# DiffMoE: Dynamic Token Selection for Scalable Diffusion Transformers

Minglei Shi [1*] Ziyang Yuan [1*] Haotian Yang [2] Xintao Wang [2†] Mingwu Zheng [2] Xin Tao [2] Wenliang Zhao [1] Wenzhao Zheng [1] Jie Zhou [1] Jiwen Lu [1†] Pengfei Wan [2] Di Zhang [2] Kun Gai [2]

**Project Page:** https://shiml20.github.io/DiffMoE/

## Abstract

Diffusion Transformers (DiTs) have emerged as the dominant architecture for visual generation tasks, yet their uniform processing of inputs across varying conditions and noise levels fails to leverage the inherent heterogeneity of the diffusion process. While recent mixture-of-experts (MoE) approaches attempt to address this limitation, they struggle to achieve significant improvements due to their restricted token accessibility and fixed computational patterns. We present **DiffMoE**, a novel MoE-based architecture that enables experts to access global token distributions through a **batch-level global token pool** during training, promoting specialized expert behavior. To unleash the full potential of inherent heterogeneity, DiffMoE incorporates a **capacity predictor** that dynamically allocates computational resources based on noise levels and sample complexity. Through comprehensive evaluation, DiffMoE achieves state-of-the-art performance among diffusion models on ImageNet benchmark, substantially outperforming both dense architectures with $3\times$ activated parameters and existing MoE approaches while maintaining $1\times$ activated parameters. The effectiveness of our approach extends beyond class-conditional generation to more challenging tasks such as text-to-image generation, demonstrating its broad applicability across different diffusion model applications.

## 1. Introduction

The Mixture-of-Experts (MoE) framework (Shazeer et al., 2017; Lepikhin et al., 2020) has emerged as a powerful paradigm for enhancing overall multi-task performance while maintaining computational efficiency. This
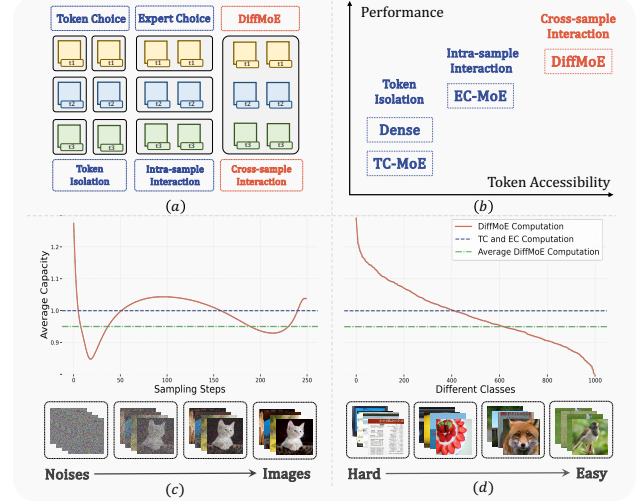


Figure 1: **Token Accessibility and Dynamic Computation.** **(a)** Token accessibility levels from token isolation to cross-sample interaction. Colors represent tokens in different samples, $t_i$ indicates noise levels. **(b)** Performance-accessibility analysis across architectures. **(c)** Computational dynamics during diffusion sampling, showing adaptive computation from noise to image. **(d)** Class-wise computation allocation from hard (technical diagrams) to easy (natural photos) tasks. Results from DiffMoE-L-E16-Flow (700K).

is achieved by combining multiple expert networks, each focusing on a distinct task, with their outputs integrated through a gating mechanism. In language modeling, MoE has achieved performance comparable to dense models of $2\times - 3\times$ activated parameters (DeepSeek-AI et al., 2024; MiniMax et al., 2025; Muennighoff et al., 2024). The current MoE primarily follows two gating paradigms: Token-Choice (TC), where each token independently selects a subset of experts for processing; and Expert-Choice (EC), where each expert selects a subset of tokens from the sequence for processing.

Diffusion (Ho et al., 2020; Rombach et al., 2022; Podell et al., 2023; Song et al., 2021) and flow-based (Ma et al., 2024; Esser et al., 2024b; Liu et al., 2023) models inherently represent multi-task learning frameworks, as they process varying token distributions across different noise levels and conditional inputs. While this heterogeneity characteristic

---
[*]Equal contribution [1]Tsinghua University, Beijing, China [2]Kuaishou Technology, Beijing, China. Correspondence to: Xintao Wang <xintao.wang@kuaishou.com>, Jiwen Lu <lujiwen@tsinghua.edu.cn>.

naturally aligns with the MoE framework's ability for multi-task handling, existing attempts (Fei et al., 2024; Sun et al., 2024a; Yatharth Gupta, 2024; Sehwag et al., 2024) to integrate MoE with diffusion models have yielded suboptimal results, failing to achieve the remarkable improvements observed in language models. Specifically, Token-choice MoE (TC-MoE) (Fei et al., 2024) often underperforms compared to conventional dense architectures under the same number of activations; Expert-choice MoE (EC-MoE) (Sehwag et al., 2024; Sun et al., 2024a) shows marginal improvements over dense models, but only when trained for much longer.

We are curious about what fundamentally limits MoE's effectiveness in diffusion models. Our key finding reveals that **global token distribution accessibility is crucial for MoE success in diffusion models, necessitating the model learn and dynamically process the tokens from different noise levels and conditions**, as illustrated in Figure 1. Previous approaches have neglected this crucial component, resulting in compromised performance. Specifically, Dense models and TC-MoE **isolates tokens**, preventing them from interacting with others during expert selection, while EC-DiT restricts **intra-sample token interaction** , which fails to access other samples with different noise levels and conditions. These limitations hinder the model's ability to capture the full spectrum of the heterogeneity inherent in diffusion processes.

To address these limitations, we introduce **DiffMoE**, a novel architecture that features a **batch-level global token pool** for enhanced **cross-sample token interaction** during training, as illustrated in Figure 2. This approach approximates the complete token distribution across different noise levels and samples, facilitating more specialized expert learning through comprehensive global token information access. Our empirical analysis demonstrates that the global token pool accelerates loss convergence, surpassing dense models with equivalent activation parameters. Though some concurrent works in large language models (DeepSeek-AI et al., 2024; Qiu et al., 2025) show similar principles in global batch, we argue that these principles are particularly crucial for diffusion transformers due to their inherently more complex heterogeneity nature.

However, conventional MoE inference strategies, which maintain fixed computational resource allocation across different noise levels and conditions, fail to fully leverage the potential of DiffMoE's batch-level global token pool. To optimize token selection during inference, we propose a **capacity predictor** that dynamically adjusts resource allocation. This adaptive mechanism learns from training-time token routing patterns, efficiently distributing computational resources between complex and simple cases. Furthermore, we implement a **dynamic threshold** at inference time to

achieve flexible performance-computation trade-offs.

By integrating the global token pool and capacity predictor, **DiffMoE achieves superior performance over dense models with $3\times$ activated parameters** while maintaining efficient scaling properties (See Table 5). Our approach offers extra several advantages over existing methods: it eliminates the potentially detrimental load balancing losses present in TC-MoE and overcomes the intra-sample token selection constraints of EC-MoE, resulting in enhanced flexibility and scalability. Extensive empirical evaluations demonstrate DiffMoE's superior scaling efficiency and performance improvements across diverse diffusion applications.

Our contributions can be summarized as follows: (1) We identify the fundamental importance of global token distribution accessibility in facilitating dynamic token selection for MoE-based diffusion models; (2) We introduce DiffMoE, a novel framework incorporating a global token pool and capacity predictor to enable efficient model scaling; (3) We achieve state-of-the-art performance on ImageNet benchmark among diffusion models. through dynamic computation allocation while preserving computational efficiency; and (4) We conduct comprehensive experiments that validate our approach's effectiveness across diverse diffusion applications.

## 2. Method

### 2.1. Preliminaries

**Diffusion Models.** Diffusion models (Ho et al., 2020; Rombach et al., 2022; Sohl-Dickstein et al., 2015; Song et al., 2021) are a powerful family of generative models, which can transform the noise distribution $p_1(\mathbf{x})$ to the data distribution $p_0(\mathbf{x})$. The diffusion process can be represented as: $\mathbf{x}_t = \alpha_t \mathbf{x}_0 + \sigma_t \boldsymbol{\epsilon}, \quad t \in [0, 1], \quad \boldsymbol{\epsilon} \sim \mathcal{N}(0, \mathbf{I})$, Where $\alpha_t$ and $\sigma_t$ are monotonically decreasing and increasing functions of $t$, respectively. The marginal distribution $p_1(\mathbf{x})$ converges to $\mathcal{N}(0, \mathbf{I})$, when $\alpha_1 = \sigma_0 = 0, \alpha_0 = \sigma_1 = 1$.

To train a diffusion model, we can use the denoising score matching method (Song et al., 2021) which constructs a score prediction model $\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)$ to estimate the scaled score function $-\sigma_t \nabla_{\mathbf{x}} \log p_t(\mathbf{x}_t)$ with training objective formulated in Eq. 22. Sampling from a diffusion model can be achieved by solving the reverse-time SDE or the corresponding diffusion ODE (Song et al., 2021) in an iterative manner. Recently, flow-based models (Liu et al., 2022a; Lipman et al., 2022; Esser et al., 2024a) have shown superior performance through alternative training objective formulated in Eq. 31 while maintaining the same architecture as DiT (Peebles & Xie, 2023a). Sampling from a flow-based model can be achieved by solving the probability flow ODE.

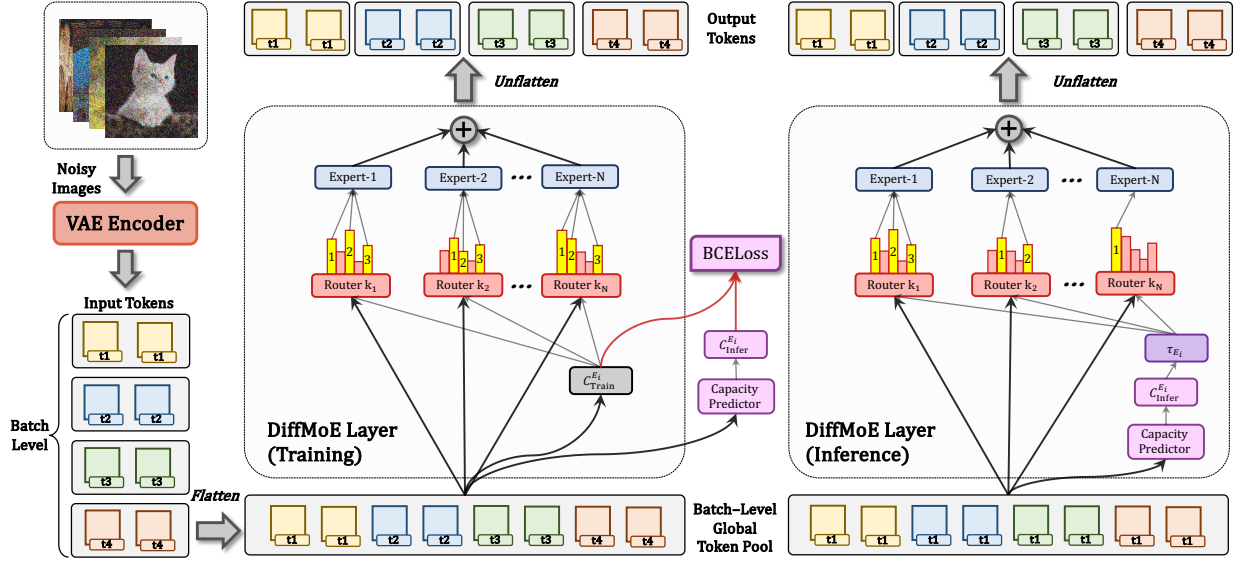**Mixture of Experts.** Mixture of Experts (MoE) (Shazeer

Figure 2: **DiffMoE Architecture Overview.** DiffMoE flattens tokens into a batch-level global token pool, where each expert maintains a fixed training capacity of $C_{\text{train}}^{E_i} = 1$. During inference, a dynamic capacity predictor adaptively routes tokens across different sampling steps and conditions. Different colors denote tokens from distinct samples, while $t_i$ represents corresponding noise levels.

et al., 2017; Cai et al., 2024) is based on a fundamental insight: different parts of a model can specialize in handling distinct tasks. By selectively activating only relevant components, MoE enables efficient scaling of model capacity while maintaining computational efficiency.

MoE layers generally consist of $N$ experts, each implemented as a Feed-Forward Network (FFN) with identical architecture, denoted by $E_1(\mathbf{x}), \ldots, E_N(\mathbf{x})$ with input $\mathbf{x}$. A routing matrix $\mathbf{W}_r \in \mathbb{R}^{D \times N}$ is used to calculate token-expert affinity matrix:

$$\mathbf{M} = \text{softmax}_E(\mathbf{x}\mathbf{W}_r), \quad \mathbf{x} \in \mathbb{R}^{B \times S \times D}, \quad (1)$$

where $B$ is the batch size, $S$ is the token length of one sample, $D$ is the hidden dimension, $\text{softmax}_E$ denotes the $\text{softmax}$ operation along the expert axis. There are two common gating paradigms: Token-Choice (TC) (Shazeer et al., 2017; Fei et al., 2024) and Expert-Choice (EC) (Zhou et al., 2022; Sun et al., 2024a). For TC, each token of each sample individually selects $\text{top-}K$ experts via a gating function, the gating function and output of TC-MoE layers are defined as follows:

$$\mathbf{G}_{s,i}^{TC} = \begin{cases} \mathbf{M}_{s,i}, & \mathbf{M}_{s,i} \in \text{top-}K(\{\mathbf{M}_{s,i}\}_{i=1}^N) \\ 0, & \text{otherwise} \end{cases} \quad (2)$$

$$\mathbf{y}_s = \sum_{i=1}^{K} \mathbf{G}_{s,i}^{TC} E(\mathbf{x}_s), \mathbf{x}_s \in \mathbb{R}^{1 \times D}, s \in \{1, \ldots, S\}. \quad (3)$$

Different from TC, EC makes every expert selects $K'$ tokens

from each sample of the input $\mathbf{x} \in \mathbb{R}^{B \times S \times D}$. The gating function $G_{s,i}^{EC}$ can be implemented analogously to Eq. 2, with the modification that the $\text{top}$ operation selects $K'$ tokens along the token length dimension, *i.e.* $S$. Similar to Eq. 3, the output for a token $\mathbf{x}_s$ of EC-MoE layers can be calculated as: $\mathbf{y}_s = \sum_{i=1}^{N} \mathbf{G}_{s,i}^{EC} E(\mathbf{x}_s), \mathbf{x}_s \in \mathbb{R}^{1 \times D}$. Both TC and EC struggle to achieve significant improvements comparing with dense models due to their restricted token accessibility and fixed computational patterns.

### 2.2. DiffMoE: Dynamic Token Selection

**Batch-level Global Token Pool.** Since MoE architectures replace FFN layers, both TC and EC paradigms in diffusion models are inherently limited to processing tokens within individual samples, where gating mechanisms operate exclusively on tokens sharing identical conditions and noise levels. This architectural constraint not only prevents experts from learning crucial contrastive patterns but, more fundamentally, restricts their access to the global token distribution that characterizes the full spectrum of the diffusion process. To capture this essential global context, we introduce a *Batch-level Global Token Pool* for DiffMoE by flattening batch and token dimensions, enabling experts to access a comprehensive token distribution spanning different noise levels and conditions. This design, which simulate the true token distribution of the entire dataset during training, can be formulated as follows:

$$\mathbf{x} \in \mathbb{R}^{B \times S \times D} \rightarrow \mathbf{x}_{\text{pool}} \in \mathbb{R}^{BS \times D}. \quad (4)$$
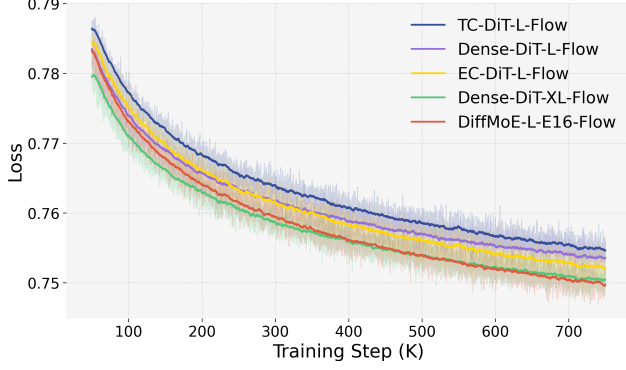
3

Figure 3: **Training Loss Curves of Different Flow-based Models.** DiffMoE with batch-level global token pool achieves consistently lower diffusion losses than baselines without batch-level global token pool.

