DrDiff: Dynamic Routing Diffusion with Hierarchical Attention for Breaking the Efficiency-Quality Trade-off

Jusheng Zhang¹, Yijia Fan¹, Kaitong Cai¹, Zimeng Huang¹, Xiaofei Sun², Jian Wang³, Chengpei Tang¹, Keze Wang^{1,†}

¹Sun Yat-sen University ²Alibaba Group ³Snap Inc

[†]Corresponding author: kezewang@gmail.com

Abstract

This paper introduces DrDiff, a novel framework for long-text generation that overcomes the efficiency-quality trade-off through three core technologies. First, we design a dynamic expert scheduling mechanism that intelligently allocates computational resources during the diffusion process based on text complexity, enabling more efficient handling of text generation tasks of varying difficulty. Second, we introduce a Hierarchical Sparse Attention (HSA) mechanism that adaptively adjusts attention patterns according to a variety of input lengths, reducing computational complexity from $O(n^2)$ to O(n) while maintaining model performance. Finally, we propose a Semantic Anchor States (SAS) module that combines with DPM-solver++ to reduce diffusion steps, significantly improving generation speed. Comprehensive experiments on various long-text generation benchmarks demonstrate the superiority of our DrDiff over the existing SOTA methods.

1 Introduction

Large Language Models (LLMs) (Radford et al., 2019; Brown et al., 2020a; Touvron et al., 2023; OpenAI, 2024) have demonstrated remarkable capabilities in knowledge encoding and contextual understanding during their pretraining phase, achieving significant success across a variety of natural language processing tasks. However, despite these advanced abilities, LLMs encounter substantial bottlenecks when generating ultra-long texts (e.g., exceeding 10,000 tokens) (Krishna et al., 2023; Liu et al., 2024). These challenges primarily manifest in maintaining long-range coherence, managing quadratically increasing computational complexity, and ensuring contextual consistency (Bai et al., 2024b; Liu et al., 2024; Wu et al., 2025). While directly fine-tuning on long sequences might appear as a straightforward solution, this method demands

prohibitive computational resources and proves difficult to optimize effectively for documents spanning tens of thousands of tokens or more.

To address these challenges, the academic community has proposed various methods to enhance the efficiency and effectiveness of LLMs in processing long sequences. These can be broadly categorized into two main types: (1) approaches based on optimizing the attention mechanism, such as sparse attention (e.g., Longformer (Beltagy et al., 2020), etc. (Bertsch et al., 2023; Kitaev et al., 2020; Child et al., 2019a)), which reduce complexity by modifying the attention computation pattern; and (2) explorations into emerging generation paradigms tailored for long sequences, for instance, applying diffusion models (Tang et al., 2023; Becker et al., 2025) to text generation or specifically optimizing LLM training and inference strategies (Chen et al., 2023) for extended contexts. Although these methods have achieved certain progress in specific scenarios, they typically employ relatively fixed resource allocation and information processing flows when handling extremely long sequences. Neither the fixed sparse patterns of sparse attention nor the generally consistent iterative approach throughout the denoising process in diffusion models adequately considers the heterogeneous requirements of different text generation stages or varying text segments, thereby limiting their adaptability.

This inherent "fixedness" becomes particularly pronounced in scenarios involving ultra-long text generation (e.g., over 10,000 tokens), directly leading to three critical, yet unresolved, limitations: (1) *Decay of Long-Range Feature Representation*: As sequence length increases, the model struggles to compress and retain historical information, causing early-input semantic features to progressively decay or "dilute." This degradation in feature quality critically impairs the model's ability to capture long-range dependencies and is a primary cause of content repetition and logical inco-

herence in long-text generation. (2) Suboptimal computational resource allocation: Applying a uniform computational intensity to all text segments or generation stages results in resource wastage on structurally simple or information-sparse parts, while potentially providing insufficient computational power for complex or critical semantic junctures, ultimately impacting overall efficiency and performance. (3) Significant degradation in generation quality with increasing length: As sequence length drastically increases, models often struggle to maintain a unified narrative thread, exhibiting a higher tendency for content repetition, logical discontinuities, and even "forgetting" important information generated earlier. Consequently, existing methods find it challenging to achieve an ideal balance among efficiency, dynamic resource scheduling, and high-quality long-text generation.

