuan Lu, and Zhenguo Li. Pixart-$\alpha$: Fast training of diffusion transformer for photorealistic text-to-image synthesis, 2023. 1, 3, 6

[8] Ricky TQ Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. Neural ordinary differential equations. *Advances in neural information processing systems*, 31, 2018. 3

[9] Mehdi Cherti, Romain Beaumont, Ross Wightman, Mitchell Wortsman, Gabriel Ilharco, Cade Gordon, Christoph Schuhmann, Ludwig Schmidt, and Jenia Jitsev. Reproducible scaling laws for contrastive language-image learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2818–2829, 2023. 6

[10] Cheng Chi, Zhenjia Xu, Siyuan Feng, Eric Cousineau, Yilun Du, Benjamin Burchfiel, Russ Tedrake, and Shuran Song. Diffusion policy: Visuomotor policy learning via action diffusion. *The International Journal of Robotics Research*, page 02783649241273668, 2023. 1

[11] DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Haowei Zhang, Honghui Ding, Huajian Xin, Huazuo Gao,

Hui Li, Hui Qu, J. L. Cai, Jian Liang, Jianzhong Guo, Jiaqi Ni, Jiashi Li, Jiawei Wang, Jin Chen, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, Junxiao Song, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Lei Xu, Leyi Xia, Liang Zhao, Litong Wang, Liyue Zhang, Meng Li, Miaojun Wang, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Mingming Li, Ning Tian, Panpan Huang, Peiyi Wang, Peng Zhang, Qiancheng Wang, Qihao Zhu, Qinyu Chen, Qiushi Du, R. J. Chen, R. L. Jin, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, Runxin Xu, Ruoyu Zhang, Ruyi Chen, S. S. Li, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shaoqing Wu, Shengfeng Ye, Shengfeng Ye, Shirong Ma, Shiyu Wang, Shuang Zhou, Shuiping Yu, Shunfeng Zhou, Shuting Pan, T. Wang, Tao Yun, Tian Pei, Tianyu Sun, W. L. Xiao, Wangding Zeng, Wanjia Zhao, Wei An, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, X. Q. Li, Xiangyue Jin, Xianzu Wang, Xiao Bi, Xiaodong Liu, Xiaohan Wang, Xiaojin Shen, Xiaokang Chen, Xiaokang Zhang, Xiaosha Chen, Xiaotao Nie, Xiaowen Sun, Xiaoxiang Wang, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xingkai Yu, Xinnan Song, Xinxia Shan, Xinyi Zhou, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, Y. K. Li, Y. Q. Wang, Y. X. Wei, Y. X. Zhu, Yang Zhang, Yanhong Xu, Yanhong Xu, Yanping Huang, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Li, Yaohui Wang, Yi Yu, Yi Zheng, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Ying Tang, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yu Wu, Yuan Ou, Yuchen Zhu, Yuduan Wang, Yue Gong, Yuheng Zou, Yujia He, Yukun Zha, Yunfan Xiong, Yunxian Ma, Yuting Yan, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Z. F. Wu, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhen Huang, Zhen Zhang, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhibin Gou, Zhicheng Ma, Zhigang Yan, Zhihong Shao, Zhipeng Xu, Zhiyu Wu, Zhongyu Zhang, Zhuoshu Li, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Ziyi Gao, and Zizheng Pan. Deepseek-v3 technical report, 2024. 3

[12] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018. 5

[13] Arthur Douillard, Qixuan Feng, Andrei A Rusu, Rachita Chhaparia, Yani Donchev, Adhiguna Kuncoro, Marc'Aurelio Ranzato, Arthur Szlam, and Jiajun Shen. Diloco: Distributed low-communication training of language models. *arXiv preprint arXiv:2311.08105*, 2023. 4

[14] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024. 1

[15] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first International Conference on Machine Learning*, 2024. 9, 12

[16] William Fedus, Barret Zoph, and Noam Shazeer. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *Journal of Machine Learning Research*, 23(120):1–39, 2022. 3, 5

[17] Trevor Gale, Deepak Narayanan, Cliff Young, and Matei Zaharia. Megablocks: Efficient sparse training with mixture-of-experts. *Proceedings of Machine Learning and Systems*, 5: 288–304, 2023. 3

[18] Geoffrey Hinton. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015. 6

[19] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance, 2022. 12

[20] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 3

[21] Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text degeneration, 2020. 7