During the training phase, we push expert $E$ to select $K_{\text{train}}^E$ tokens, forcing each expert to capture the characteristics of tokens from different conditional information and noise levels, while keeping expert load balance during training. The corresponding batch-level global token-expert affinity matrix will be calculated as follows:

$$\mathbf{M}^{Dy} = \mathbf{x}_{\text{pool}}\mathbf{W}_r, \mathbf{x}_{\text{pool}} \in \mathbb{R}^{BS \times D}, \mathbf{W}_r \in \mathbb{R}^{D \times N}. \quad (5)$$

Then, using $\mathbf{M}^{Dy} \in \mathbb{R}^{BS \times N}$, the gating value of MoE and the output of DiffMoE layers can be computed as follows:

$$\mathbf{G}_{s,i}^{Dy} = \begin{cases} \mathbf{M}_{s,i}^{Dy}, & \mathbf{M}_{s,i}^{Dy} \in \texttt{top-}K_{\text{train}}^E(\{\mathbf{M}_{s,i}^{Dy}\}_{s=1}^{BS}) \\ 0, & \text{otherwise} \end{cases} \quad (6)$$

$$\mathbf{y}_s = \sum_{i=1}^{N} \mathbf{G}_{s,i}^{Dy} E(\mathbf{x}_s). \quad (7)$$

As shown in Figure 3, DiffMoE consistently achieves lower diffusion losses than all baselines.

**Capacity and Computational Cost.** To establish a rigorous and fair comparison framework, we define the capacity $C^E$ for a single expert $E$, which serves as a standardized metric for quantifying computational costs. This capacity metric enables fair comparisons between DiffMoE and baseline models by accurately measuring the computational resources utilized by each expert:

$$C^E = \frac{N \times \text{\# tokens processed by } E}{\text{\# all input tokens}} = \frac{NK^E}{BS}, \quad (8)$$

where $K^E$ denotes the number of tokens assigned to expert $E$, $N$ is the number of experts, and $BS$ represents the size of global token pool. Here, we define the capacity $C$ for one forward process for both training and inference phases:

$$C = \frac{1}{LN} \sum_{l=1}^{L} \sum_{i=1}^{N} C^{E_i^l}, \quad (9)$$

where $E_i^l$ denotes $i_{\text{th}}$ expert in the $l_{\text{th}}$ MoE layer.

During training phase, $C_{\text{train}}$ is fixed to 1 across all the MoE models, indicating that they keep the same computational cost as dense models. Specifically, DiffMoE keeps $K_{\text{train}}^E = BS/N$ for $C_{\text{train}} = 1$. TC-DiT selects the top-1 expert, while EC-DiT selects top-$(S/N)$ tokens per sample in a batch to ensure the same computation. During the inference phase, we compute the global average inference capacity $C_{\text{infer}}^{\text{avg}}$ by averaging over all timesteps: $C_{\text{infer}}^{\text{avg}} = \frac{1}{T} \sum_{t=1}^{T} C_{\text{infer}}^t$, where $C_{\text{infer}}^t$ represents the inference capacity at sampling step $t$.

**Capacity Predictor.** Although batch-level global token routing enables efficient model training, conventional MoE inference strategies with fixed computational resource allocation fail to fully leverage its potential. This limitation stems from the static resource distribution across different noise levels and conditional information during inference. To optimize token selection, we propose a *capacity predictor*—a lightweight structure that dynamically determines token selection per expert through a two-layer MLP with SiLU activations. This adaptive mechanism learns from training-time token routing patterns, efficiently distributing computational resources between complex and simple cases. Formally, let $\text{CP}(\mathbf{x}_{\text{pool}}) \in \mathbb{R}^{BS \times N}$ denote the predictor's output for input $\mathbf{x}_{\text{pool}} \in \mathbb{R}^{BS \times D}$:

$$\text{CP}(\mathbf{x}_{\text{pool}}) = \mathbf{W}_2 \sigma_{\text{SiLU}}(\mathbf{W}_1 \mathbf{x}_{\text{pool}}). \quad (10)$$

We can optimize the capacity predictor by minimizing the object function:

$$\mathcal{L}_{\text{CP}} = BCELoss(\mathbf{O}, \text{CP}(\texttt{sg}[\mathbf{x}_{\text{pool}}])) \quad (11)$$

$$= -\frac{1}{LBSN} \sum_{l=1}^{L} \sum_{i=1}^{BS} \sum_{j=1}^{N} \{\mathbf{O}_{i,j}^l \log(\text{CP}(\texttt{sg}[\mathbf{x}_{\text{pool}}])_{i,j}) + (1 - \mathbf{O}_{i,j}^l) \log(1 - \text{CP}(\texttt{sg}[\mathbf{x}_{\text{pool}}])_{i,j})\}, \quad (12)$$

where sg denotes the stop-gradient operation, and $\mathbf{O} \in \mathbb{R}^{L \times BS \times N}$ is defined as follows:

$$\mathbf{O}_{s,i}^l = \begin{cases} 1, & \text{if } \mathbf{x}_{\text{pool},s} \text{ is processed by } E_i^l \\ 0, & \text{otherwise} \end{cases}. \quad (13)$$

We employ the stop-gradient technique to train the capacity predictor, ensuring it focuses solely on the input features at the current layer while preventing it from interfering with the training of the main diffusion transformers. Therefore, $\mathcal{L}_{\text{CP}}$ will not affect actual diffusion loss. During inference, the capacity predictor determines the inference capacity $C_{\text{infer}}^{E_i^l, t}$ for each expert $E_i^l$ at timestep $t$ based on a threshold. Let $\tau_{E_i^l}$ denote the threshold of $E_i^l$. Using $\mathcal{T} = \{\tau_{E_i^l} \mid i \in \{1, \dots, N\}, l \in \{1, \dots, L\}\}$, the model

achieves an adaptive $C_{\text{infer}}^{E_i^l,t}$ allocation at sampling step $t$ tailored to different input tokens as follows:

$$C_{\text{infer}}^{E_i^l,t}(\tau_{E_i^l}) = \sum_{s=1}^{S} \delta_{s,i}^l$$

$$\text{where } \delta_{s,i}^l = \begin{cases} 1, & \text{if } \text{CP}(\mathbf{x}_{\text{pool}})_{s,i} > \tau_{E_i^l} \\ 0, & \text{otherwise} \end{cases}. \quad (14)$$

Then we can calculate $C_{\text{infer}}^{\text{avg}}$ w.r.t. $\mathcal{T}$

$$C_{\text{infer}}^{\text{avg}}(\mathcal{T}) = \frac{1}{TNL} \sum_{t=1}^{T} \sum_{i=1}^{N} \sum_{l=1}^{L} C_{\text{infer}}^{E_i^l,t}(\tau_{E_i^l}). \quad (15)$$

**Dynamic Threshold.** We can set the threshold $\mathcal{T}$ to control the number of tokens processed by each expert during inference. It is evident that #tokens processed during inference, decreases as $\tau_{E_i^l}$ increases for all expert. We can adjust $\mathcal{T}$ flexibly to achieve a better trade-off between computational complexity and generation quality. We employ two distinct approaches for threshold $\mathcal{T}$ determination: *Interval Search* and *Dynamic Threshold*. The interval search method addresses an optimization problem formulated as follows:

$$\min_{\mathcal{T}} \text{FID}(\mathcal{T}) \quad \text{subject to } C_{\text{infer}}^{\text{avg}}(\mathcal{T}) \leq 1. \quad (16)$$

To simplify the optimization problem, we assume that $\tau = \gamma < 1, \forall \tau \in \mathcal{T}$, where $\gamma$ is a constant in our experiments.

However, interval search method is labor-intensive and time-consuming, making it impractical for real-world applications. To address this limitation, we propose a dynamic threshold method that automatically maintains thresholds (denoted as $\mathcal{T}^{Dy} = \{\tau_{E_i^l}^{Dy} \mid i \in \{1,\dots,N\}, l \in \{1,\dots,L\}\}$) for all experts during the training phase. To ensure the inference computational cost approximates the training cost (*i.e.* $C_{\text{infer}}^{\text{avg}} \approx 1$), we employ the Exponential Moving Average (EMA) technique as follows:

$$\text{Quantile}_{E_i^l} \leftarrow \text{CP}(\mathbf{x}_{\text{pool}})_{s_k,i},$$
$$\tau_{E_i^l}^{Dy} \leftarrow \alpha \cdot \tau_{E_i^l}^{Dy} + (1-\alpha) \cdot \text{Quantile}_{E_i^l}, \quad (17)$$

where $s_k$ denotes the $k_{\text{th}}$ value in descending order, $\alpha$ is a constant which is equals to 0.95 in our experiments.

# 3. Experiments

We evaluate DiffMoE on class-conditional image generation across three key aspects: (1) **Training and Inference Performance**, (2) **Dynamic Computation**, and (3) **Scalability**. Our experiments demonstrate DiffMoE's effectiveness through extensive analysis. Additionally, we verify its adaptability on text-to-image generation tasks.
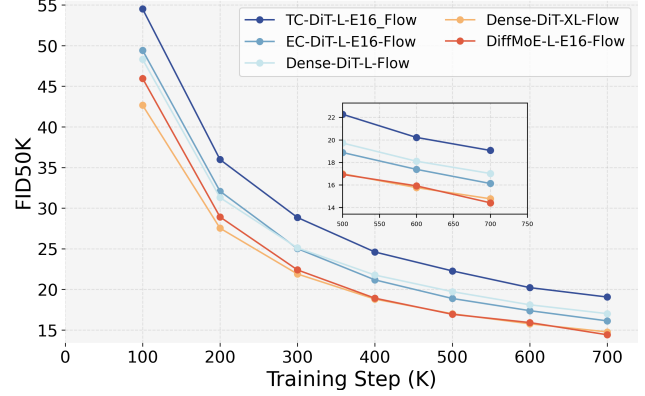


Figure 4: **Comparisons with the Baseline Models.** We compare TC, EC, and Dense Models. DiffMoE-L-E16-Flow even surpasses the DenseDiT-XL-Flow (1.5x params) by achieving the best quality (14.41 FID50K w/o CFG at 700K). The results of the DDPM method remain consistent with those provided in the Appendix D.1.

## 3.1. Experiment Setup

**Baseline and Model architecture.** We compare with Dense-DiT trained by denoising score matching (Peebles & Xie, 2023a) and flow matching (Ma et al., 2024), TC-DiT (Fei et al., 2024), and EC-DiT (Sun et al., 2024a). For fair comparison, we reimplement TC-DiT and EC-DiT based on their public repositories and pseudo-codes while maintaining identical activated computation. Models are named as: **[Model]-[Size]-[# Experts]-[Training Type]**. For class-conditional generation, we replace even FFN layers with MoE layers containing $N$ identical FFN components (Lepikhin et al., 2020), while maintaining the original DiT architecture (Peebles & Xie, 2023a), details of architecture are shown in Table 1. For text-to-image generation, we add cross-attention modules (Rombach et al., 2022), where DiffMoE-E16-T2I-Flow activates 1.2B parameters (matching Dense-DiT-T2I-Flow) from total 4.6B parameters. Full details are in Appendix B.1.

**Evaluation.** We evaluated DiffMoE through both quantitative and qualitative metrics. Quantitatively, we used FID50K (Heusel et al., 2017) with 250 DDPM/Euler steps for class-conditional generation, and compared with SiT (Ma et al., 2024) using Heun sampler at 125 steps. For text-to-image generation, we employed GenEval metrics (Ghosh et al., 2023). Training losses were analyzed to validate the model's convergence behavior. Additionally, we assessed the model's performance through visual inspection of samples generated from diverse prompts.

## 3.2. Main Results: Class-conditional Image Generation

Class-conditional image generation is a task of synthesizing images based on specified class labels.

Table 1: **DiffMoE Model Configurations**. Hyperparameter settings and computational specifications for class-conditional models. See Appendix C for activated parameter calculations.

| Model Config | #Avg. Activated Params ($C_{\text{infer}}^{\text{avg}} = 1$). | #Total Params. | #Blocks L | #Hidden dim. D | #Head n | #Experts | $C_{\text{train}}$ |
|---|---|---|---|---|---|---|---|
| DiffMoE-S-E16 | 32M | 139M | 12 | 384 | 6 | 16 | 1 |
| DiffMoE-B-E16 | 130M | 555M | 12 | 768 | 12 | 16 | 1 |
| DiffMoE-L-E8 | 458M | 1.176B | 24 | 1024 | 16 | 8 | 1 |
| DiffMoE-L-E16 | 458M | 1.982B | 24 | 1024 | 16 | 16 | 1 |

Table 2: **State-of-the-art Comparison**. Evaluation on ImageNet $256 \times 256$ class-conditional generation. DiffMoE achieves better FID with fewer parameters. -G/-U denotes with/without guidance (Ho & Salimans, 2022). [†]: results from (Ma et al., 2024) (DDPM) and (Peebles & Xie, 2023a) (Flow). [*]: our reproduction. **Bold** indicates best performance in each cell.

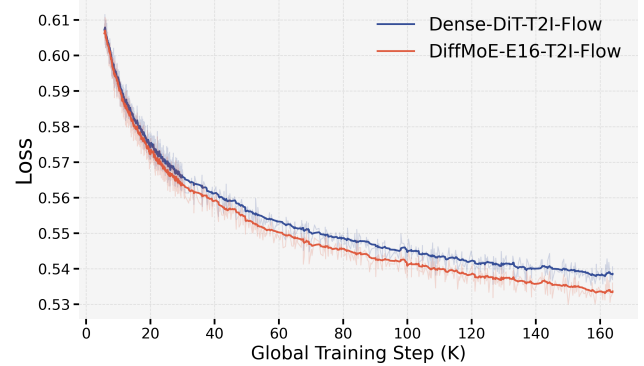| Diffusion Models (7000K) | # Avg. Activated Params. | FID↓ | IS↑ | Precision↑ | Recall↑ |
|---|---|---|---|---|---|
| Dense-DiT-XL-Flow-U[†] (Peebles & Xie, 2023a) | 675M | 9.35 | 126.06 | 0.67 | 0.68 |
| Dense-DiT-XL-Flow-U[*] (Peebles & Xie, 2023a) | 675M | **9.47** | 115.58 | 0.67 | 0.67 |
| DiffMoE-L-E8-Flow-U | 458M | 9.60 | 131.46 | 0.67 | 0.67 |
| Dense-DiT-XL-Flow-G[†] (cfg=1.5, ODE) (Ma et al., 2024) | 675M | 2.15 | 254.9 | 0.81 | 0.60 |
| Dense-DiT-XL-Flow-G[*] (cfg=1.5, ODE) (Ma et al., 2024) | 675M | 2.19 | 272.30 | 0.83 | 0.58 |
| DiffMoE-L-E8-Flow-G (cfg=1.5, ODE) | 458M | **2.13** | 274.39 | 0.81 | 0.60 |
| Dense-DiT-XL-DDPM-U[†] (Peebles & Xie, 2023a) | 675M | 9.62 | 121.50 | 0.67 | 0.67 |
| Dense-DiT-XL-DDPM-U[*] (Peebles & Xie, 2023a) | 675M | 9.62 | 123.19 | 0.66 | 0.68 |
| DiffMoE-L-E8-DDPM-U | 458M | **9.17** | 131.10 | 0.67 | 0.67 |
| Dense-DiT-XL-DDPM-G[†] (cfg=1.5) (Peebles & Xie, 2023a) | 675M | 2.27 | 278.2 | 0.83 | 0.57 |
| Dense-DiT-XL-DDPM-G[*] (cfg=1.5) (Peebles & Xie, 2023a) | 675M | 2.32 | 279.18 | 0.83 | 0.57 |
| DiffMoE-L-E8-DDPM-G (cfg=1.5) | 458M | **2.30** | 284.78 | 0.82 | 0.59 |

Table 3: **Baseline Model Comparisons.** DiffMoE-L-E16-Flow achieves best FID50K (14.41 w/o CFG) among TC, EC, and Dense variants. DDPM results in Appendix 11.

| Model (700K) | # Avg. Activated Params. | $C_{\text{infer}}^{\text{avg}}$ | FID50K ↓ |
|---|---|---|---|
| TC-DiT-L-E16-Flow | 458M | 1 | 19.06 |
| EC-DiT-L-E16-Flow | 458M | 1 | 16.12 |
| Dense-DiT-L-Flow | 458M | 1 | 17.01 |
| Dense-DiT-XL-Flow | 675M | 1 | <u>14.77</u> |
| DiffMoE-L-E16-Flow | 454M | 0.95 | **14.41** |

Table 4: **Ablation of Capacity Predictor.** We employ the capacity predictor to perform dynamic token selection, comparing it with the fixed TopK token selection method.

| Model (700K) | Capacity Predictor | $C_{\text{infer}}^{\text{avg}}$ | FID50K ↓ |
|---|---|---|---|
| DiffMoE-L-E16-Flow | w/o | 0.9 | 16.63 |
| DiffMoE-L-E16-Flow | w/o | 0.95 | 15.99 |
| DiffMoE-L-E16-Flow | w/o | 1 | <u>15.25</u> |
| DiffMoE-L-E16-Flow | w | 0.95 | **14.41** |



Figure 5: **Text-to-Image Generation Loss Curves.** Training loss comparison between DiffMoE-E16-T2I-Flow and Dense-DiT-T2I-Flow models over 160K steps. DiffMoE consistently achieves lower loss values, demonstrating superior convergence efficiency compared to the dense baseline.