To tackle these fundamental issues, we move beyond attempting localized improvements on existing fixed architectures and propose DrDiff, a novel dynamic generation framework. DrDiff's core principle is the dynamic adjustment of its internal processing mechanisms. To avoid sacrificing efficiency through indiscriminate complexity increases, we introduce a strategy that combines Hierarchical Sparse Attention (HSA) with a dynamically routed diffusion mechanism. The central idea of HSA is to adaptively select and combine different attention patterns (local, dilated, global) based on the current text length and content characteristics, ensuring effective dependency capture at various scales while optimizing computational complexity from $O(N^2)$ to nearly O(N). Furthermore, we deeply integrate this dynamic attention mechanism with the denoising process of diffusion models. Through Dynamic Expert Scheduling (DES), the model can allocate different computational resources (i.e., expert networks) to text segments or generation steps of varying complexities. Additionally, semantic anchor state guidance is employed to optimize diffusion paths and attention allocation, which, combined with efficient solvers such as DPM-solver++, further enhances generation efficiency. Essentially, this design ensures that computational resources and attention focus are intelligently allocated based on real-time demands, rather than adhering to predefined patterns. Specifically, DrDiff introduces three key innovations:

 Dynamic Hierarchical Sparse Attention: Adaptively selects different attention mecha-

- nisms based on text complexity and length, significantly enhancing long-range dependency modeling capabilities while reducing computational complexity.
- Dynamic Expert Scheduling and Diffusion: Dynamically allocates computational resources during the generation process, efficiently integrating dynamic attention with diffusion models to improve generation quality and efficiency.
- Semantic Anchor State Optimized Inference Path: Leverages a semantic anchor mechanism to optimize diffusion paths and attention allocation, further improving the coherence and efficiency for ultra-long text generation.

2 Related Work

Long-Text Generation and Transformer Ar-Since its introduction, the Transformer (Raffel et al., 2023; Brown et al., 2020b; Tay et al., 2022a; Dai et al., 2019; Wang et al., 2020; Child et al., 2019b; Zhang et al., 2025b,a) architecture has demonstrated remarkable performance in natural language processing (NLP) tasks, particularly in long-text generation (Bai et al., 2024c; Guan et al., 2021). However, despite its advantages in maintaining textual coherence and contextual consistency, its computational complexity of $O(n^2)$ (Child et al., 2019b; Beltagy et al., 2020) when handling extremely long sequences presents severe computational resource bottlenecks (Ashkboos et al., 2024; Chamberlain et al., 2008). More recently, researchers have explored various optimization techniques to enhance the efficiency and quality of Transformers in long-text generation, e.g., models such as BERT and GPT perform exceptionally well in generating short to medium-length texts. However, their performance significantly degrades when handling generation tasks exceeding 10,000 tokens, limiting their practical applicability. Sparse Attention Mechanisms. To improve the efficiency of Transformers (Vaswani et al., 2023) in processing long sequences, sparse attention mechanisms have emerged as effective solutions (Child et al., 2019c; Lan et al., 2020; Tay et al., 2022b; Clark et al., 2020; Chen et al., 2021; Rasooli and Tetreault, 2015). Longformer, for example, restricts each token's attention span to a local neighborhood, reducing the computational complexity

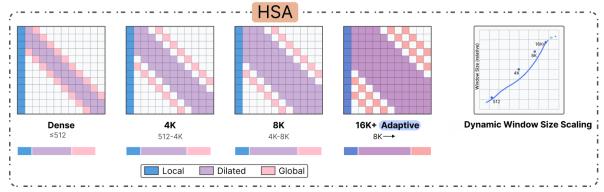


Figure 1: The Hierarchical Sparse Attention (HSA) mechanism dynamically adjusts attention patterns for texts of different lengths: using dense attention for short texts within 512 tokens, local + sparse dilated attention combination for 512–4K, sparse dilated + global attention for 4K–8K texts, and global attention with larger windows for texts above 8K. This hierarchical strategy reduces computational complexity from $\mathcal{O}(n^2)$ to $\mathcal{O}(n)$ while maintaining model performance.

from $O(n^2)$ to O(nw) (Vaswani et al., 2023), where w represents the window size. BigBird (Devlin et al., 2019) further integrates local attention, random attention, and global attention, maintaining linear complexity while capturing global dependencies more effectively. Although sparse attention enhances long-text processing efficiency, these methods still struggle to model global semantic dependencies and maintain coherence as sequence length increases. Prior studies (Gao et al., 2024; Tang et al., 2024) have introduced hierarchical attention structures and dynamic (Shazeer et al., 2017) attention patterns to strengthen sparse attention's capability, providing new avenues for improving long-text generation quality.