[22] Emiel Hoogeboom, Jonathan Heek, and Tim Salimans. Simple diffusion: End-to-end diffusion for high resolution images, 2023. 12

[23] Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024. 3

[24] Hamidreza Kamkari, Brendan Leigh Ross, Rasa Hosseinzadeh, Jesse C Cresswell, and Gabriel Loaiza-Ganem. A geometric view of data complexity: Efficient local intrinsic dimension estimation with diffusion models. *arXiv preprint arXiv:2406.03537*, 2024. 5

[25] Margaret Li, Suchin Gururangan, Tim Dettmers, Mike Lewis, Tim Althoff, Noah A Smith, and Luke Zettlemoyer. Branch-train-merge: Embarrassingly parallel training of expert language models. *arXiv preprint arXiv:2208.03306*, 2022. 4

[26] Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022. 3, 4, 12

[27] Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. *arXiv preprint arXiv:2209.03003*, 2022. 3

[28] Jiawei Ma, Po-Yao Huang, Saining Xie, Shang-Wen Li, Luke Zettlemoyer, Shih-Fu Chang, Wen-Tau Yih, and Hu Xu. Mode: Clip data experts via clustering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26354–26363, 2024. 6

[29] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*, pages 1273–1282. PMLR, 2017. 3

[30] Mark Morey. Data center owners turn to nuclear as potential electricity source - u.s. energy information administration (eia), 2024. 3

[31] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023. 6

[32] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4195–4205, 2023. 5, 6, 9

[33] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952*, 2023. 1

[34] Adam Polyak, Amit Zohar, Andrew Brown, Andros Tjandra, Animesh Sinha, Ann Lee, Apoorv Vyas, Bowen Shi, Chih-Yao Ma, Ching-Yao Chuang, et al. Movie gen: A cast of media foundation models. *arXiv preprint arXiv:2410.13720*, 2024. 1

[35] 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. 6

[36] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67, 2020. 9

[37] Sashank Reddi, Zachary Charles, Manzil Zaheer, Zachary Garrett, Keith Rush, Jakub Konečnỳ, Sanjiv Kumar, and H Brendan McMahan. Adaptive federated optimization. *arXiv preprint arXiv:2003.00295*, 2020. 4

[38] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 1

[39] 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. *Advances in neural information processing systems*, 35:36479–36494, 2022. 1

[40] Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. *arXiv preprint arXiv:2202.00512*, 2022. 6

[41] 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. *Advances in Neural Information Processing Systems*, 35:25278–25294, 2022. 6

[42] Vikash Sehwag, Xianghao Kong, Jingtao Li, Michael Spranger, and Lingjuan Lyu. Stretching each dollar: Diffusion training from scratch on a micro-budget. *arXiv preprint arXiv:2407.15811*, 2024. 3

[43] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International conference on machine learning*, pages 2256–2265. PMLR, 2015. 3

[44] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *Advances in neural information processing systems*, 32, 2019.

[45] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020. 3, 4, 12

[46] Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency models. *arXiv preprint arXiv:2303.01469*, 2023. 6

[47] Tao Sun, Dongsheng Li, and Bao Wang. Decentralized federated averaging. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(4):4289–4301, 2022. 4

[48] Peng Wang, Huijie Zhang, Zekai Zhang, Siyi Chen, Yi Ma, and Qing Qu. Diffusion models learn low-dimensional distributions via subspace clustering. *arXiv preprint arXiv:2409.02426*, 2024. 5

[49] Mitchell Wortsman, Gabriel Ilharco, Samir Ya Gadre, Rebecca Roelofs, Raphael Gontijo-Lopes, Ari S Morcos, Hongseok Namkoong, Ali Farhadi, Yair Carmon, Simon Kornblith, et al. Model soups: averaging weights of multiple fine-tuned models improves accuracy without increasing inference time. In *International conference on machine learning*, pages 23965–23998. PMLR, 2022. 4

[50] Tianwei Yin, Michaël Gharbi, Richard Zhang, Eli Shechtman, Fredo Durand, William T Freeman, and Taesung Park. One-step diffusion with distribution matching distillation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6613–6623, 2024. 6

[51] Sihyun Yu, Sangkyung Kwak, Huiwon Jang, Jongheon Jeong, Jonathan Huang, Jinwoo Shin, and Saining Xie. Representation alignment for generation: Training diffusion transformers is easier than you think. *arXiv preprint arXiv:2410.06940*, 2024. 3