**Comparison with Baseline.** DiffMoE-L-E16 demonstrates superior efficiency by outperforming Dense-DiT-XL (with $1.5\times$ parameters) after 700K steps, as shown in Table 3. It consistently achieves lower training loss compared to variants TC-DiT-L-E16, EC-DiT-L-E16 and Dense-DiT-L (Figure 4, 11, 10). These improvements hold across both DDPM and Flow Matching paradigms while maintaining equivalent activated parameters. With more training time

computation, DiffMoE-L-E16 can outperform Dense-DiT-XXXL (with $3\times$ parameters) as shown in Table 5.

**Comparison with SOTA.** After 7000K steps, DiffMoE-L-E8 achieves state-of-the-art FID50K scores of DDPM/Flow (2.30/2.13) with cfg=1.5, surpassing Dense-DiT-XL (2.32/2.19) as shown in Table 2. All evaluations follow DiT's (Peebles & Xie, 2023a) protocol. Generated C2I samples are shown in Figure 15 and 16.
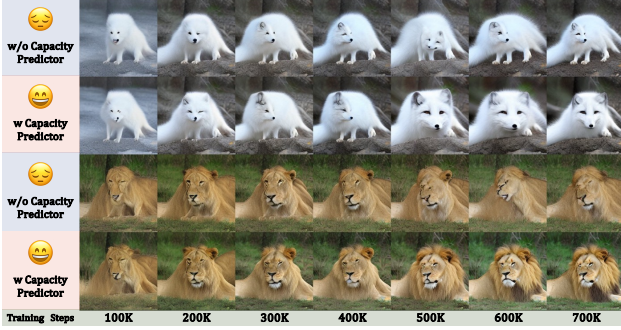
Figure 6: **Training and Inference Gap.** Comparison of sampling strategies with DiffMoE-L-E16-Flow (Batch Size = 1). For each group, **Top:** Sampling w/o Capacity Predictor (with Fixed TopK Method.) **Bottom:** Sampling with Capacity Predictor.



Figure 7: **Different Threshold Methods.** We employ two distinct approaches for threshold determination: dynamic threshold (red point) and interval search (blue points). Visualization using DiffMoE- L-E16-Flow (700K).

### 3.3. Main Results: Text-to-Image Generation

In Text-to-Image (T2I) generation, DiffMoE demonstrates superior performance over the Dense Model across multiple metrics. Without/With supervised fine-tuning (SFT), DiffMoE achieves GenEval scores of 0.44/0.51, outperforming the Dense Model's 0.38/0.49, with improvements across nearly all sub-metrics (Appendix Table 7). And as shown in Figure 5, DiffMoE maintains lower training losses while using the same activated parameters. Qualitative results in Figure 17 validate DiffMoE's ability to generate higher-quality images without SFT. Additional samples from SFT-enhanced models are shown in Figure 18.

### 3.4. Dynamic Computation Analysis

For the convenience of elaboration, we use the flow matching training method to do the following analysis while DDPM results are also provided in the Appendix D.1.

**Analysis of Inference Capacity.** DiffMoE-L-E16-Flow demonstrates superior parameter efficiency with its inference capacity ($C_{\text{infer}}^{\text{avg}}$) being 1 less than TC-DiT and EC-DiT, while achieving better performance, as shown in Table 3. Notably, with only 454M average activated parameters, our model outperforms Dense-DiT-XL-Flow (675M parameters), highlighting the effectiveness of dynamic expert allocation. Detailed analysis of average activated parameters is provided in the Appendix C.

**Ablation of Capacity Predictor.** Dynamic token selection through our capacity predictor demonstrates superior performance over traditional static topK token selection, as shown in Table 4. This improvement stems from the predictor's ability to intelligently allocate more computational resources to challenging tasks. The capacity predictor plays a crucial role in unleashing DiffMoE's full potential by dynamically adjusting resource allocation, which is particularly important for optimizing inference efficiency
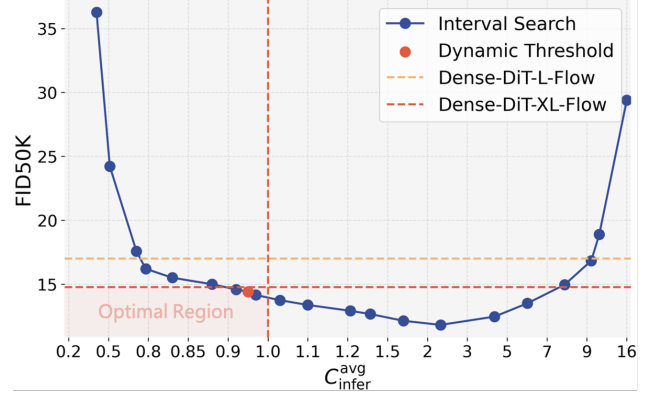
(Section 2.2). Without such adaptive mechanism, DiffMoE suffers from severe quality degradation due to sub-optimal resource utilization, as illustrated in Figure 6.

**Interval Search vs. Dynamic Threshold.** Both interval search and dynamic threshold methods achieve optimal performance in DiffMoE-L-E16-Flow, with the dynamic threshold ($\mathcal{T}^{Dy}$) emerging as our preferred approach due to its elegance and efficiency. Through interval search from 0.0 to 0.999, we identify an optimal threshold ($\gamma \approx 0.4$) that minimizes FID while maintaining $C_{\text{infer}}^{\text{avg}} \leq 1$. Meanwhile, the dynamic threshold automatically maintains $C_{\text{infer}}^{\text{avg}} \approx 1$ during inference, achieving comparable FID scores within the optimal region, as shown in Figure 7 and Table 10. Our experiments reveal a U-shaped relationship between FID and $C_{\text{infer}}^{\text{avg}}$, indicating that both over-activation and under-activation of parameters degrade performance. Both methods successfully identify thresholds within the optimal region, but the dynamic threshold's straightforward implementation and computational efficiency make it our default choice throughout this paper.

**Harder Work Needs More Computation.** Figure 1 demonstrates that different classes require varying computational resources during generation. By analyzing 1K class labels and ranking their $C_{\text{infer}}^{\text{avg}}$, we observe distinct patterns in computational demands. The most challenging cases typically involve objects with precise details, complex materials, structural accuracy, and specific viewing angles (e.g., technical instruments, detailed artifacts). In contrast, natural subjects like common animals (birds, dogs, cats) generally require less computation. Figures 13 and 14 display the top-10 most and least computationally intensive classes for both flow-based and DDPM models, respectively.

Table 5: **Parameter Scaling Behavior** of Diffusion Models on ImageNet 256×256 Class-Conditional Generation. This table evaluates the impact of parameter scaling on model performance. DiffMoE demonstrates superior FID scores with fewer parameters after 3000K training steps, highlighting its efficiency. Models with -G employ classifier-free guidance (Ho & Salimans, 2022). The base parameter configuration is fixed at 458 million (458M). **Bold** indicates the best performance in each metric, while underline denotes the second-best performance.

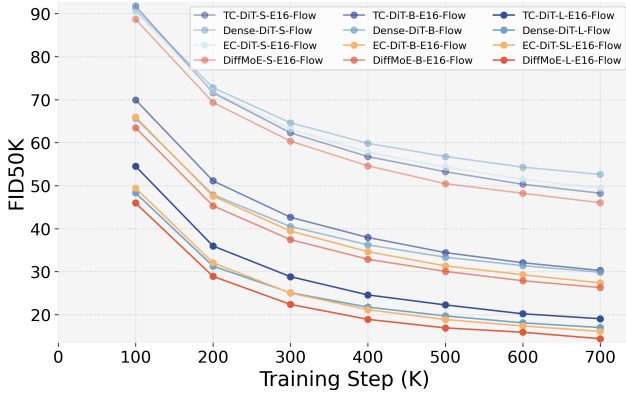| Diffusion Models (3000K) | # Avg. Activated Params. | FID↓ | IS↑ | Precision↑ | Recall↑ |
|---|---|---|---|---|---|
| Dense-DiT-XL-FlowG (cfg=1.5, ODE) | 675M (1.5x) | 2.52 | 273.78 | 0.84 | 0.56 |
| Dense-DiT-XXL-Flow-G (cfg=1.5, ODE) | 951M (2x) | 2.41 | 281.96 | 0.84 | 0.57 |
| Dense-DiT-XXXL-Flow-G (cfg=1.5, ODE) | 1353M (3x) | 2.37 | 291.29 | 0.84 | 0.57 |
| DiffMoE-L-E8-Flow-G (cfg=1.5, ODE) | 458M (1x) | 2.40 | 280.30 | 0.83 | 0.57 |
| DiffMoE-L-E16-Flow-G (cfg=1.5, ODE) | 458M (1x) | 2.36 | 287.26 | 0.83 | 0.58 |
| DiffMoE-XL-E16-Flow-G (cfg=1.5, ODE) | 675M (1.5x) | **2.30** | 291.23 | 0.83 | 0.58· |



Figure 8: **Scaling Model Size.** We analyze the impact of model size scaling by plotting FID50K scores across training steps. DiffMoE consistently outperforms the corresponding baseline models across all scales (S/B/L).
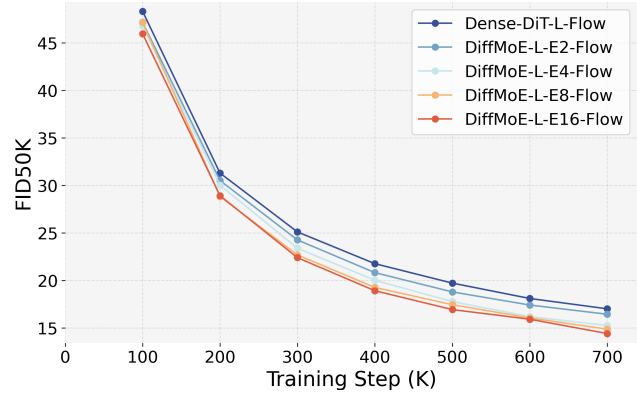


Figure 9: **Scaling Number of Experts.** Comparison of FID50K scores during training between Dense-DiT-L-Flow (E1) and models with increasing expert counts (E2, E4, E8, E16).

### 3.5. Scaling Behavior

**Scaling the Model Size.** DiffMoE demonstrates consistent performance improvements across small (S), base (B), and large (L) configurations, with activated parameters of 32M, 130M, and 458M respectively (Figure 8).

**Scaling Number of Experts.** As shown in Figure 9, model performance improves consistently when scaling experts from 2 to 16, with diminishing returns between E8 and E16. Based on this analysis, we trained DiffMoE-L-E8 for 7000K iterations, achieving optimal performance-efficiency trade-off and state-of-the-art results.

**Scaling Parameter Behavior.** To explore the upper limits of DiffMoE and quantify its performance efficiency, we scaled the model to larger sizes and trained them for 3000K steps. As illustrated in Table 5, DiffMoE-L-E16-Flow achieves the best performance among the evaluated models. Notably, DiffMoE-L-E16 surpasses the performance of Dense-DiT-XXXL-Flow, which uses 3x the parameters, while operating with only 1x the parameters. This highlights the exceptional parameter efficiency and scalability of DiffMoE.

## 4. Related Works

**Diffusion Models.** Diffusion models (Ho et al., 2020; Podell et al., 2023; Peebles & Xie, 2023a; Esser et al., 2024b) have emerged as the dominant paradigm in visual generation in recent years. These models transform gaussian distribution into target data distribution through iterative processes, with two primary training paradigms: Denoising Diffusion Probabilistic Models (DDPM) trained via score-matching (Ho et al., 2020; Song et al., 2021), which learns the inverse of a diffusion process and Rectified Flow approaches optimized through flow-matching (Lipman et al., 2022; Ma et al., 2024; Esser et al., 2024b), which is a more generic modeling techinique and can construct a straight probaility path connecting data and noise. We implement DiffMoE using both paradigms, demonstrating its versatility across these complementary training methodologies.

**Mixture of Experts.** Mixture of Experts (MoE) (Shazeer et al., 2017; Lepikhin et al., 2020) enables efficient model scaling through conditional computation by selectively activating expert subsets. This approach has demonstrated remarkable success in Large Language Models (LLMs), as

evidenced by cutting-edge implementations like DeepSeek-V3 (DeepSeek-AI et al., 2024), minimax-01 (MiniMax et al., 2025), and OLMOE (Muennighoff et al., 2024).Recent works have explored incorporating MoE architectures into diffusion models, but face several limitations. MEME (Lee et al., 2023), eDiff-I (Balaji et al., 2022), and ERNIE-ViLG 2.0 (Feng et al., 2023) restrict experts to specific timestep ranges. SegMoE (Yatharth Gupta, 2024) and DiT-MoE (Fei et al., 2024) suffer from expert utilization imbalance due to isolated token processing. While EC-DiT (Sehwag et al., 2024; Sun et al., 2024a) recognizes complex tokens' need for additional computation, it constrains token selection within individual samples and requires longer training for marginal improvements. These approaches, by limiting global token distribution across noise levels and conditions, fail to capture diffusion processes' inherent heterogeneity. DiffMoE addresses these challenges through batch-level global token pool for training, and dynamically adapting computation to both noise levels and sample complexity for inference.

## 5. Conclusion

In this work, we introduce DiffMoE, a simple yet powerful approach to scaling diffusion models efficiently through dynamic token selection and global token accessibility. Our method effectively addresses the uniform processing limitation in diffusion transformers by leveraging specialized expert behavior and dynamic resource allocation. Extensive experimental results demonstrate that DiffMoE substantially outperforms existing TC-MoE and EC-MoE methods, as well as dense models with 1.5× parameters, while maintaining comparable computational costs. The effectiveness of our approach not only validates its utility in class-conditional generation but also positions DiffMoE as a key enabler for advancing large-scale text-to-image and text-to-video generation tasks. While we excluded modern MoE enhancements for fair comparisons, integrating advanced techniques like Fine-Grained Expert (Yang et al., 2024) and Shared Expert (Dai et al., 2024) presents compelling opportunities for future work. These architectural improvements, combined with DiffMoE's demonstrated scalability and efficiency, could significantly advance the development of more powerful world simulators in the AI landscape.

## Impact Statement

DiffMoE represents a significant advancement in visual content generation, demonstrating exceptional scalability to billion-parameter architectures while maintaining computational efficiency. However, like other powerful AI models, it raises important ethical considerations that warrant careful attention. These include potential misuse for generating misleading content, privacy concerns regarding training data and generated content, environmental impact of large-scale model training, and broader societal implications of automated content generation. We are committed to addressing these challenges through robust safety measures, transparent guidelines, and continuous engagement with stakeholders to ensure responsible development and deployment of this technology.

## Acknowledgements

## References

Albergo, M. S. and Vanden-Eijnden, E. Building normalizing flows with stochastic interpolants, 2023. URL https://arxiv.org/abs/2209.15571.

Balaji, Y., Nah, S., Huang, X., Vahdat, A., Song, J., Zhang, Q., Kreis, K., Aittala, M., Aila, T., Laine, S., Catanzaro, B., Karras, T., and Liu, M.-Y. ediff-i: Text-to-image diffusion models with ensemble of expert denoisers. *arXiv preprint arXiv:2211.01324*, 2022.

Bao, F., Nie, S., Xue, K., Cao, Y., Li, C., Su, H., and Zhu, J. All are worth words: A vit backbone for diffusion models, 2023. URL https://arxiv.org/abs/2209.12152.

Cai, W., Jiang, J., Wang, F., Tang, J., Kim, S., and Huang, J. A survey on mixture of experts, 2024. URL https://arxiv.org/abs/2407.06204.

Chang, H., Zhang, H., Jiang, L., Liu, C., and Freeman, W. T. Maskgit: Masked generative image transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11315–11325, 2022.

Chen, R. T., Rubanova, Y., Bettencourt, J., and Duvenaud, D. K. Neural ordinary differential equations. *Advances in neural information processing systems*, 31, 2018.

Dai, D., Deng, C., Zhao, C., Xu, R. X., Gao, H., Chen, D., Li, J., Zeng, W., Yu, X., Wu, Y., Xie, Z., Li, Y. K., Huang, P., Luo, F., Ruan, C., Sui, Z., and Liang, W. Deepseekmoe: Towards ultimate expert specialization in mixture-of-experts language models, 2024. URL https://arxiv.org/abs/2401.06066.

DeepSeek-AI, Liu, A., Feng, B., Xue, B., Wang, B., Wu, B., Lu, C., Zhao, C., Deng, C., Zhang, C., Ruan, C., Dai, D., Guo, D., et al. Deepseek-v3 technical report, 2024. URL https://arxiv.org/abs/2412.19437.

Dhariwal, P. and Nichol, A. Diffusion models beat gans on image synthesis. *NeurIPS*, 34:8780–8794, 2021.

Esser, P., Kulal, S., Blattmann, A., Entezari, R., Müller, J., Saini, H., Levi, Y., Lorenz, D., Sauer, A., Boesel, F., et al. Scaling rectified flow transformers for high-resolution image synthesis. *arXiv preprint arXiv:2403.03206*, 2024a.

Esser, P., Kulal, S., Blattmann, A., Entezari, R., Müller, J., Saini, H., Levi, Y., Lorenz, D., Sauer, A., Boesel, F., Podell, D., Dockhorn, T., English, Z., Lacey, K., Goodwin, A., Marek, Y., and Rombach, R. Scaling rectified flow transformers for high-resolution image synthesis, 2024b. URL https://arxiv.org/abs/2403.03206.