Application of Diffusion Models in Text Generation. Diffusion models (Ho et al., 2020; Yuan et al., 2023) have shown great potential in structured text generation (Lin et al., 2023; Xu et al., 2025; Mirbeygi and Beigy, 2025), leveraging an iterative denoising process to produce high-quality (Nichol and Dhariwal, 2021), complex text while capturing fine-grained semantics (Beltagy et al., 2020). However, in the long-sequence generation, they suffer from slow convergence and high training costs, especially when handling sequences exceeding 10,000 tokens, where computational demands escalate sharply, reducing efficiency. Ensuring stability and quality control in long-text generation remains challenging. Researchers have proposed various optimizations, such as DiffSeq (Gong et al., 2023), which significantly improve efficiency and quality. Yet, a key challenge persists: enhancing coherence and semantic consistency while minimizing computational costs. To achieve this, we propose our **DrDiff** framework, which integrates

dynamic sparse diffusion routing and hierarchical attention mechanisms to fundamentally improve efficiency and quality in long-text generation.

3 Methodology

3.1 Framework Overview

DrDiff is architected as a diffusion model-based framework for text generation, designed to tackle the aforementioned challenges inherent in producing ultra-long texts. As illustrated in Figure 2, an input text sequence, denoted as \mathbf{x} , is first tokenized and processed through an embedding layer to yield an initial token embedding sequence $Z_0 \in \mathbb{R}^{N \times d}$, where N is the sequence length and d is the embedding dimension.

The forward process of the diffusion model is defined as a Markov chain that gradually introduces Gaussian noise to Z_0 over T steps:

$$q(Z_t \mid Z_{t-1}) = \mathcal{N}\left(Z_t; \underbrace{\sqrt{1-\beta_t}Z_{t-1}}_{\text{attenuated previous state}}, \underbrace{\beta_t \mathbf{I}}_{\text{injected noise variance}}\right) \tag{1}$$

where $t \in \{1,...,T\}$, and $\beta_t \in (0,1)$ are predefined noise schedule parameters. Using $\alpha_t = 1 - \beta_t$ and $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$, the noisy state Z_t at any timestep t can be directly derived from Z_0 :

$$Z_t = \underbrace{\sqrt{\bar{\alpha}_t} Z_0}_{\text{original data component}} + \underbrace{\sqrt{1 - \bar{\alpha}_t} \pmb{\epsilon}}_{\text{cumulative noise component}}$$

where $\epsilon \sim \mathcal{N}(0, \mathbf{I})$. The reverse denoising process involves learning a neural network $\epsilon_{\theta}(Z_t, t)$, parameterized by θ , to predict the noise ϵ added at timestep t given Z_t . The core innovation of DrDiff lies within its denoising network ϵ_{θ} , which deeply integrates *Hierarchical Sparse Attention (HSA)*, a *Dynamic Routing Diffusion* mechanism based on Mixture of Experts (MoE), and employs *Semantic Anchor States (SAS)* for explicit path guidance

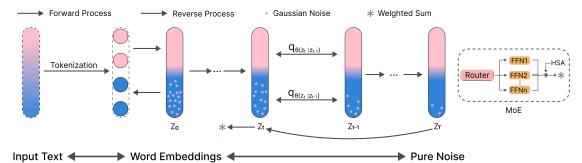


Figure 2: This diagram illustrates the diffusion process in text generation, transforming input text through tokenization into word embeddings (Z_0) , then following a forward process that gradually adds Gaussian noise to create pure noise (Z_t) . The reverse process uses a denoising model with $q_{\phi}(Z_t|Z_{t-1})$ transitions to reconstruct the text, enhanced by a Mixture-of-Experts (MoE) architecture with a Router directing computation to specialized FFN experts, and incorporating Hierarchical Sparse Attention (HSA) for efficient processing of varying text lengths.

and optimization. These components operate synergistically to achieve efficient and high-quality generation of ultra-long sequences.