[52] Yanli Zhao, Andrew Gu, Rohan Varma, Liang Luo, Chien-Chin Huang, Min Xu, Less Wright, Hamid Shojanazeri, Myle Ott, Sam Shleifer, et al. Pytorch fsdp: experiences on scaling fully sharded data parallel. *arXiv preprint arXiv:2304.11277*, 2023. 3

[53] Yanqi Zhou, Tao Lei, Hanxiao Liu, Nan Du, Yanping Huang, Vincent Zhao, Andrew M Dai, Quoc V Le, James Laudon, et al. Mixture-of-experts with expert choice routing. *Advances in Neural Information Processing Systems*, 35:7103–7114, 2022. 3

## A. Score Matching Derivation

We provide an alternative derivation of Decentralized Flow Matching based on score matching [45] rather than flow matching [26]. We begin with the score, which is the gradient of the log likelihood, $p_t(x_t)$.

$$\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t) \tag{9}$$

By applying the chain rule, this can be expressed in terms of the derivative of the likelihood itself:

$$\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t) = \frac{1}{p_t(\mathbf{x}_t)} \cdot \nabla_{\mathbf{x}_t} p_t(\mathbf{x}_t). \tag{10}$$

Let us assume our data, $x_0$, is generated in a bi-level fashion:

$$\mathbf{x}_0 \sim p_0(\mathbf{x}_0|\mathbf{v}), \quad \mathbf{v} \sim p_\mathbf{v}(\mathbf{v}), \tag{11}$$

where $v$ is the cluster label discussed in the DDM method and $p_\mathbf{v}(\mathbf{v})$ follows a distribution defined by the clustering procedure. The marginal likelihood $p_t(\mathbf{x}_t)$ can then be expressed by integrating over $\mathbf{v}$:

$$\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t) = \frac{1}{p_t(\mathbf{x}_t)} \cdot \nabla_{\mathbf{x}_t} \sum_\mathbf{v} p(\mathbf{v}) \cdot p_t(\mathbf{x}_t|\mathbf{v}). \tag{12}$$

By linearity of differentiation, we distribute the gradient over the summation:

$$= \frac{1}{p_t(\mathbf{x}_t)} \cdot \sum_\mathbf{v} p(\mathbf{v}) \cdot \nabla_{\mathbf{x}_t} p_t(\mathbf{x}_t|\mathbf{v}). \tag{13}$$

Since the gradient of a log probability can be expressed as the probability multiplied by the gradient of its log,

$$= \frac{1}{p_t(\mathbf{x}_t)} \cdot \sum_\mathbf{v} p(\mathbf{v}) p_t(\mathbf{x}_t|\mathbf{v}) \cdot \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|\mathbf{v}). \tag{14}$$

Finally, we invoke Bayes' Theorem:

$$= \sum_\mathbf{v} p_t(\mathbf{v}|\mathbf{x}_t) \cdot \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|\mathbf{v}) \tag{15}$$

This result mirrors the flow matching derivation, showing that the score prediction for the overall data distribution can be recast as a linear combination of score predictions for each data cluster. Each learned expert predicts its own conditional score, which ensemble according to the posterior probability of the expert label given the latent $\mathbf{x}_t$.

## B. Additional Training and Evaluation Details

Our split of LAION Aesthetics contains 153.6 million image-caption pairs. We pretrain all diffusion models on 256x256 square crop images encoded through Huggingface's fine-tuned Stable Diffusion VAE (sd-vae-ft-mse). This encoder employs an 8× spatial downsampling factor. Throughout training, we maintain a patch size of 2 for best quality, resulting in a pretraining context length of 256.

For high-resolution finetuning, we choose five aspect ratio buckets to handle varying image dimensions while maintaining consistent tokenized sequence length (3600). The bucket are as follows:

- **1280 × 720** (16:9 landscape)
- **1200 × 768** ( 3:2 landscape)
- **960 × 960** (square)
- **768 × 1200** ( 2:3 portrait)
- **720 × 1280** (9:16 portrait)

Images are mapped to their nearest matching bucket by aspect ratio. Following best practices from Stable Diffusion 3 [15], we adjust the timestep schedule for high-resolution training and inference by applying a log-SNR shift of 3 [22]. We also modify the Rotary Positional Embedding (RoPE) in the MMDiT architecture by interpolating RoPE inputs within the central square region and extrapolating for peripheral areas.