Fei, Z., Fan, M., Yu, C., Li, D., and Huang, J. Scaling diffusion transformers to 16 billion parameters. *arXiv preprint*, 2024.

Feng, Z., Zhang, Z., Yu, X., Fang, Y., Li, L., Chen, X., Lu, Y., Liu, J., Yin, W., Feng, S., Sun, Y., Chen, L., Tian, H., Wu, H., and Wang, H. Ernie-vilg 2.0: Improving text-to-image diffusion model with knowledge-enhanced mixture-of-denoising-experts, 2023. URL https://arxiv.org/abs/2210.15257.

Ghosh, D., Hajishirzi, H., and Schmidt, L. Geneval: An object-focused framework for evaluating text-to-image alignment, 2023. URL https://arxiv.org/abs/2310.11513.

Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., and Hochreiter, S. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In Guyon, I., Luxburg, U. V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. URL https://proceedings.neurips.cc/paper_files/paper/2017/file/8a1d694707eb0fefe65871369074926d-Paper.pdf.

Ho, J. and Salimans, T. Classifier-free diffusion guidance, 2022. URL https://arxiv.org/abs/2207.12598.

Ho, J., Jain, A., and Abbeel, P. Denoising diffusion probabilistic models. *NeurIPS*, 33:6840–6851, 2020.

Kingma, D., Salimans, T., Poole, B., and Ho, J. Variational diffusion models. *NeurIPS*, 34:21696–21707, 2021.

Lee, Y., Kim, J.-Y., Go, H., Jeong, M., Oh, S., and Choi, S. Multi-architecture multi-expert diffusion models, 2023. URL https://arxiv.org/abs/2306.04990.

Lepikhin, D., Lee, H., Xu, Y., Chen, D., Firat, O., Huang, Y., Krikun, M., Shazeer, N., and Chen, Z. Gshard: Scaling giant models with conditional computation and automatic

sharding, 2020. URL https://arxiv.org/abs/2006.16668.

Lipman, Y., Chen, R. T., Ben-Hamu, H., Nickel, M., and Le, M. Flow matching for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022.

Lipman, Y., Chen, R. T. Q., Ben-Hamu, H., Nickel, M., and Le, M. Flow matching for generative modeling, 2023. URL https://arxiv.org/abs/2210.02747.

Liu, X., Gong, C., and Liu, Q. Flow straight and fast: Learning to generate and transfer data with rectified flow. *arXiv preprint arXiv:2209.03003*, 2022a.

Liu, X., Gong, C., and Liu, Q. Flow straight and fast: Learning to generate and transfer data with rectified flow, 2022b. URL https://arxiv.org/abs/2209.03003.

Liu, X., Zhang, X., Ma, J., Peng, J., et al. Instaflow: One step is enough for high-quality diffusion-based text-to-image generation. In *The Twelfth International Conference on Learning Representations*, 2023.

Loshchilov, I. and Hutter, F. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.

Ma, N., Goldstein, M., Albergo, M. S., Boffi, N. M., Vanden-Eijnden, E., and Xie, S. Sit: Exploring flow and diffusion-based generative models with scalable interpolant transformers. *arXiv preprint arXiv:2401.08740*, 2024.

MiniMax, Li, A., Gong, B., Yang, B., Shan, B., Liu, C., Zhu, C., Zhang, C., Guo, C., Chen, D., Li, D., Jiao, E., Li, G., Zhang, G., Sun, H., Dong, H., Zhu, J., Zhuang, J., Song, J., Zhu, J., Han, J., Li, J., Xie, J., Xu, J., Yan, J., Zhang, K., Xiao, K., Kang, K., et al. Minimax-01: Scaling foundation models with lightning attention, 2025. URL https://arxiv.org/abs/2501.08313.

Muennighoff, N., Soldaini, L., Groeneveld, D., Lo, K., Morrison, J., Min, S., Shi, W., Walsh, P., Tafjord, O., Lambert, N., Gu, Y., Arora, S., Bhagia, A., Schwenk, D., Wadden, D., Wettig, A., Hui, B., Dettmers, T., Kiela, D., Farhadi, A., Smith, N. A., Koh, P. W., Singh, A., and Hajishirzi, H. Olmoe: Open mixture-of-experts language models, 2024. URL https://arxiv.org/abs/2409.02060.

Pan, J., Sun, K., Ge, Y., Li, H., Duan, H., Wu, X., Zhang, R., Zhou, A., Qin, Z., Wang, Y., Dai, J., Qiao, Y., and Li, H. Journeydb: A benchmark for generative image understanding, 2023.

Peebles, W. and Xie, S. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 4195–4205, 2023a.

Peebles, W. and Xie, S. Scalable diffusion models with transformers, 2023b. URL https://arxiv.org/abs/2212.09748.

Podell, D., English, Z., Lacey, K., Blattmann, A., Dockhorn, T., Müller, J., Penna, J., and Rombach, R. Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952*, 2023.

Qiu, Z., Huang, Z., Zheng, B., Wen, K., Wang, Z., Men, R., Titov, I., Liu, D., Zhou, J., and Lin, J. Demons in the detail: On implementing load balancing loss for training specialized mixture-of-expert models. *arXiv preprint arXiv:2501.11873*, 2025.

Rombach, R., Blattmann, A., Lorenz, D., Esser, P., and Ommer, B. High-resolution image synthesis with latent diffusion models. In *CVPR*, pp. 10684–10695, 2022.

Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A. C., and Fei-Fei, L. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015. doi: 10.1007/s11263-015-0816-y.

Sauer, A., Schwarz, K., and Geiger, A. Stylegan-xl: Scaling stylegan to large diverse datasets. In *ACM SIGGRAPH 2022 conference proceedings*, pp. 1–10, 2022.

Sehwag, V., Kong, X., Li, J., Spranger, M., and Lyu, L. Stretching each dollar: Diffusion training from scratch on a micro-budget, 2024. URL https://arxiv.org/abs/2407.15811.

Shazeer, N., Mirhoseini, A., Maziarz, K., Davis, A., Le, Q., Hinton, G., and Dean, J. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer, 2017. URL https://arxiv.org/abs/1701.06538.

Sohl-Dickstein, J., Weiss, E., Maheswaranathan, N., and Ganguli, S. Deep unsupervised learning using nonequilibrium thermodynamics. In *ICML*, pp. 2256–2265. PMLR, 2015.

Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Ermon, S., and Poole, B. Score-based generative modeling through stochastic differential equations. In *ICLR*, 2021.

Sun, H., Lei, T., Zhang, B., Li, Y., Huang, H., Pang, R., Dai, B., and Du, N. Ec-dit: Scaling diffusion transformers with adaptive expert-choice routing, 2024a. URL https://arxiv.org/abs/2410.02098.

Sun, P., Jiang, Y., Chen, S., Zhang, S., Peng, B., Luo, P., and Yuan, Z. Autoregressive model beats diffusion: Llama for scalable image generation, 2024b. URL https://arxiv.org/abs/2406.06525.

Tian, K., Jiang, Y., Yuan, Z., Peng, B., and Wang, L. Visual autoregressive modeling: Scalable image generation via next-scale prediction. 2024.

Yang, Y., Qi, S., Gu, W., Wang, C., Gao, C., and Xu, Z. Xmoe: Sparse models with fine-grained and adaptive expert selection, 2024. URL https://arxiv.org/abs/2403.18926.

Yatharth Gupta, Vishnu V Jaddipal, H. P. Segmoe: Segmind mixture of diffusion experts. https://github.com/segmind/segmoe, 2024.

Zhou, Y., Lei, T., Liu, H., Du, N., Huang, Y., Zhao, V., Dai, A., Chen, Z., Le, Q., and Laudon, J. Mixture-of-experts with expert choice routing, 2022. URL https://arxiv.org/abs/2202.09368.

# A. Generative Modeling.

In this section, we will provide a detailed bachground of generative modeling of both DDPM (Ho et al., 2020; Song et al., 2021; Rombach et al., 2022) and Rectified Flow (Lipman et al., 2023; Ma et al., 2024; Esser et al., 2024b) which is helpful to understand the difference and relationship between them.

Generative modeling essentially defines a mapping between $\mathbf{x}_1$ from a noise distribution $p_1(\mathbf{x})$ to $\mathbf{x}_0$ from a data distribution $p_0(\mathbf{x})$ leads to time-dependent processes represented as below

$$\mathbf{x}_t = \alpha_t \mathbf{x}_0 + \sigma_t \boldsymbol{\epsilon}, t \in [0, 1], \tag{18}$$

where $\alpha_t$ is a decreasing function of $t$ and $\sigma_t$ is an increasing function of $t$. We set $\alpha_0 = 1, \sigma_0 = 0$ and $\alpha_1 = 0, \sigma_0 = 1$ to make the marginals $p_t(\mathbf{x}_t) = \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(0, \mathrm{I})} p_t(\mathbf{x}_t | \boldsymbol{\epsilon})$ are consistent with data $p_0(\mathbf{x})$ and noise $p_1(\mathbf{x})$ distirbution. $p_1(\mathbf{x})$ usually be chosen as gaussion distribution $\mathcal{N}(0, 1)$.

Different forward path from data to noise leads to different training object which significantly affect the performance of the model. Next we will introduce DDPM and Rectified Flow.

## A.1. Denosing Diffusion Probabilistic Models (DDPM)

In DDPM the choice for $\alpha_t$ and $\sigma_t$ is referred to as the noise schedule and the signal-to-noise-ratio (SNR) $\alpha_t^2 / \sigma_t^2$ is strictly decreasing w.r.t $t$ (Kingma et al., 2021). Moreover, (Kingma et al., 2021) prove that the following stochastic differential Eq. (SDE) has same transition distribution as $p_t(\mathbf{x}_t | \mathbf{x}_0)$ for any $t \in [0, 1]$:

$$\mathrm{d}\mathbf{x}_t = f(t)\mathbf{x}_t \mathrm{d}t + g(t)\mathrm{d}\mathbf{w}_t, t \in [0, 1], \mathbf{x}_0 \sim p_0(\mathbf{x}_0), \tag{19}$$

where $\mathbf{w}_t \in \mathbb{E}^D$ is the standard Wiener process, and

$$f(t) = \frac{\mathrm{d} \log \alpha_t}{\mathrm{d}t}, \quad g^2(t) = \frac{\mathrm{d}\sigma_t^2}{\mathrm{d}t} - 2\frac{\mathrm{d} \log \alpha_t}{\mathrm{d}t}\sigma_t^2. \tag{20}$$

(Song et al., 2021) proved that the forward path in Eq. 19 has an equivalent reverse process from time 1 to 0 under some regularity conditions, starting with $p_T(\mathbf{x}_T)$

$$\mathrm{d}\mathbf{x}_t = [f(t)\mathbf{x}_t - g(t)^2 \nabla_{\mathbf{x}} \log p_t(\mathbf{x}_t)]\mathrm{d}t + g(t)\mathrm{d}\bar{\mathbf{w}}_t, \quad \mathbf{x}_T \sim p_T(\mathbf{x}_T), \tag{21}$$

where $\bar{\mathbf{w}}_t \in \mathbb{E}^D$ is the standard Wiener process. We can esitimate the score term $\nabla_{\mathbf{x}} \log p_t(\mathbf{x}_t)$ as each time $t$ to iterativtely solve the reverse process, then get the gernerated target. DPMs train a neural network $\boldsymbol{\epsilon}_\theta(\mathbf{x}, t)$ parameterized by $\theta$ to esitimated the scaled score function $-\sigma \nabla_{\mathbf{x}} \log p_t(\mathbf{x}_t)$. To optimize $\boldsymbol{\epsilon}_\theta$, we minimize the following objective (Ho et al., 2020; Song et al., 2021; Ma et al., 2024)

$$\mathcal{L}_{\mathrm{DDPM}}(\theta) = \mathbb{E}_{t, p_0(\mathbf{x}_0), p(\mathbf{x}_t | \mathbf{x}_0)} \left[ \lambda(t) \| \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) + \sigma_t \nabla_{\mathbf{x}} \log p_t(\mathbf{x}_t) \|_2^2 \right], \tag{22}$$

where $\lambda_t$ is a time-dependent coeffect. $\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)$ can be interpreted as predicting the Gaussian noise added to $\mathbf{x}_t$, thus it is commonly referred to as a noise prediction model. Consequently, the diffusion model is known as a denoising diffusion probabilistic model. Substitute score term in (21) with $-\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)/\sigma_t$, we can solve the reverse process and generate samples from DPMs with numerical solvers. To further accelerate the sampling process, Song *et al.* (Song et al., 2021) proved that the equvivalent probability flow ODE is

$$\frac{\mathrm{d}\mathbf{x}_t}{\mathrm{d}t} = \mathbf{v}_\theta(\mathbf{x}_t, t) := f(t)\mathbf{x}_t + \frac{g^2(t)}{2\sigma_t} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t), \quad \mathbf{x}_1 \sim \mathcal{N}(0, \mathbf{I}). \tag{23}$$

Thus samples can be also generated by solving the ODE from time 1 to 0.

## A.2. Rectified Flow Models (Flow)

Recified flow models (Liu et al., 2022b; Albergo & Vanden-Eijnden, 2023; Lipman et al., 2023) connects data $\mathbf{x}_0$ and noise $\boldsymbol{\epsilon}$ on a straight line as follows

$$\mathbf{x}_t = (1 - t)\mathbf{x}_0 + t\boldsymbol{\epsilon}, \quad t \in [0, 1]. \tag{24}$$

To precisely express the relationship between $\mathbf{x}_t$, $\mathbf{x}_0$, and $\boldsymbol{\epsilon}$, we first construct a time-dependent vector field $u : [0, 1] \times \mathbb{R}^D \to \mathbb{R}^D$. This vector field $u_t$ can be used to construct a time-dependent diffeomorphic map, known as a flow $\phi : [0, 1] \times \mathbb{R}^D \to \mathbb{R}^D$, through the following ODE:

$$\frac{\mathrm{d}}{\mathrm{d}t}\phi_t(\mathbf{x}_0) = u_t(\phi_t(\mathbf{x}_0)) \tag{25}$$

$$\phi_0(\mathbf{x}_0) = \mathbf{x}_0. \tag{26}$$

The vector field $u_t$ can be modeled as a neural network $\mathbf{v}_\theta$ (Chen et al., 2018) which leads to a deep parametric model of the flow $\phi_t$, called a *Continuous Normalizing Flow* (CNF). We can using *conditional flow matching* (CFM) technique (Lipman et al., 2022) to training a CNF. Now we can define the flow we need as follows:

$$\psi_t(\cdot|\boldsymbol{\epsilon}) : \mathbf{x}_0 \mapsto \alpha_t \mathbf{x}_0 + \sigma_t \boldsymbol{\epsilon}. \tag{27}$$

The corresponding velocity vector field of the flow $\psi_t$ can be represented as:

$$u_t(\psi_t(\cdot|\boldsymbol{\epsilon})|\boldsymbol{\epsilon}) = \frac{\mathrm{d}}{\mathrm{d}t}\psi_t(\mathbf{x}_0|\boldsymbol{\epsilon}) = \dot{\alpha}_t \mathbf{x}_0 + \dot{\sigma}_t \boldsymbol{\epsilon} = \boldsymbol{\epsilon} - \mathbf{x}_0. \tag{28}$$

Using conditional flow matching technique, $\mathbf{v}(\mathbf{x}_t, t)$ in Eq. (23) can be modeled as a neural network $\mathbf{v}_\theta(\mathbf{x}_t, t)$ by minimziing the following objective

$$\mathcal{L}_{\text{Flow}}(\theta) = \mathbb{E}_{t, p_0(\mathbf{x}_0), p_1(\boldsymbol{\epsilon})} \|\mathbf{v}_\theta(\mathbf{x}_t, t) - \frac{\mathrm{d}}{\mathrm{d}t}\psi_t(\mathbf{x}_0|\boldsymbol{\epsilon})\|_2^2 \tag{29}$$

$$= \mathbb{E}_{t, p_0(\mathbf{x}_0), p(\boldsymbol{\epsilon})} \|\mathbf{v}_\theta(\mathbf{x}_t, t) - (\dot{\alpha}_t \mathbf{x}_0 + \dot{\sigma}_t \boldsymbol{\epsilon})\|_2^2 \tag{30}$$

$$= \mathbb{E}_{t, p_0(\mathbf{x}_0), p(\boldsymbol{\epsilon})} \|\mathbf{v}_\theta(\mathbf{x}_t, t) - (\boldsymbol{\epsilon} - \mathbf{x}_0)\|_2^2. \tag{31}$$

Samples can be generated by sovling the probability flow ODE below with learned velocity using numerical sovler like Euler, Heun, Runge-Kutta method.

$$\mathrm{d}\mathbf{x}_t = \mathbf{v}_\theta(\mathbf{x}_t, t)\mathrm{d}t, \quad \mathbf{x}_1 = \boldsymbol{\epsilon} \sim \mathcal{N}(0, 1). \tag{32}$$