3.2 Hierarchical Sparse Attention (HSA)

To overcome the $O(N^2)$ computational bottleneck of standard self-attention and to dynamically adapt to the dependency modeling requirements of texts with varying lengths, DrDiff incorporates $Hierarchical\ Sparse\ Attention\ (HSA)$ within the Transformer modules of its denoising network ϵ_{θ} . The central concept of HSA is the dynamic construction of a sparse attention mask $M^{HSA}(N)$ based on the input sequence length N, which dictates the scope of attention computations. As depicted in Figure 1, standard self-attention is computed as Attention $(Q,K,V)=\operatorname{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$. HSA modifies the attention scores (prior to softmax) using the mask:

$$S_{ij}^{HSA} = \begin{cases} \frac{Q_i K_j^T}{\sqrt{d_k}} & \text{if } (M^{HSA}(N))_{ij} = 1\\ -\infty & \text{if } (M^{HSA}(N))_{ij} = 0 \end{cases}$$
(3)

The construction strategy for $M^{HSA}(N)$ hierarchically combines several base sparse patterns, with its parameters (such as window size w, dilation rate δ , number/proportion of global nodes, and length thresholds N_1, N_2, N_3) determined through analysis of target task data characteristics and preliminary hyperparameter tuning. This strategy unfolds as follows: for short sequences $(N \leq N_1)$, dense attention $(M^{HSA}(N) = M^{dense})$ is employed to capture all local details; for medium sequences $(N_1 < N \leq N_2)$, a combination of local attention with a fixed window w_1 and dilated attention with a dilation rate δ_1 $(M^{HSA}(N) = M^{local}(w_1) \vee M^{dilated}(\delta_1))$ is used to effectively cover short to medium-range dependencies; for

longer sequences $(N_2 < N \le N_3)$, dilated attention with a larger dilation rate δ_2 is combined with global attention based on a pre-selected set of global nodes G_1 (e.g., tokens selected at regular intervals or based on learnable importance scores) $(M^{HSA}(N) = M^{dilated}(\delta_2) \vee M^{global}(G_1));$ and for ultra-long sequences $(N > N_3)$, the mechanism primarily relies on global attention $(M^{HSA}(N))$ $M^{global}(G_2, W_q(N)))$, where the global node selection strategy G_2 may be more dynamic (e.g., based on saliency from previous layer outputs), and its effective "window size" or attention span $W_a(N)$ employs a dynamic scaling strategy, such as $W_a(N) = \min(N, c \cdot N/\log N)$ or another smooth growth function, to adapt to inputs of arbitrary length and ensure global information capture. The innovation of HSA, distinct from existing fixed sparse patterns, lies in its hierarchical, lengthbased dynamic pattern switching and parameter adjustment logic, coupled with an adaptive scaling mechanism for the global attention range in ultra-long sequences. This approach significantly reduces the average computational complexity, approaching linear complexity, while retaining the ability to capture critical long-range dependencies.

3.3 Dynamic Routing Diffusion

To address the heterogeneous computational requirements of different stages and content segments during text generation, DrDiff incorporates a *Dynamic Routing Diffusion* mechanism based on Mixture of Experts (MoE) within the Feed-Forward Network (FFN) components of its ϵ_{θ} Transformer modules. The input hidden state h_t (from the HSA layer) is routed to one or more expert networks for processing.

Routing Decision. A routing network R_{qate} ,

typically a small MLP, computes logits $\mathbf{s}(h_t)$ for each of the M experts based on h_t . We posit that h_t , having been processed by HSA, already encodes the complexity and dependency characteristics of the current context, enabling the routing network to learn a mapping from these features to an optimal selection of experts.

$$\mathbf{s}(h_t) = [s_1(h_t), \dots, s_M(h_t)] = \underbrace{W_{gate} \cdot \operatorname{StopGradient}(h_t) + b_{gate}}_{\text{raw logits for expert selection}}$$
 DrDiff employs Top-k sparse gating (where k is

DrDiff employs Top-k sparse gating (where k is typically 1 or 2) to select experts. The gating weight $G_j(h_t)$ for expert j is:

$$G_j(h_t) = \begin{cases} \frac{\exp(s_j(h_t))}{\sum\limits_{l \in \mathsf{TopK}(\mathbf{s}(h_t))} \exp(s_l(h_t))} & \text{if expert } j \in \underbrace{\mathsf{TopK}(\mathbf{s}(h_t))}_{k \text{ experts selected based on logits}} \\ 0 & \text{otherwise} \end{cases}$$