For evaluation, we use standard classifier-free guidance scales: 7.5 for LAION text-conditional generation and 3 for ImageNet class-conditional generation. All evaluations use 50 sampling steps to ensure consistent comparisons. We compute FID, CLIP-FID and DINO-FID metrics on fixed evaluation splits to standardize evaluations.

## C. Additional Quantitative Analysis

Our additional quantitative analysis explores key test-time DDM hyperparameters. We test various ensemble combinations in Table 2, including nucleus sampling which is common in LLM decoding. Top-1 sampling consistently delivers the best performance while being the most computationally efficient. This finding holds across different expert counts, router temperatures, and threshold probabilities.

Our classifier-free guidance (CFG) [19] scale experiments (Figure 8) show that DDMs respond similarly to monolith models, suggesting that standard CFG scales can be directly applied to DDMs. Additionally, our training efficiency analysis (Figures 8e and 8f) confirms that distillation achieves comparable generation quality (FID) with only one-quarter of the monolith's batch size (256 vs 1024).

## D. Additional Qualitative Results

We provide additional selected samples from our largest DDM ensemble in Figure 9 as well as random samples for different text prompts in Figures 10 through 17.
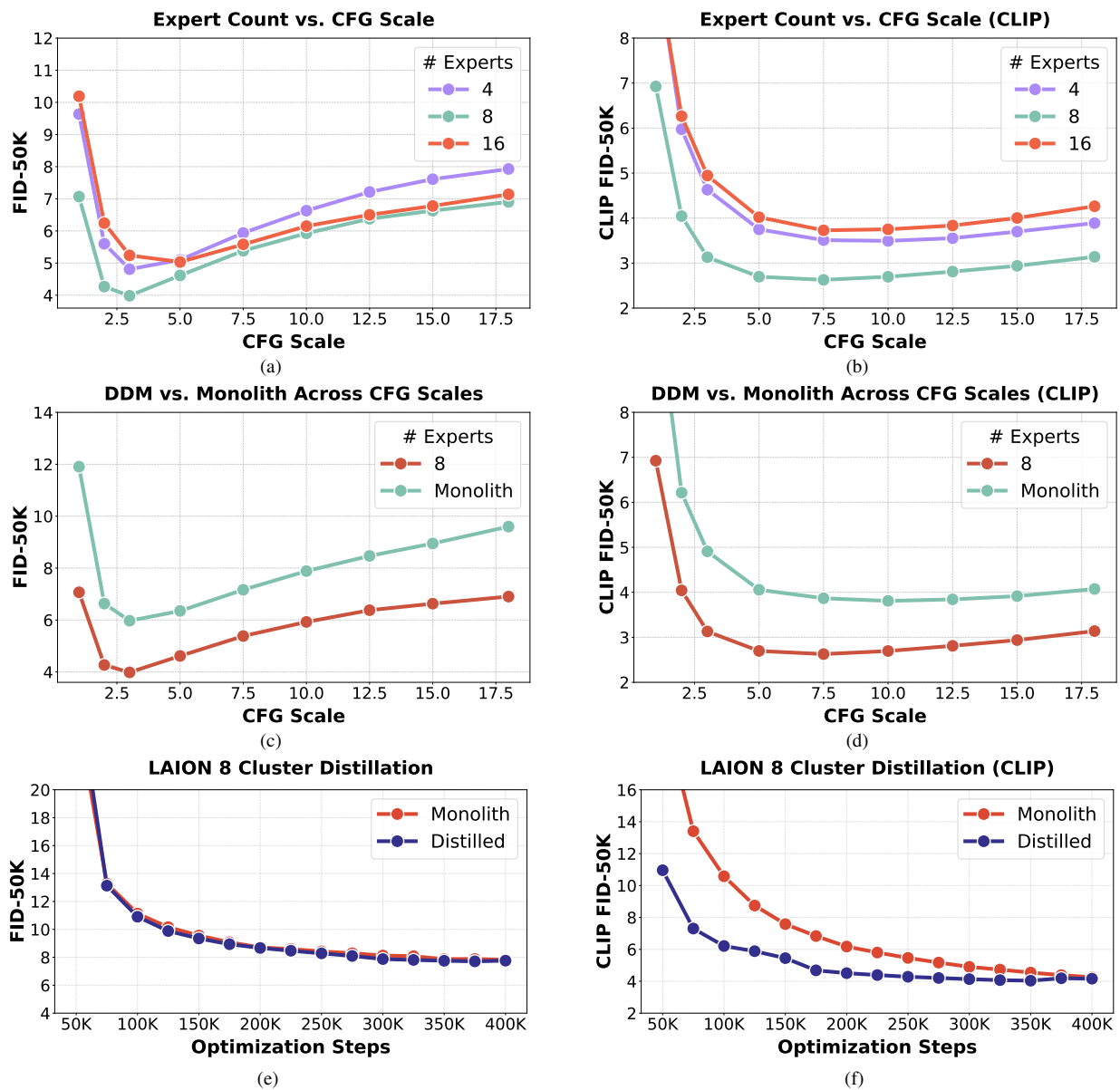
Figure 8. **Additional Quantitative Analysis.** We sweep CFG scales across decentralized and monolith diffusion models trained on LAION Aesthetics (a, b, c, d), finding that optimal CFG scales are consistent across models. Distillation matches the performance of training a monolith from raw data at a fraction of the FLOP-cost (e, f).