## A.3. Relationship DDPM and Flow

There exists a straightforward connection between $\mathbf{v}_\theta(\mathbf{x}_t, t)$ and the score term $-\sigma_t \nabla_{\mathbf{x}} \log p_t(\mathbf{x}_t)$ can be derived as follows

$$\mathbf{v}_\theta(\mathbf{x}_t, t) = f(t)\mathbf{x}_t + \frac{g^2(t)}{2\sigma_t}\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \tag{33}$$

$$\approx \frac{\dot{\alpha}_t}{\alpha_t}\mathbf{x}_t + \left(\dot{\sigma}_t - \frac{\dot{\alpha}_t}{\alpha_t}\right)(-\sigma_t \nabla_{\mathbf{x}} \log p_t(\mathbf{x}_t)). \tag{34}$$

Let $\zeta_t = \dot{\sigma}_t - \frac{\dot{\alpha}_t}{\alpha_t}$, and we have $\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \approx -\sigma_t \nabla_{\mathbf{x}} \log p_t(\mathbf{x}_t$, then we can get $\mathbf{v}_\theta(\mathbf{x}_t, t) = \frac{\dot{\alpha}_t}{\alpha_t}\mathbf{x}_t + \zeta_t \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)$

By plugging (34) into the loss $\mathcal{L}_{\text{Flow}}$ in Eq. (31) we have:

Table 6: **Comprehensive State-of-the-art Comparison**. Evaluation on ImageNet $256 \times 256$ class-conditional generation. DiffMoE achieves better FID with fewer parameters. -G/-U denotes with/without guidance (Ho & Salimans, 2022). $^\dagger$: results from (Ma et al., 2024) (DDPM) and (Peebles & Xie, 2023a) (Flow). $^*$: our reproduction. **Bold** indicates best performance in each cell.

| Diffusion Models | # Avg. Activated Params. | FID↓ | IS↑ | Precision↑ | Recall↑ |
|---|---|---|---|---|---|
| *GAN* | | | | | |
| StyleGAN-XL (Sauer et al., 2022) | 166M | 2.30 | 265.1 | 0.78 | 0.53 |
| *Masked Modeling* | | | | | |
| Mask-GIT (Chang et al., 2022) | 227M | 6.18 | 182.1 | 0.80 | 0.51 |
| *Autoregressive Model* | | | | | |
| LlamaGen-3B (Sun et al., 2024b) | 3000M | 10.94 | 101.0 | 0.69 | 0.63 |
| VAR-$d$20 (Tian et al., 2024) | 600M | 2.57 | 302.6 | 0.83 | 0.56 |
| *Diffusion Model* | | | | | |
| ADM-U (Dhariwal & Nichol, 2021) | 554M | 10.94 | 101.0 | 0.69 | 0.63 |
| ADM-G (Dhariwal & Nichol, 2021) | 608M | 4.59 | 186.7 | 0.83 | 0.53 |
| LDM-4-G (Rombach et al., 2022) | 400M | 3.95 | 247.7 | 0.87 | 0.48 |
| U-ViT-H/2-G (Bao et al., 2023) | 501M | 2.29 | 263.88 | 0.82 | 0.57 |
| Dense-DiT-XL-Flow-U$^\dagger$ (Peebles & Xie, 2023a) | 675M | 9.35 | 126.06 | 0.67 | 0.68 |
| Dense-DiT-XL-Flow-G$^\dagger$ (cfg=1.5, ODE) (Ma et al., 2024) | 675M | 2.15 | 254.9 | 0.81 | 0.60 |
| Dense-DiT-XL-Flow-U$^*$ (Peebles & Xie, 2023a) | 675M | **9.47** | 115.58 | 0.67 | 0.67 |
| Dense-DiT-XL-Flow-G$^*$ (cfg=1.5, ODE) (Ma et al., 2024) | 675M | 2.19 | 272.30 | 0.83 | 0.58 |
| DiffMoE-L-E8-Flow-U | 458M | **9.60** | 131.46 | 0.67 | 0.67 |
| DiffMoE-L-E8-Flow-G (cfg=1.5, ODE) | 458M | **2.13** | 274.39 | 0.81 | 0.60 |
| Dense-DiT-XL-DDPM-U$^\dagger$ (Peebles & Xie, 2023a) | 675M | 9.62 | 121.50 | 0.67 | 0.67 |
| Dense-DiT-XL-DDPM-G$^\dagger$ (cfg=1.5) (Peebles & Xie, 2023a) | 675M | 2.27 | 278.2 | 0.83 | 0.57 |
| Dense-DiT-XL-DDPM-U$^*$ (Peebles & Xie, 2023a) | 675M | 9.62 | 123.19 | 0.66 | 0.68 |
| Dense-DiT-XL-DDPM-G$^*$ (cfg=1.5) (Peebles & Xie, 2023a) | 675M | 2.32 | 279.18 | 0.83 | 0.57 |
| DiffMoE-L-E8-DDPM-U | 458M | **9.17** | 131.10 | 0.67 | 0.67 |
| DiffMoE-L-E8-DDPM-G (cfg=1.5) | 458M | **2.30** | 284.78 | 0.82 | 0.59 |

$$\mathcal{L}_{\text{Flow}}(\theta) = \mathbb{E}_{t,p_0(\mathbf{x}_0),p_1(\boldsymbol{\epsilon})} \|\mathbf{v}_\theta(\mathbf{x}_t, t) - (\dot{\alpha}_t \mathbf{x}_0 + \dot{\sigma}_t \boldsymbol{\epsilon})\|_2^2 \tag{35}$$

$$= \mathbb{E}_{t,p_0(\mathbf{x}_0),p(\boldsymbol{\epsilon})} \|\frac{\dot{\alpha}_t}{\alpha_t}\mathbf{x}_t + \zeta_t \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) - \dot{\sigma}_t \boldsymbol{\epsilon}\|_2^2 \tag{36}$$

$$= \mathbb{E}_{t,p_0(\mathbf{x}_0),p(\boldsymbol{\epsilon})} \|\zeta_t \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) - \zeta_t \boldsymbol{\epsilon}\|_2^2 \tag{37}$$

$$= \mathbb{E}_{t,p_0(\mathbf{x}_0),p(\boldsymbol{\epsilon})} \left[ \zeta_t^2 \|\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) - \boldsymbol{\epsilon}\|_2^2 \right]. \tag{38}$$

Considering Eq. (18), we have $\mathbf{x}_t \sim \mathcal{N}(\alpha_t \mathbf{x}_0, \sigma_t \mathbf{I})$ and $\nabla_{\mathbf{x}} \log p(\mathbf{x}_t) = \sigma_t^{-1}(\mathbf{x}_t - \alpha_t \mathbf{x}_0) = \sigma_t^{-1}(\sigma_t \boldsymbol{\epsilon}) = \boldsymbol{\epsilon}$. Then, we can get the equivilant loss function as below:

$$\mathcal{L}_{\text{Flow}}(\theta) = \mathbb{E}_{t,p_0(x_0),p(\mathbf{x}_t|\mathbf{x}_0)} \left[ \zeta_t^2 \|\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) + \sigma_t \nabla_{\mathbf{x}} \log p_t(\mathbf{x}_t)\|_2^2 \right]. \tag{39}$$

Recall that $\mathcal{L}_{\text{DDPM}}(\theta) = \mathbb{E}_{t,p_0(\mathbf{x}_0),p(\mathbf{x}_t|\mathbf{x}_0)} \left[ \lambda(t) \|\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) + \sigma_t \nabla_{\mathbf{x}} \log p_t(\mathbf{x}_t)\|_2^2 \right]$.

We can find that $\mathcal{L}_{\text{DDPM}}$ and $\mathcal{L}_{\text{Flow}}$ have the same form, with the only difference being their time-dependent weighting functions, which lead to different trajectories and properties.

## B. More Implementation Details

### B.1. Training setup.

We train class-conditional DiffMoE and baseline models at 256x256 image resolution on the ImageNet dataset (Russakovsky et al., 2015) a highly-competitive generative benchmark, which contains 1281167 training images. We use horizontal

Table 7: **Performance Comparison on Text-to-Image Generation Tasks**. Evaluation results on GenEval Benchmark (Ghosh et al., 2023) at $256 \times 256$ resolution across six different categories: single object generation, two-object composition, object counting, color recognition, spatial positioning, and color attribute understanding. DiffMoE-E16-T2I-Flow demonstrates superior performance over Dense-DiT-T2I-Flow, particularly in object generation and spatial tasks. **Bold** indicates best overall performance in each cell. # A.A.P. denotes # Avg. Act. Params. # T.P. denotes # Total Params.

| Diffusion Models | # A.A.P. | # T.P. | Single Obj. ↑ | Two Obj.↑ | Counting Obj.↑ | Colors ↑ | Position ↑ | Color Attri. ↑ | Overall ↑ |
|---|---|---|---|---|---|---|---|---|---|
| Dense-DiT-T2I-Flow (w/o SFT) | 1.2B | 1.2B | 0.80 | 0.33 | 0.24 | 0.53 | 0.10 | 0.26 | 0.38 |
| DiffMoE-E16-T2I-Flow (w/o SFT) | 1.2B | 4.6B | 0.90 | 0.42 | 0.24 | 0.59 | 0.19 | 0.28 | **0.44** |
| Dense-DiT-T2I-Flow (w SFT) | 1.2B | 1.2B | 0.93 | 0.48 | 0.50 | 0.73 | 0.07 | 0.20 | 0.49 |
| DiffMoE-E16-T2I-Flow (w SFT) | 1.2B | 4.6B | 0.96 | 0.53 | 0.46 | 0.78 | 0.13 | 0.20 | **0.51** |

Table 8: **Performance comparison of different diffusion models with varying token interaction strategies**. All models are trained with Flow Matching for 700K steps. The interaction levels (L1/L2/L3) represent: L1 for isolated token processing, L2 for local token routing within samples, and L3 for global token routing across samples. Our DiffMoE-L-Flow with Dynamic Global CP achieves the best FID score of 14.41 while maintaining parameter efficiency and reduced computational cost. # A.A.P. denotes # Avg. Act. Params.

| Model | # A.A.P. | Training Strategy | Inference Strategy | FID50K ↓ |
|---|---|---|---|---|
| TC-DiT-L-E16-Flow | 458M | L1: Isolated | L1: Isolated | 19.06 |
| Dense-DiT-L-Flow | 458M | L1: Isolated | L1: Isolated | 17.01 |
| EC-DiT-L-E16-Flow | 458M | L2: Local | L2: Local Static TopK Routing | 16.12 |
| EC-DiT-L-E16-Flow | 458M | L2: Local | L2: Local Dynamic Intra-sample Routing | 23.74 |
| DiffMoE-L-E16-Flow | 458M | L3: Global | L3: Global Static TopK Routing | 15.25 |
| Dense-DiT-XL-Flow | 675M | L1: Isolated | L1: Isolated | <u>14.77</u> |
| DiffMoE-L-E16-Flow | 454M | L3: Global | L3: Global Dynamic Cross-sample Routing | **14.41** |

flips as the only data augmentation. We train all models with AdamW (Loshchilov & Hutter, 2017). We use a constant learning rate of $1 \times 10^{-4}$, no weight decay and a fixed global batch size of 256, following (Peebles & Xie, 2023b). We also maintain an exponential moving average(EMA) of DiffMoE and baseline model weights over training with a decay of 0.9999. All results reported use the EMA mode. For most experiments, we utilized 4 NVIDIA H800 GPUs during training. To achieve state-of-the-art results, we extended the training process with 8 NVIDIA H800 GPUs for improved efficiency. For the text-to-image model pre-training, we employ 32 NVIDIA H800 GPUs with our internal datasets, and then conduct supervised fine-tuning (SFT) on the JourneyDB dataset (Pan et al., 2023).

### B.2. Implementation Algorithms.

We provide a detailed illustration of the DiffMoE layer during training and inference in Algorithm 3 and 4, respectively. We also implemented the same EC-DiT layer in Algorithm 1 as (Sun et al., 2024a). We implemented TC-DiT in layer in Algorithm 2 similar to (Dai et al., 2024; Fei et al., 2024)

## C. Calculation of Average Activated Parameters and Average Capacity $C_{\text{infer}}^{\text{avg}}$

### C.1. Computing $C_{\text{infer}}^{\text{avg}}$

To compute the global average capacity ($C_{\text{infer}}^{\text{avg}}$), we analyze 50K samples across all experts and sampling steps. For a quick approximation, sampling 1K samples is sufficient due to DiffMoE's stable performance characteristics.

### C.2. Estimating Average Activated Parameters using $C_{\text{infer}}^{\text{avg}}$

We will introduce the relationship between average activated parameters and average capacity in detail. Let #Module denote the number of parameters a certain Module, $N_E$ denote the number of experts. Under approximate conditions, for DiT (Peebles & Xie, 2023a) models, we have

Table 9: **Module Parameters and Percentage.** We have counted the number of parameters (M) of various modules to facilitate our analysis.

| Model | FFN | Attention | AdaLN | Others | Total |
|---|---|---|---|---|---|
| Dense-DiT-L | 201.44(44.0%) | 100.7(22.0%) | 151(33.0%) | 4.7(1.0%) | 457.84 |
| DiffMoE-L-E2 | 314.8(55.1%) | 100.7(17.6%) | 151(26.4%) | 4.7(0.8%) | 571.2 |
| DiffMoE-L-E4 | 516.3(66.8%) | 100.7(13.0%) | 151(19.5%) | 4.7((0.6%)) | 772.7 |
| DiffMoE-L-E8 | 919.3(78.2%) | 100.7(8.6%) | 151(12.8%) | 4.7(0.4%) | 1175.7 |
| DiffMoE-L-E16 | 1725.3(87.1%) | 100.7(5.1%) | 151(7.6%)) | 4.7(0.2%) | 1981.7 |

Table 10: **Different Threshold Method.** We use both interval search and dynamic threshold method to find out the optimal $\tau_{E_i}$. We find that the dynamic threshold makes a good balance between $C_{\text{infer}}^{\text{avg}}$ and performance.

| $\mathcal{T}$ | $C_{\text{infer}}^{\text{avg}}$ | FID50K ↓ |
|---|---|---|
| 0.999 | 0.41 | 36.28 |
| 0.99 | 0.51 | 24.22 |
| 0.9 | 0.71 | 17.59 |
| 0.8 | 0.78 | 16.21 |
| 0.7 | 0.83 | 15.51 |
| 0.6 | 0.88 | 15.00 |
| 0.5 | 0.92 | 14.59 |
| 0.4 | 0.97 | **14.16** |
| 0.3 | 1.03 | 13.75 |
| 0.2 | 1.10 | 13.38 |
| 0.1 | 1.22 | 12.92 |
| Dynamic | 0.95 | 14.41 |

| $\mathcal{T}$ | $C_{\text{infer}}^{\text{avg}}$ | FID50K ↓ |
|---|---|---|
| 0.2 | 1.10 | 13.38 |
| 0.1 | 1.22 | 12.92 |
| 1E-2 | 1.71 | 12.14 |
| 1E-3 | 2.33 | **11.82** |
| 1E-4 | 3.17 | 11.87 |
| 1E-5 | 4.36 | 12.47 |
| 1E-6 | 6.01 | 13.51 |
| 1E-7 | 7.88 | 13.96 |
| 1E-8 | 9.67 | 16.85 |
| 1E-9 | 11.18 | 18.90 |
| 0.0 | 16 | 29.39 |

Table 11: **Comparisons with the Baseline Models. (DDPM)** We compare TC, EC and Dense Model and show the average activated parameters of all the experts across all the sampling steps.

| Model (700K) | # Avg. Activated Params. | $C_{\text{infer}}^{\text{avg}}$ | FID50K ↓ |
|---|---|---|---|
| TC-DiT-L-E16-DDPM | 458M | 1 | 20.81 |
| EC-DiT-L-E16-DDPM | 458M | 1 | 17.65 |
| Dense-DiT-L-DDPM | 458M | 1 | 17.87 |
| Dense-DiT-XL-DDPM | 675M | 1 | 15.28 |
| DiffMoE-L-E16-DDPM | 458M | 1 | **14.60** |

Table 12: **Decoder Ablation Study**. Evaluation of various pre-trained VAE decoder weights. †: results from (Ma et al., 2024) (DDPM) and (Peebles & Xie, 2023a) (Flow). *: our reproduction. All other results are from our experiments. In general, with VAE decoder EMA version, the FID score is consistently lower than MSE version.

| Model | Training Steps | VAE-Decoder | Sampler | Batch Size | FID50K ↓ |
|---|---|---|---|---|---|
| Dense-DiT-XL-Flow | 400K | ft-MSE | Euler | 125 | 18.80 |
| Dense-DiT-XL-Flow | 400K | ft-EMA | Euler | 125 | 18.74 |
| Dense-DiT-XL-Flow | 400K | ft-MSE | Dopri5 | 125 | 18.63 |
| Dense-DiT-XL-Flow | 400K | ft-EMA | Dopri5 | 125 | 18.45 |
| Dense-DiT-XL-Flow* | 7000K | ft-MSE | Heun | 125 | 9.66 |
| Dense-DiT-XL-Flow* | 7000K | ft-EMA | Heun | 125 | 9.63 |
| Dense-DiT-XL-Flow* | 7000K | ft-MSE | Dopri5 | 125 | 9.51 |
| Dense-DiT-XL-Flow* | 7000K | ft-EMA | Dopri5 | 125 | 9.48 |
| Dense-DiT-XL-DDPM-G † | 7000K | ft-MSE | DDPM | 125 | 2.30 |
| Dense-DiT-XL-DDPM-G † | 7000K | ft-EMA | DDPM | 125 | 2.27 |

$$\text{\# Average Activated Parameters} \approx \left(\frac{1 + C_{\text{infer}}^{\text{avg}}}{1 + N_E}\right) \text{\# FFN} + \text{\# Attention} + \text{\# AdaLN} + \text{\# Other Modules}. \quad (40)$$

.

Table 9 displays the parameters and their corresponding percentages of the main modules of large-size DiffMoE.

## D. More DiffMoE Analysis

### D.1. Additional DiffMoE DDPM Class-conditional Generation Qualitative and Quantitative Results

We present comprehensive evaluations of DiffMoE-L-E16-DDPM series models. Table 11 shows the experimental results, while Figure 10 illustrates the diffusion loss comparison against the baseline model, revealing substantial performance improvements. Furthermore, Figure 11 demonstrates the scaling capabilities of our DiffMoE-DDPM architecture.

### D.2. Additional DiffMoE Text-to-Image Qualitative and Quantitative Results

To further validate the capability of DiffMoE of more challenging task, we conducted quantitative and qualitative evaluations of Dense-DiT-T2I-Flow and DiffMoE-E16-T2I-Flow. The experimental results are presented in Table 7, demonstrating the comparative performance metrics of both models. Figure 5 provides a detailed visualization of the diffusion loss comparison between the dense model and DiffMoE, highlighting significant performance improvements achieved by our approach. The superior text-to-image generation capabilities of DiffMoE are further illustrated through qualitative examples in Figure 17 and Figure 15. It is worth noting that both text-to-image models were only trained for about 2 days, which is limited for this task, suggesting significant room for further performance improvements.

### D.3. Analysis of Token Interaction Strategies.

As shown in Figure 1. The interaction levels (L1/L2/L3) represent: L1 for isolated token processing, L2 for local token routing within samples, and L3 for global token routing across sample. Table 8 presents a comprehensive comparison of different token interaction strategies in diffusion models. The baseline models with L1 strategy (TC-DiT-L-Flow and
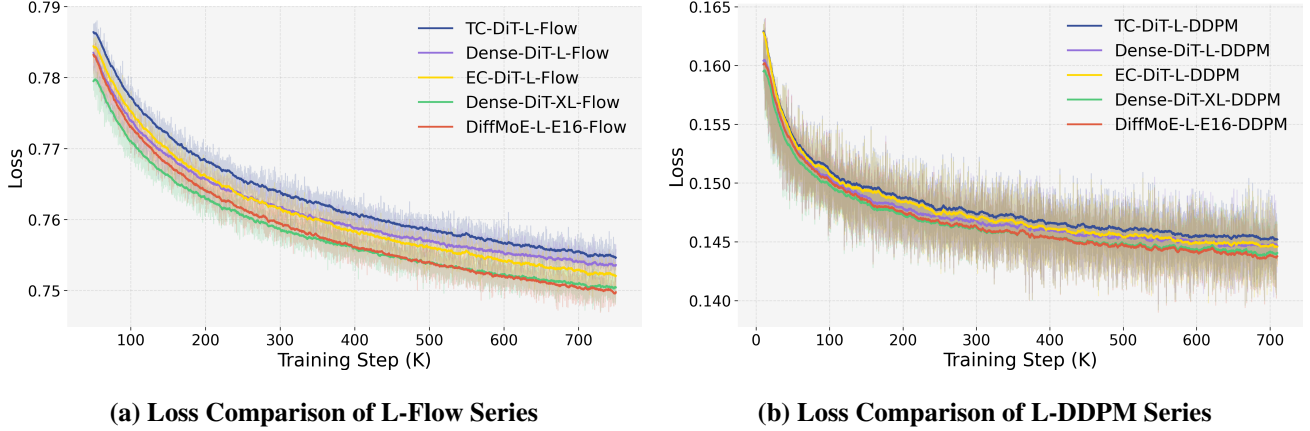
(a) Loss Comparison of L-Flow Series

(b) Loss Comparison of L-DDPM Series

Figure 10: **Loss Comparison of L-Flow and L-DDPM Series**. The relative losses illustrated in subfigures (a) and (b) demonstrate the exceptional training dynamics of DiffMoE, consistently outperforming all baseline models.

Dense-DiT-L-Flow) process tokens independently, resulting in limited performance (FID: 19.06 and 17.01). The L2 strategy, implemented in EC-DiT-L-Flow, enables local token routing within samples, showing improved performance (FID: 16.12) with the same parameter count. Our proposed L3 strategy in DiffMoE-L-Flow introduces cross-sample token routing, achieving superior results (FID: 14.41) even compared to the 1.5x larger Dense-DiT-XL-Flow (675M parameters). Notably, when combined with Dynamic Global CP, our model not only achieves the best FID score but also reduces the computational capacity to 0.95x, demonstrating both effectiveness and efficiency.

### D.4. Dynamic Conditional Computation: Harder Work needs More Computation

Figure 1 demonstrates that different classes require varying computational resources during generation. To analyze this variation, we sample 1K different class labels in a batch and rank their $C_{\text{infer}}^{\text{avg}}$ in descending order, revealing the computational complexity of generation across classes. The top-10 most computationally intensive classes for both flow-based and DDPM models are displayed in Figure 13. The top-10 least computationally intensive classes for both flow-based and DDPM models are displayed in Figure 14.
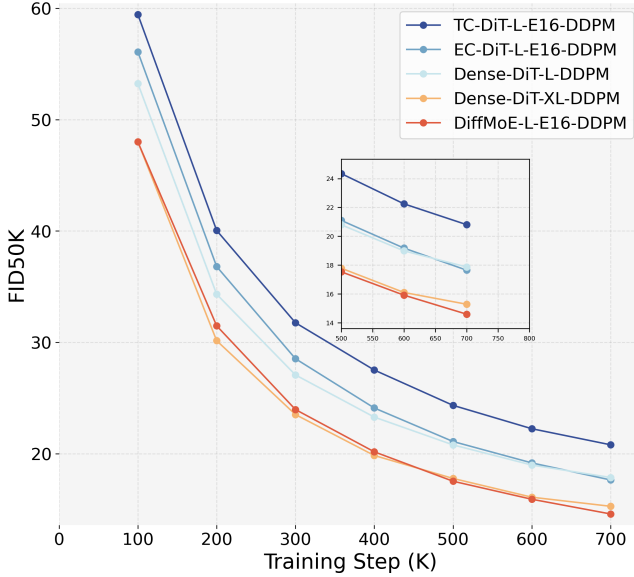
### D.5. Dynamic Token Selection across Network Layers

As illustrated in Figure 12, our analysis reveals distinctive expert utilization patterns across different network depths. The shallow Layer 1 exhibits pronounced fluctuations and sharp capacity spikes, indicating intensive early-stage feature extraction. Moving to intermediate Layer 7, we observe more stabilized capacity patterns, suggesting balanced processing of mid-level features. Layer 13 demonstrates gradual, long-term capacity transitions, while the deep Layer 19 shows notably uniform expert utilization. This systematic progression from volatile to stable expert engagement reflects the natural specialization of experts: from low-level feature detection in early layers to refined semantic processing in deeper layers. Such hierarchical organization of expert behaviors aligns with the progressive nature of diffusion-based generation.
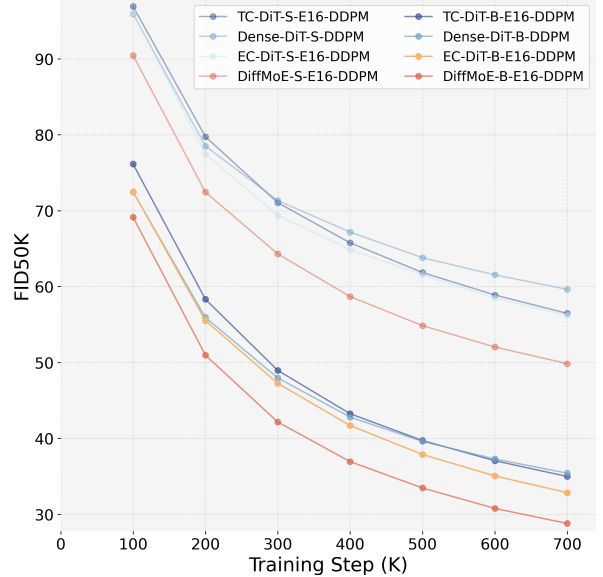
## E. FID Sensibility and Ablation Study

FID scores are sensitive to implementation details, necessitating careful ablation studies to understand the differences between various implementations. Through these studies, we aim to provide the academic community with clearer insights for fair comparisons between diffusion models.

The observed FID degradation at higher CFG scales is well-documented (Ma et al., 2024), primarily due to ImageNet's diverse image quality distribution. When generating high-quality samples, the deviation from ImageNet's mixed-quality dataset can lead to increased FID scores, despite improved visual quality.

For FID calculation, we follow the implementations from SiT (Ma et al., 2024) and DiT (Peebles & Xie, 2023a). Results marked with $^\dagger$ are directly quoted from (Ma et al., 2024) (DDPM) and (Peebles & Xie, 2023a) (Flow). For results marked with $^*$, we reproduce the experiments using officially released checkpoints under identical evaluation conditions.

(a) **Comparisons with the Baseline Size-L (DDPM).**

(b) **Comparisons with the Baseline Size-S/B (DDPM).**

Figure 11: **Comparisons with the Baseline Models. (a)** We compare TC, EC, and Dense Models and show the average activated parameters of all experts across all sampling steps. DiffMoE-L-E16-DDPM even surpasses DenseDiT-XL-DDPM (1.5x params). **(b)** We also examine S/B size DiffMoE models to further demonstrate the scalability.

### E.1. VAE Decoder Ablations

Following (Peebles & Xie, 2023a), throughout our experiments, we employed pre-trained VAE models. Specifically, we utilized fine-tuned versions (ft-MSE and ft-EMA) of the original LDM "f8" model, where only the decoder weights were fine-tuned. For the analysis presented in Experiments section 3 and Tables 11 and 10, we tracked metrics using the ft-MSE decoder, while the final metrics reported in Table 2 was obtained using the ft-EMA decoder. In this section, we examine the impact of two distinct VAE decoders for our experiments - the two fine-tuned variants employed in Stable Diffusion. Since all models share identical encoders, we can interchange decoders without necessitating diffusion model retraining. As demonstrated in Table 12, the DiffMoE model maintains its superior performance over existing diffusion models.

### E.2. Flow ODE-Sampler Ablations

Higher-order ODE samplers generally achieve better FID scores. As shown in Table 15, the black-box dopri5 sampler outperforms heun (NFE=250), which in turn surpasses euler (NFE). For fair comparison with baseline models, we employ the euler sampler in flow-based experiments. However, to benchmark against SiT-XL (Dense-DiT-XL-Flow) (Ma et al., 2024), we use the heun sampler to achieve SOTA results.

### E.3. Classifier-Free-Guidance Ablations

We evaluate different models with varying classifier-free guidance (CFG) scales in Table 14, and discover that the CFG scale of 1.5 adopted in DiT (Peebles & Xie, 2023a) and SiT (Ma et al., 2024) studies may not be universally optimal. Our analysis reveals the best CFG scale approximates to 1.43 through comprehensive comparisons. However, different models exhibit distinct characteristics that lead to varying optimal CFG scales, suggesting that fixing a uniform scale across all models could introduce evaluation bias. To ensure relatively fair comparisons while maintaining consistency with established practices, we ultimately adopt CFG 1.5 as the default setting in our experiments. This decision aligns with the well-documented trade-off in diffusion models: higher CFG scales (*e.g.*, 4.0) typically enhance image fidelity at the cost of increased FID scores, while lower scales (*e.g.*, 1.5) yield better FID metrics despite reduced perceptual quality. This phenomenon primarily stems from FID's sensitivity to distributional coverage - higher guidance scales tend to produce samples with reduced diversity that more closely match the training distribution statistics, paradoxically resulting in worse FID scores despite improved individual sample quality.

Table 13: **Batch Size Ablation Study:** FID scores under different batch sizes using fine-tuned EMA VAE decoder and Heun sampler. **Bold** indicates best performance in its cell.

| Model | Training Steps | CFG | Batch Size | FID50K ↓ |
|---|---|---|---|---|
| DiffMoE-L-E8-Flow | 7000K | 1.0 | 10 | **9.60** |
| DiffMoE-L-E8-Flow | 7000K | 1.0 | 15 | 9.62 |
| DiffMoE-L-E8-Flow | 7000K | 1.0 | 32 | 9.77 |
| DiffMoE-L-E8-Flow | 7000K | 1.0 | 50 | 9.98 |
| DiffMoE-L-E8-Flow | 7000K | 1.0 | 75 | 9.90 |
| DiffMoE-L-E8-Flow | 7000K | 1.0 | 100 | 9.76 |
| DiffMoE-L-E8-Flow | 7000K | 1.0 | 125 | 9.78 |
| Dense-DiT-XL-Flow* | 7000K | 1.0 | 10 | 9.57 |
| Dense-DiT-XL-Flow* | 7000K | 1.0 | 15 | **9.47** |
| Dense-DiT-XL-Flow* | 7000K | 1.0 | 50 | 9.84 |
| Dense-DiT-XL-Flow* | 7000K | 1.0 | 75 | 9.67 |
| Dense-DiT-XL-Flow* | 7000K | 1.0 | 100 | 9.65 |
| Dense-DiT-XL-Flow* | 7000K | 1.0 | 125 | 9.64 |
| DiffMoE-L-E8-Flow | 7000K | 1.5 | 10 | 2.19 |
| DiffMoE-L-E8-Flow | 7000K | 1.5 | 32 | 2.16 |
| DiffMoE-L-E8-Flow | 7000K | 1.5 | 50 | 2.16 |
| DiffMoE-L-E8-Flow | 7000K | 1.5 | 75 | **2.13** |
| DiffMoE-L-E8-Flow | 7000K | 1.5 | 100 | 2.17 |
| DiffMoE-L-E8-Flow | 7000K | 1.5 | 125 | 2.18 |
| Dense-DiT-XL-Flow* | 7000K | 1.5 | 10 | 2.23 |
| Dense-DiT-XL-Flow* | 7000K | 1.5 | 15 | 2.23 |
| Dense-DiT-XL-Flow* | 7000K | 1.5 | 50 | 2.21 |
| Dense-DiT-XL-Flow* | 7000K | 1.5 | 75 | **2.19** |
| Dense-DiT-XL-Flow* | 7000K | 1.5 | 100 | 2.22 |
| Dense-DiT-XL-Flow* | 7000K | 1.5 | 125 | 2.21 |
| DiffMoE-L-E8-DDPM | 7000K | 1.5 | 50 | **2.30** |
| DiffMoE-L-E8-DDPM | 7000K | 1.5 | 75 | 2.33 |
| DiffMoE-L-E8-DDPM | 7000K | 1.5 | 100 | 2.32 |
| DiffMoE-L-E8-DDPM | 7000K | 1.5 | 125 | 2.32 |

Table 14: **CFG Scale Ablation Study:** FID scores across different CFG scales using fine-tuned EMA VAE decoder. **Bold** indicates best performance in its cell.

| Model | Training Steps | Sampler | CFG | Batch Size | FID50K ↓ |
|---|---|---|---|---|---|
| DiffMoE-L-E8-Flow | 4900K | Heun | 1.0 | 125 | **9.21** |
| Dense-DiT-XL-Flow* | 7000K | Heun | 1.0 | 125 | 9.64 |
| DiffMoE-L-E8-Flow | 7000K | Heun | 1.0 | 125 | 9.78 |
| DiffMoE-L-E8-Flow | 4900K | Heun | 1.43 | 125 | 2.14 |
| Dense-DiT-XL-Flow* | 7000K | Heun | 1.43 | 125 | **2.08** |
| DiffMoE-L-E8-Flow | 7000K | Heun | 1.43 | 125 | 2.13 |
| DiffMoE-L-E8-Flow | 4900K | Heun | 1.5 | 125 | 2.28 |
| Dense-DiT-XL-Flow* | 7000K | Heun | 1.5 | 125 | 2.21 |
| DiffMoE-L-E8-Flow | 7000K | Heun | 1.5 | 125 | **2.18** |
| DiffMoE-L-E8-DDPM | 6500K | DDPM | 1.0 | 125 | 9.39 |
| Dense-DiT-XL-DDPM* | 7000K | DDPM | 1.0 | 125 | 9.63 |
| DiffMoE-L-E8-DDPM | 7000K | DDPM | 1.0 | 125 | **9.17** |
| DiffMoE-L-E8-DDPM | 6500K | DDPM | 1.5 | 125 | **2.27** |
| Dense-DiT-XL-DDPM* | 7000K | DDPM | 1.5 | 125 | 2.32 |
| DiffMoE-L-E8-DDPM | 7000K | DDPM | 1.5 | 125 | 2.32 |

Table 15: **Flow ODE Sampler Ablation Study:** FID scores across different ODE samplers with CFG scale 1.0 and fine-tuned EMA VAE decoder. **Bold** indicates best performance in its cell.

| Model | Training Steps | Sampler | Batch Size | FID50K ↓ |
|---|---|---|---|---|
| DiffMoE-L-E8-Flow | 4900K | Euler | 125 | 9.37 |
| DiffMoE-L-E8-Flow | 4900K | Heun | 125 | 9.21 |
| DiffMoE-L-E8-Flow | 4900K | Euler | 250 | 9.39 |
| DiffMoE-L-E8-Flow | 4900K | Dopri5 | 250 | **9.06** |
| DiffMoE-L-E8-Flow | 7000K | Euler | 125 | 9.86 |
| DiffMoE-L-E8-Flow | 7000K | Heun | 125 | 9.78 |
| DiffMoE-L-E8-Flow | 7000K | Euler | 250 | 9.94 |
| DiffMoE-L-E8-Flow | 7000K | Dopri5 | 250 | **9.56** |

### E.4. Batch Sizes Ablations

The batch sizes ablation study reveals critical insights into the interplay between batch size and classifier-free guidance (CFG) scales for the DiffMoE-L-E8-Flow model. At CFG=1.0, FID scores remain elevated (9.60–9.98), with smaller batch sizes (*e.g.*, bs=10) marginally outperforming larger configurations, exhibiting a U-shaped trend. However, elevating CFG to 1.5 drastically reduces FID to 2.13–2.19, achieving optimal performance at bs=75 (**2.13**), while demonstrating remarkable robustness to batch size variations ($\Delta$=0.06 vs. $\Delta$=0.38 at CFG=1.0).