| Inference Strategy | Expert Count | Temp. | Active Experts | $p$ | GFLOPs ↓ | FID ↓ | CLIP FID ↓ | DINO FID ↓ |
|---|---|---|---|---|---|---|---|---|
| Monolith | 1 | - | 1 | 0.00 | 308 | 12.81 | 5.58 | 343.96 |
| Full | 4 | - | 4 | 0.00 | 1245 | 12.75 | 6.82 | 386.3 |
| Top-1 | 4 | - | 1 | 0.00 | **334** | **12.54** | 6.72 | **378.4** |
| Top-2 | 4 | - | 2 | 0.00 | 642 | 12.76 | 6.75 | 384.12 |
| Top-3 | 4 | - | 3 | 0.00 | 950 | 12.88 | 6.8 | 385.79 |
| Sample | 4 | 0.5 | 1 | 0.00 | **334** | 117.02 | 36.25 | 1321.17 |
| Sample | 4 | 1.0 | 1 | 0.00 | **334** | 89.14 | 28.01 | 1042.09 |
| Sample | 4 | 2.0 | 1 | 0.00 | **334** | 15.08 | 7.53 | 425.29 |
| Sample | 4 | 0.5 | 2 | 0.00 | 642 | 12.67 | **6.71** | 380.85 |
| Sample | 4 | 1.0 | 2 | 0.00 | 642 | 12.67 | 6.73 | 382.93 |
| Sample | 4 | 2.0 | 2 | 0.00 | 642 | 13.51 | 7.05 | 400.08 |
| Sample | 4 | 0.5 | 3 | 0.00 | 950 | 12.57 | 6.76 | 381.44 |
| Sample | 4 | 1.0 | 3 | 0.00 | 950 | 12.84 | 6.81 | 385.18 |
| Sample | 4 | 2.0 | 3 | 0.00 | 950 | 13.59 | 7.08 | 400.23 |
| Nucleus | 4 | 0.5 | 1 | 0.90 | **334** | 15.74 | 9.99 | 412.49 |
| Nucleus | 4 | 1.0 | 1 | 0.90 | **334** | 15.72 | 10.04 | 411.31 |
| Nucleus | 4 | 2.0 | 1 | 0.90 | **334** | 17.35 | 10.31 | 432.46 |
| Threshold | 4 | 1.0 | - | 0.01 | - | 12.82 | 6.81 | 385.44 |
| Threshold | 4 | 1.0 | - | 0.05 | - | 12.67 | 6.73 | 382.54 |
| Threshold | 4 | 1.0 | - | 0.10 | - | 12.63 | 6.76 | 382.86 |
| Full | 8 | - | 4 | 0.00 | 2490 | 10.52 | 5.85 | 354.15 |
| Top-1 | 8 | - | 1 | 0.00 | **334** | **9.85** | **5.54** | **339.56** |
| Top-2 | 8 | - | 2 | 0.00 | 642 | 10.33 | 5.73 | 349.28 |
| Top-3 | 8 | - | 3 | 0.00 | 950 | 10.45 | 5.77 | 351.91 |
| Sample | 8 | 1.0 | 1 | 0.00 | **334** | 190.95 | 59.03 | 2105.79 |
| Sample | 8 | 2.0 | 1 | 0.00 | **334** | 184.06 | 50.55 | 1790.24 |
| Sample | 8 | 0.5 | 2 | 0.00 | 642 | 9.93 | 5.57 | 343.51 |
| Sample | 8 | 1.0 | 2 | 0.00 | 642 | 10.28 | 5.72 | 348.39 |
| Sample | 8 | 2.0 | 2 | 0.00 | 642 | 17.11 | 8.09 | 471.18 |
| Sample | 8 | 0.5 | 3 | 0.00 | 950 | 10.04 | 5.62 | 342.86 |
| Sample | 8 | 1.0 | 3 | 0.00 | 950 | 10.42 | 5.78 | 350.91 |
| Sample | 8 | 2.0 | 3 | 0.00 | 950 | 12.06 | 6.38 | 380.54 |
| Nucleus | 8 | 0.5 | 1 | 0.90 | **334** | 188.66 | 60.09 | 2110.22 |
| Nucleus | 8 | 1.0 | 1 | 0.90 | **334** | 152.16 | 48.37 | 1609.23 |
| Nucleus | 8 | 2.0 | 1 | 0.90 | **334** | 33.9 | 14 | 682.31 |
| Threshold | 8 | 1.0 | - | 0.01 | - | 10.51 | 5.82 | 351.17 |
| Threshold | 8 | 1.0 | - | 0.05 | - | 10.32 | 5.73 | 349.86 |
| Threshold | 8 | 1.0 | - | 0.10 | - | 10.18 | 5.7 | 346.9 |
| Full | 16 | - | 4 | 0.00 | 4980 | 15.43 | 7.57 | 440.54 |
| Top-1 | 16 | - | 1 | 0.00 | **334** | **12.51** | **6.6** | **397.99** |
| Top-2 | 16 | - | 2 | 0.00 | 642 | 148.26 | 41.13 | 1535.85 |
| Top-3 | 16 | - | 3 | 0.00 | 950 | 91.76 | 29.53 | 1105.92 |
| Sample | 16 | 1.0 | 1 | 0.00 | **334** | 232.1 | 71.88 | 2557.41 |
| Sample | 16 | 2.0 | 1 | 0.00 | **334** | 259 | 81.49 | 2797.76 |
| Sample | 16 | 0.5 | 2 | 0.00 | 642 | 161.29 | 47.76 | 1732.88 |
| Sample | 16 | 1.0 | 2 | 0.00 | 642 | 174.23 | 54.41 | 1941.49 |
| Sample | 16 | 0.5 | 3 | 0.00 | 950 | 119.84 | 40.18 | 1510.78 |
| Sample | 16 | 1.0 | 3 | 0.00 | 950 | 44.62 | 20.04 | 772.02 |
| Sample | 16 | 2.0 | 3 | 0.00 | 950 | 26.01 | 10.92 | 603.31 |
| Threshold | 16 | 1.0 | - | 0.01 | - | 14.92 | 7.44 | 431.89 |
| Threshold | 16 | 1.0 | - | 0.05 | - | 12.62 | 6.61 | 399.09 |
| Threshold | 16 | 1.0 | - | 0.10 | - | 12.69 | 6.61 | 398.33 |