### E.5. Conclusion: A little thought about FID

While Fréchet Inception Distance (FID) is widely adopted for evaluating generative models, particularly on ImageNet, it exhibits several notable limitations. Our analysis reveals counterintuitive behaviors, especially when evaluating models with classifier-free guidance (CFG). For example, higher CFG scales typically enhance perceptual quality but paradoxically result in worse FID scores, despite producing visually superior images. This discrepancy stems from FID's fundamental mechanism: it measures statistical similarities between generated and real distributions in the Inception network's feature space, often failing to capture perceptual quality and fine-grained details. Moreover, FID scores are susceptible to various implementation factors, including choice of ODE samplers, hardware configurations, random seeds, and sample size for estimation. These sensitivities can impact reproducibility and comparison across different studies. Furthermore, FID's focus on distributional overlap overlooks critical aspects such as mode collapse and overfitting, as it does not explicitly evaluate sample diversity or novelty. These limitations underscore the pressing need for more robust and comprehensive metrics that can better reflect the true modeling capabilities of generative models. We advocate for developing new evaluation frameworks that combine precision-recall curves, perceptual quality metrics, and human evaluation studies, which would provide a more reliable assessment of generative model performance.

# F. Visual Generation Results

### F.1. Class-Conditional Image Generation

To demonstrate the generation capabilities of our model, we showcase diverse images sampled from DiffMoE-L-E8-Flow and DiffMoE-L-E8-DDPM, conditioned on ImageNet class labels. These visualizations illustrate the model's ability to generate high-quality, class-specific images. See Figure 15 and 16.

### F.2. Text-Conditional Image Generation

We present a collection of images generated by our DiffMoE-T2I-Flow model using various text prompts as conditioning inputs. These examples demonstrate the model's versatility in translating textual descriptions into corresponding visual representations. See Figure 18.

**(a)** Capacity of All Experts in Layer 1

**(b)** Capacity of All Experts in Layer 7

**(c)** Capacity of All Experts in Layer 13

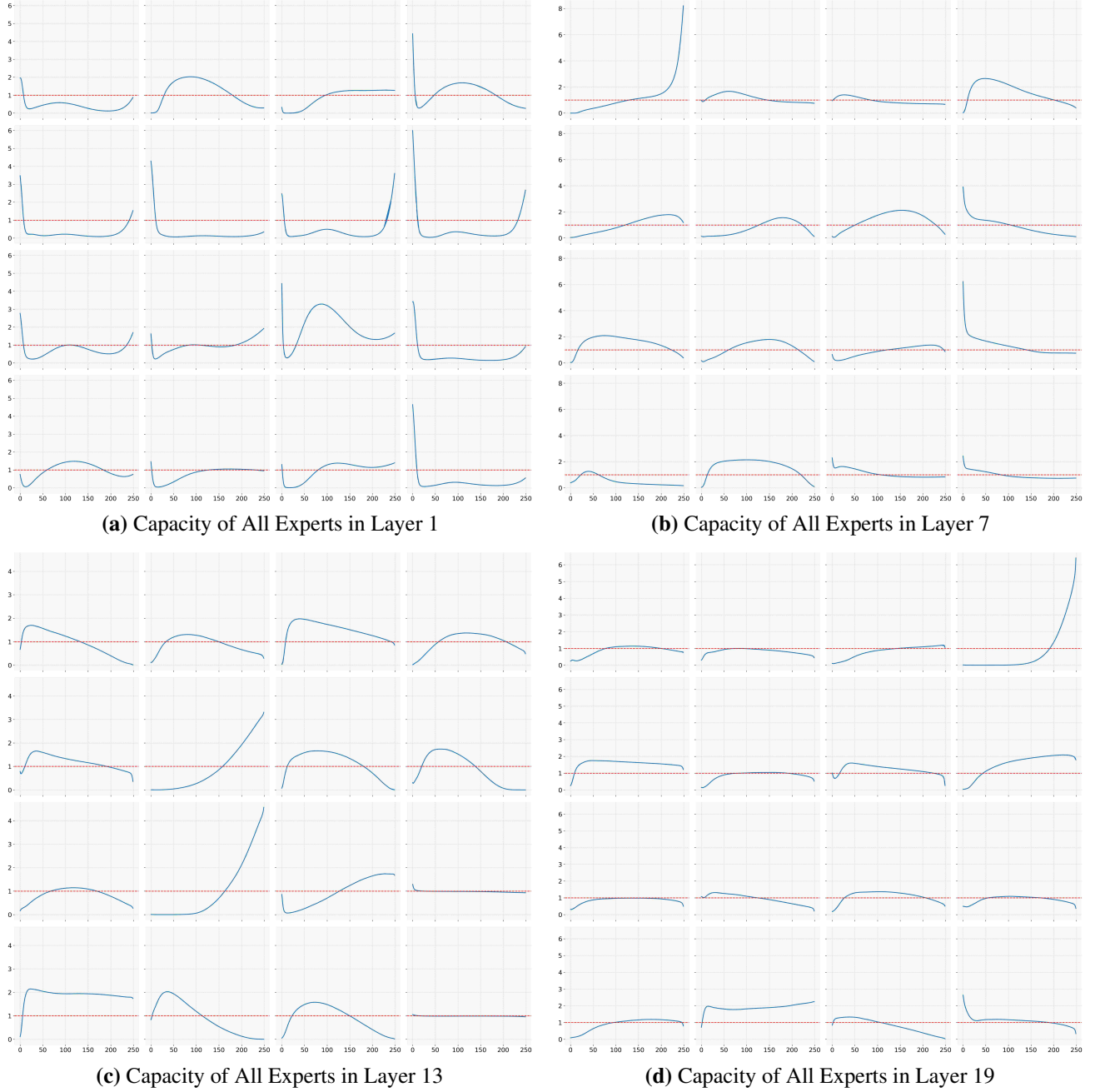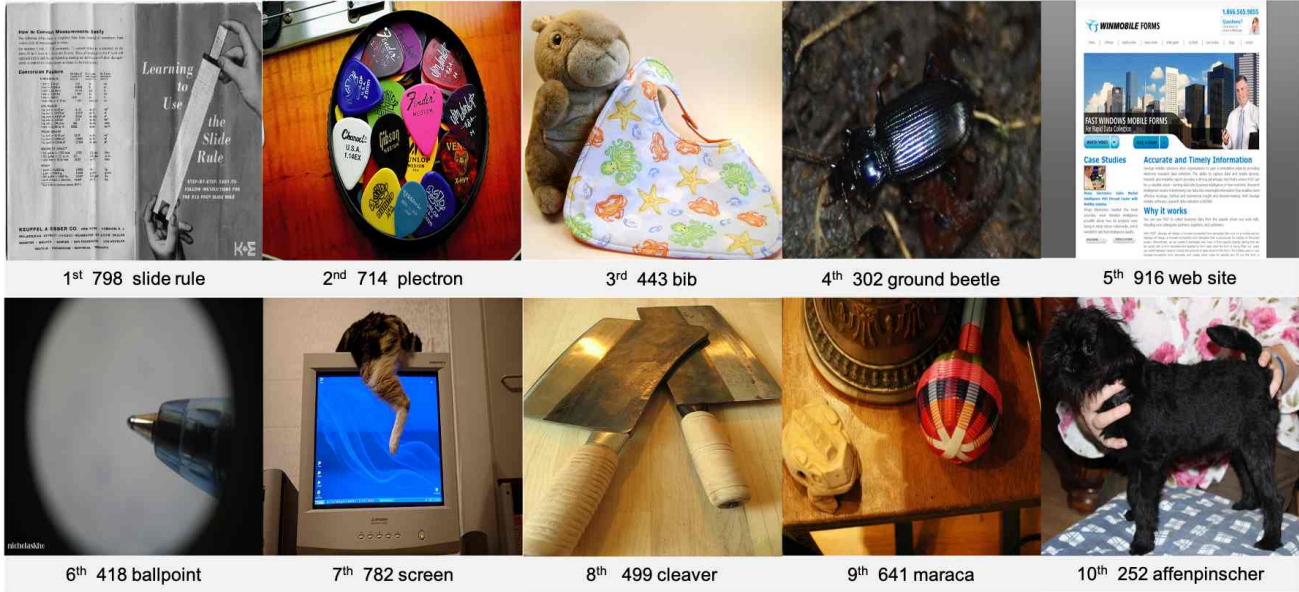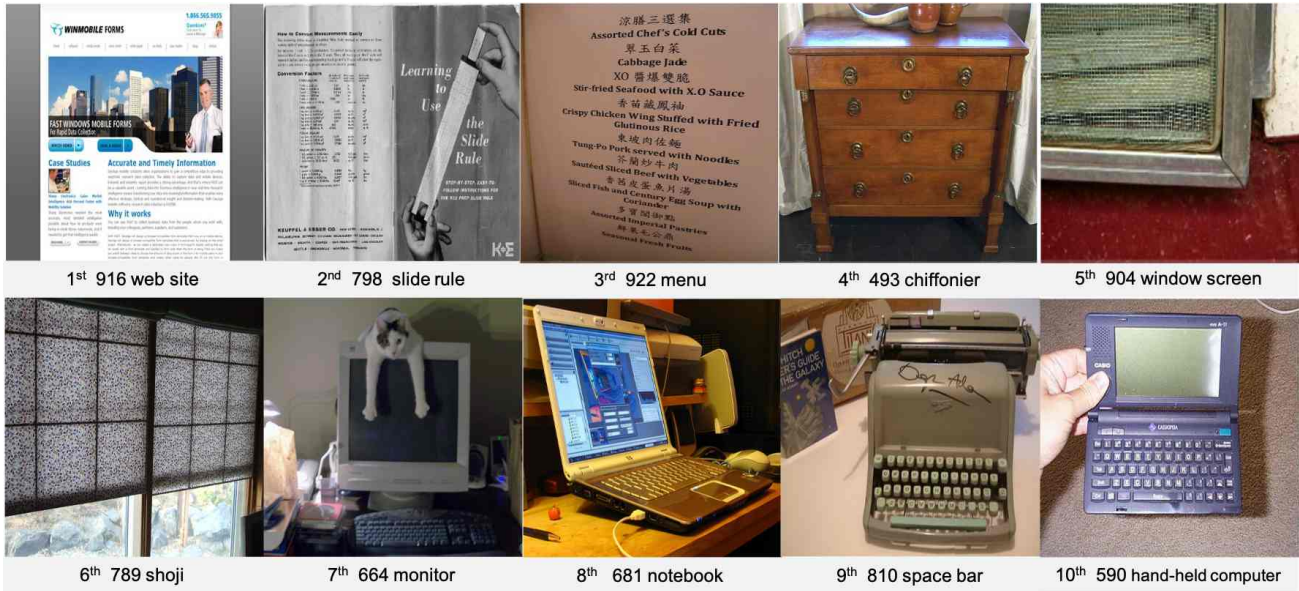**(d)** Capacity of All Experts in Layer 19

Figure 12: **Expert Dynamics across Network Layers.** Visualization of expert capacity patterns in network layers (1, 7, 13, 19). Early layers show high-amplitude fluctuations, while deeper layers exhibit increasingly stable utilization, demonstrating natural expert specialization throughout the diffusion process.
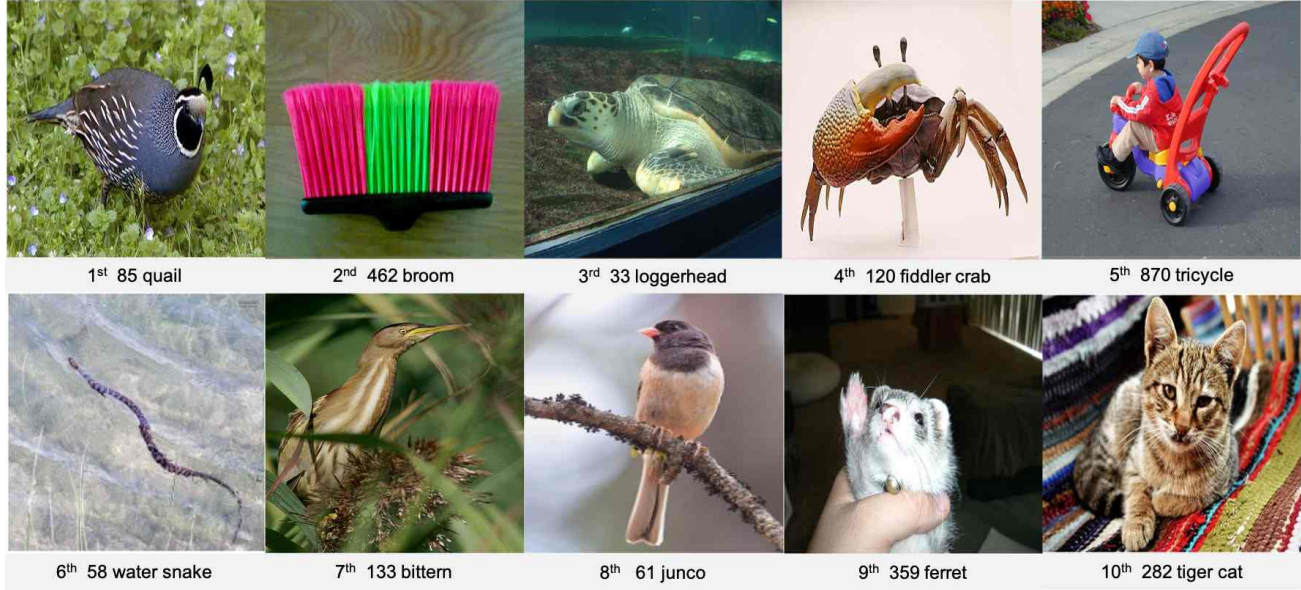
| 1st 798 slide rule | 2nd 714 plectron | 3rd 443 bib | 4th 302 ground beetle | 5th 916 web site |
| 6th 418 ballpoint | 7th 782 screen | 8th 499 cleaver | 9th 641 maraca | 10th 252 affenpinscher |

**(a) DiffMoE-L-E16-Flow (700K)**



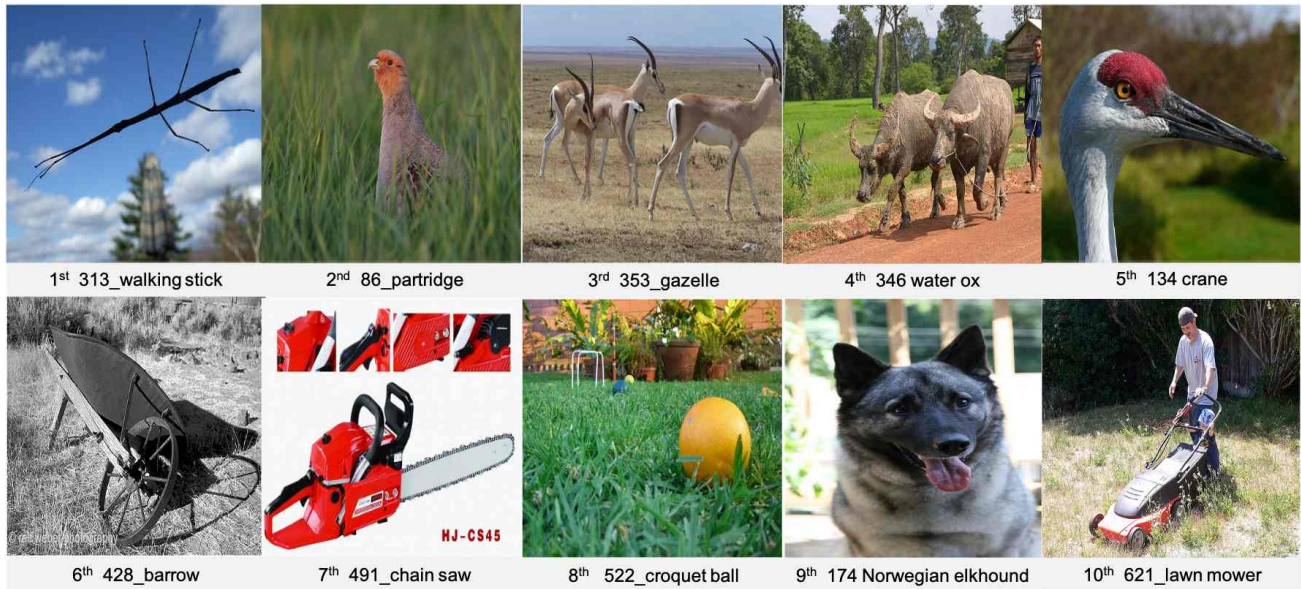| 1st 916 web site | 2nd 798 slide rule | 3rd 922 menu | 4th 493 chiffonier | 5th 904 window screen |
| 6th 789 shoji | 7th 664 monitor | 8th 681 notebook | 9th 810 space bar | 10th 590 hand-held computer |

**(b) DiffMoE-L-E16-DDPM (700K)**

Figure 13: **Top 10 Hardest Classes.** The 10 classes with the highest computational cost, sampled from the training set.

1st 85 quail  2nd 462 broom  3rd 33 loggerhead  4th 120 fiddler crab  5th 870 tricycle

6th 58 water snake  7th 133 bittern  8th 61 junco  9th 359 ferret  10th 282 tiger cat

**(a) DiffMoE-L-E16-Flow (700K)**



1st 313_walking stick  2nd 86_partridge  3rd 353_gazelle  4th 346 water ox  5th 134 crane

6th 428_barrow  7th 491_chain saw  8th 522_croquet ball  9th 174 Norwegian elkhound  10th 621_lawn mower

**(b) DiffMoE-L-E16-DDPM (700K)**

Figure 14: **Top 10 Easiest Classes.** The 10 classes with the lowest computational cost, sampled from the training set.

---

**Algorithm 1** EC-DiT Layer

---

**Input:** $x$ (input tensor)
**Variables:** $B$ (batch size), $S$ (sequence length), $d$ (hidden dim), $W_r$ (routing weights), $experts$ (list of expert FFNs)
$\qquad$ $E$ (number of experts), $C$ (expert capacity)
/* Step 1: Compute Token-Expert Affinity Matrix */
$\text{logits} \leftarrow \text{einsum}('bsd, de \rightarrow bse', x, W_r)$
$\text{scores} \leftarrow \texttt{softmax}(logits, \dim = -1).\texttt{permute}(-1, -2)$
/* Step 2: Select top-k tokens for each expert */
$\text{gating, index} \leftarrow \text{top\_k}(\text{scores}, \text{k=}C, \dim=-1)$
$\text{dispatch} \leftarrow \text{one\_hot}(\text{index}, \text{num\_classes=}S)$
/* Step 3: Process tokens through experts and combine */
$x_{in} \leftarrow \text{einsum}('becs, bsd \rightarrow becd', \text{dispatch}, x)$
$x_e \leftarrow [\text{experts}[e](x_{in}[:, e]) \text{ for } e \text{ in range}(E)]$
$x_e \leftarrow \text{stack}(x_e, \dim=1)$
$x_{out} \leftarrow \text{einsum}('becs, bec, becd \rightarrow bsd', \text{dispatch}, \text{gating}, x_e)$
**Return:** $x_{out}$

---

**Algorithm 2** TC-DiT layer

---

**Input:** $x$ (input tensor), $W_r$ (routing weights), $experts$ (list of expert FFNs)
**Variables:** $B$ (batch size), $S$ (sequence length), $d$ (hidden dim) ,$K$ (experts per token)
/* Step 1: Save original input shape */
$\text{orig\_shape} \leftarrow \texttt{shape}(x)$
/* Step 2: Compute Token-Expert Affinity Matrix */
$\text{logits} \leftarrow \text{einsum}('bsd, de \rightarrow bse', x, W_r)$
$\text{scores} \leftarrow \texttt{softmax}(logits, \dim = -1)$
/* Step 3: Select top-k tokens for each expert */
$\text{gating, index} \leftarrow \text{top\_k}(\text{scores}, \text{k=}K, \dim=-1)$
/* Step 4: Flatten $x$ and top-k indices */
$x \leftarrow \texttt{view}(x, (-1, x.shape[-1]))$
$\text{flat\_topk\_idx} \leftarrow \texttt{view}(topk\_idx, (-1))$
/* Step 4: Process tokens through experts */
$x \leftarrow \texttt{repeat\_interleave}(x, K, \dim = 0)$
$y \leftarrow \texttt{empty\_like}(x)$
**for** $i \leftarrow 1$ **to** $len(experts)$ **do**
$\quad y[\text{flat\_topk\_idx} == i] \leftarrow \texttt{expert}_i(x[\text{flat\_topk\_idx} == i])$
**end for**
$y \leftarrow \text{sum}(\texttt{view}(y, (*\text{gating}.shape, -1)) \cdot \text{gating}.unsqueeze(-1), \dim = 1)$
$y \leftarrow \texttt{view}(y, \text{orig\_shape})$
**Return:** $y$

---

---

**Algorithm 3** DiffMoE layer (Training)

---

**Input:** $x$ (input tensor)
**Variables:** $B$ (batch size), $S$ (flattened sequence length), $D$ (hidden dim), $N$ (number of experts)
          $W_r$ (routing weights), $experts$ (list of expert FFNs), $C$ (expert capacity)
/* Step 1: Batch-level token pool and compute capacity prediction */
$x \leftarrow \texttt{view}(x, (-1, D))$
$S \leftarrow \texttt{shape}(x)[0]$
$capacity\_pred \leftarrow \texttt{capacity\_predictor}(\texttt{detach}(x))$
$C_{\text{train}} \leftarrow \texttt{int}((S/N) \times C)$
/* Step 2: Compute token-expert affinity scores */
$logits \leftarrow \texttt{einsum}('sd, de \rightarrow se', x, W_r)$
$scores \leftarrow \texttt{softmax}(logits, \dim = -1).\texttt{permute}(-1, -2)$
$gating, index \leftarrow \texttt{top\_k}(scores, k = C_{\text{train}}, \dim = -1, sorted = \texttt{False})$
/* Step 3: Process tokens through experts */
$y \leftarrow \texttt{zeros\_like}(x)$
$ones \leftarrow \texttt{zeros}(N, S)$
**for** $i \leftarrow 1$ **to** N **do**
    $y[index[i], :] \leftarrow y[index[i], :] + gating[i].\texttt{unsqueeze}(-1) \times \texttt{expert}_i(x[index[i], :])$
    $ones[i][index[i]] \leftarrow 1.$
**end for**
/* Step 4: Update capacity threshold */
$\texttt{update\_threshold}(capacity\_pred)$
/* Step 5: Reshape output */
$x_{out} \leftarrow \texttt{view}(y, (B, s, D))$
**Return:** $x_{out}, ones, capacity\_pred$

---

**Algorithm 4** DiffMoE layer (Inference)

---

**Input:** $x$ (input tensor)
**Variables:** $B$ (batch size), $S$ (flattened sequence length), $D$ (hidden dim), $N$ (number of experts)
          $W_r$ (routing weights), $experts$ (list of expert FFNs), $C$ (expert capacity)
          $threshold$ (expert threshold)
/* Step 1: Reshape input and compute capacity prediction */
$x \leftarrow \texttt{view}(x, (-1, D))$
$S \leftarrow \texttt{shape}(x)[0]$
$capacity\_pred \leftarrow \texttt{sigmoid}(\texttt{capacity\_predictor}(\texttt{detach}(x)))$
/* Step 2: Compute token-expert affinity scores */
$logits \leftarrow \texttt{einsum}('sd, de \rightarrow se', x, W_r)$
$scores \leftarrow \texttt{softmax}(logits, \dim = -1).\texttt{permute}(-1, -2)$
/* Step 3: Process tokens through experts */
$y \leftarrow \texttt{zeros\_like}(x)$
**for** $i \leftarrow 1$ **to** N **do**
    $k_{pred} \leftarrow \texttt{sum}(\texttt{where}(capacity\_pred[:, i] > threshold[i], 1, 0))$
    $gating, index \leftarrow \texttt{top\_k}(scores[i], k = k_{pred}, \dim = -1, sorted = \texttt{False})$
    $y[index, :] \leftarrow y[index, :] + gating.\texttt{unsqueeze}(-1) \times \texttt{expert}_i(x[index, :])$
**end for**
/* Step 4: Reshape output */
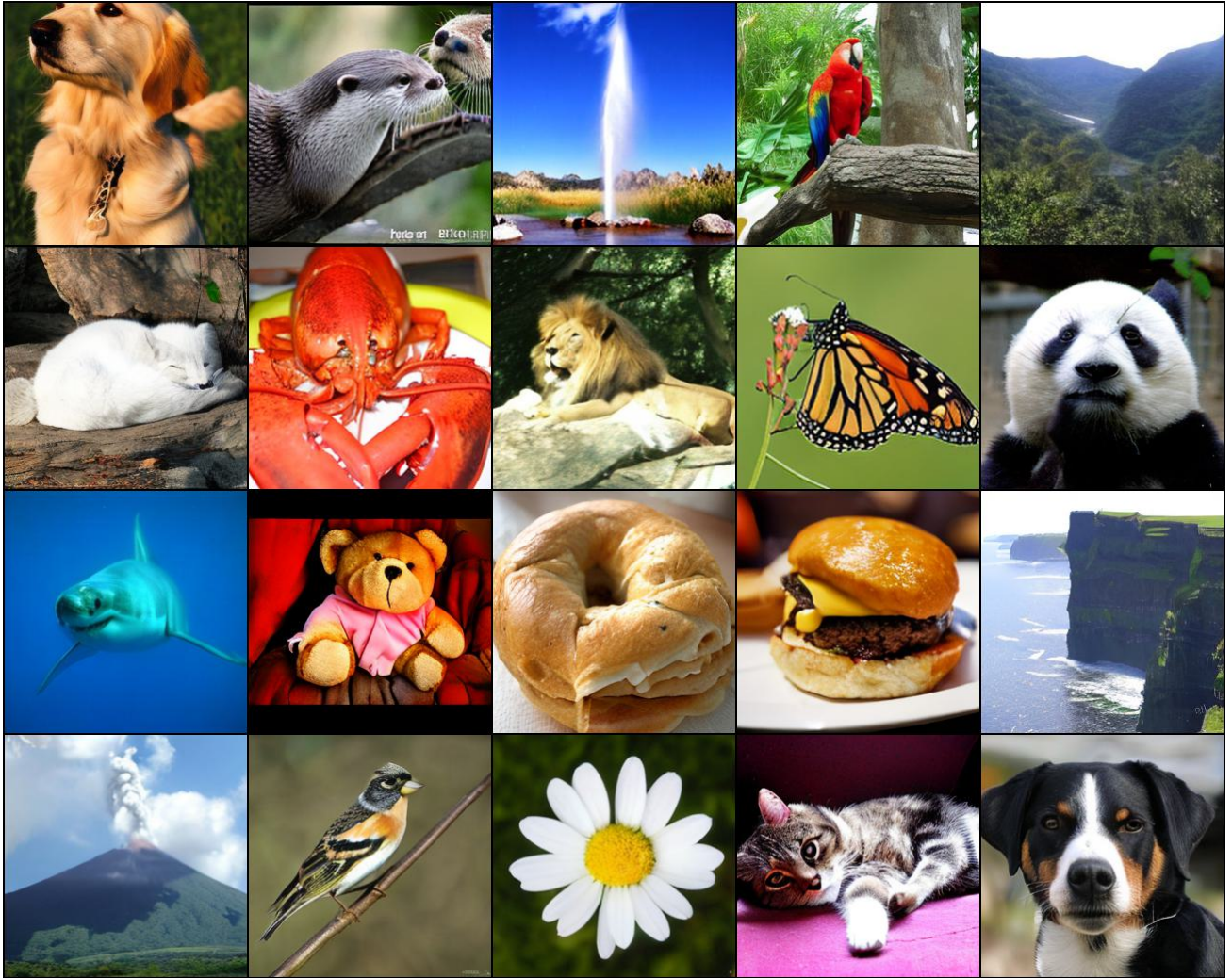$x_{out} \leftarrow \texttt{view}(y, (B, s, D))$
**Return:** $x_{out}$

---

Figure 15: **Class-conditional Generation. Uncurated** $256 \times 256$ **DiffMoE-L-E8-Flow** samples CFG scale = 4.0.
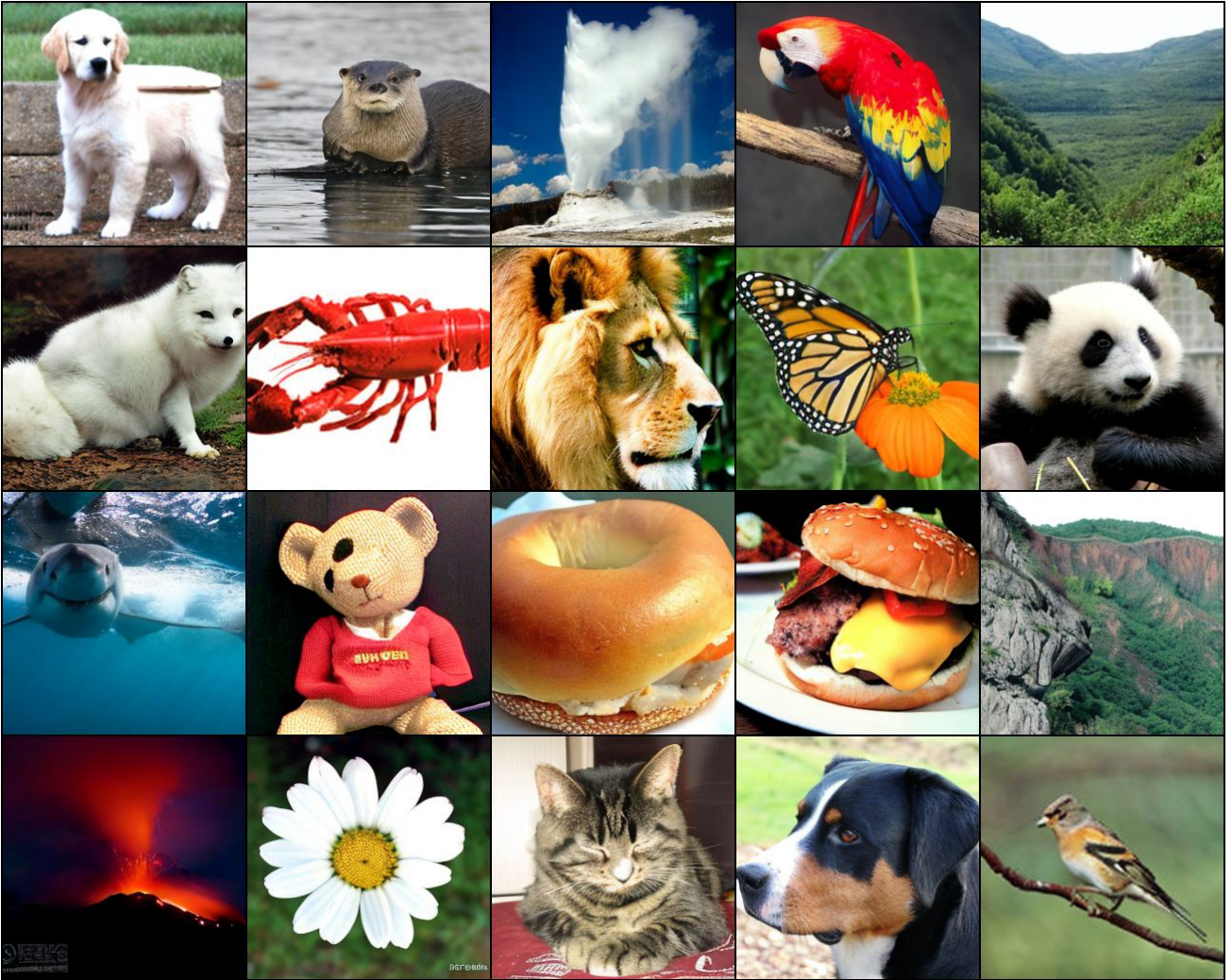
Figure 16: **Class-conditional Generation. Uncurated** $256{\times}256$ **DiffMoE-L-E8-DDPM** samples CFG scale = 4.0.

Figure 17: **Text-to-Image Generation.** Comparison of DiffMoE-E16-T2I-Flow (w/o SFT) and Dense Model (w/o SFT).

Figure 18: **Text-to-Image Generation. Uncurated** 256×256 Images generated by **DiffMoE-E16-T2I-Flow** (w SFT)