Table 2. **Test-Time Combination Strategies.** We ablate strategies and relevant hyperparameters for sampling from our ImageNet DDM ensemble at test-time. Across many experiments, we find that simply selecting the top expert outperforms more sophisticated alternatives.

Figure 9. **Additional Selected Samples.**

Figure 10. **Random Samples, Fixed Prompt.** a photo of the dolomites

Figure 11. **Random Samples, Fixed Prompt.** 1969 Polaris Colt, restored to showroom, static display in snow, winter sunrise

Figure 12. **Random Samples, Fixed Prompt.** weather research station in extreme conditions, monitoring equipment, natural elements

18

Figure 13. **Random Samples, Fixed Prompt.** ancient bristlecone pine forest, twisted trees, high-altitude light, rugged mountain backdrop

Figure 14. **Random Samples, Fixed Prompt.** deep desert slot canyon, sandstone textures, shaft of light, natural color gradients

Figure 15. **Random Samples, Fixed Prompt.** gothic cathedral spires piercing morning mist, ancient European city roofscape

Figure 16. **Random Samples, Fixed Prompt.** historic textile mill interior, preserved machinery, sunbeams through industrial windows

Figure 17. **Random Samples, Fixed Prompt.** symphony orchestra during rehearsal, conductor's perspective, historic concert hall