Expert Networks & Specialization. There are M expert networks $\{E_j\}_{j=1}^M$ operating in parallel. A key innovation in DrDiff is the differentiated design of these expert networks to promote functional specialization. This includes: i) capacity differences, where some experts employ FFNs with smaller dimensions (e.g., 0.5 times the standard FFN dimension) to process simpler or repetitive text segments more economically; ii) optionally, structural differences or task-specific fine-tuning, where some experts might feature slightly different activation functions or layer normalization strategies, or even undergo preliminary, lightweight finetuning on specific sub-tasks (such as syntactic generation or content generation for particular topics) during a pre-training phase to enhance their sensitivity to certain input types. Each activated expert E_j processes the input h_t as follows:

$$O_j(h_t) = E_j(h_t) = \underbrace{W_2^{(j)} \text{ ReLU}\Big(W_1^{(j)} h_t + b_1^{(j)} \Big)}_{\text{second linear transform}} + b_2^{(j)}$$

Output Integration & Load Balancing. The final output of the MoE layer is $h'_t = \sum_{j=1}^M G_j(h_t) \cdot O_j(h_t)$. During training, an auxiliary load balancing loss $L_{aux} = \lambda_{aux} \sum_{j=1}^M f_j \cdot P_j$ (where f_j is the fraction of tokens dispatched to expert j and P_j is the routing probability for expert j) is introduced to ensure all experts are sufficiently trained. Through this mechanism, the MoE learns to dynamically allocate computational load to the most suitable experts based on input features (which serve as implicit "complexity" signals), thereby avoiding the inefficient computation characteristic of fixed structures that treat all inputs uniformly.

3.4 Semantic Anchor States Guided Optimization and Efficient Inference

To further optimize the generation path for ultralong texts, enhance global coherence, and accelerate inference, DrDiff introduces *Semantic Anchor States (SAS)*, denoted as \hat{Z}_{t_k} , for explicit guidance at specific intermediate denoising timesteps t_k (e.g., $t_k \in \{T/4, T/2, 3T/4\}$, chosen based on a conceptual division of typical text generation phases). The core innovation of SAS lies in its target definition and acquisition mechanism.

Construction of SAS Targets \hat{Z}_{t_k} . We explore methods such as the following for constructing SAS targets: SAS is based on simplified representations, where for a given Z_0 , a core semantic summary Summ (Z_0) is extracted using a pretrained, lightweight summary generation model or a topic model. \hat{Z}_{t_k} is then defined as the ideal noisy state corresponding to this summary information at timestep t_k : $\hat{Z}_{t_k} = \sqrt{\bar{\alpha}_{t_k}} \mathcal{E}(\operatorname{Summ}(Z_0)) + \sqrt{1 - \bar{\alpha}_{t_k}} \epsilon'$, where \mathcal{E} is an embedding function. This form of SAS guides the model to first generate a structure compliant with the summary, and optionally, a more complex SAS based on clustering, which involves feature extraction and clustering of noisy states $q(Z_{t_k}|Z_0)$ from a large corpus of real texts at various timesteps t_k . The centroid of each cluster can serve as an SAS target, guiding the generation process towards typical regions of the data manifold.

SAS-guided Training Objective. In addition to the standard diffusion loss $L_{diffusion}$, an SAS guidance loss L_{SAS} is incorporated:

$$L_{\text{diffusion}} = \mathbb{E}_{t, Z_0, \epsilon} \left\| \underbrace{\epsilon}_{\text{target noise}} - \underbrace{\epsilon_{\theta}(Z_t, t)}_{\text{model's predicted noise}} \right\|^2$$
 (7)

To compute L_{SAS} , we estimate the state $\tilde{Z}_{t_k}^{\theta}(Z_t)$ that the model would generate at the target timestep t_k , given the current state Z_t and the model's predicted noise $\epsilon_{\theta}(Z_t,t)$. This can be achieved by first estimating \tilde{Z}_0 through a one-step deterministic reverse process (e.g., DDIM/DDPM):

$$\tilde{Z}_0(Z_t, \epsilon_\theta) = \frac{1}{\sqrt{\bar{\alpha}_t}} \left(Z_t - \sqrt{1 - \bar{\alpha}_t} \, \epsilon_\theta(Z_t, t) \right) \quad (8)$$

and then applying the forward diffusion process to reach t_k :

$$\tilde{Z}_{t_k}^{\theta}(Z_t) = \sqrt{\bar{\alpha}_{t_k}} \, \tilde{Z}_0(Z_t, \, \epsilon_{\theta}) + \sqrt{1 - \bar{\alpha}_{t_k}} \, \epsilon' \quad (9)$$

The SAS loss is then defined as:

$$L_{SAS} = \sum_{k=1}^{K_{SAS}} \lambda_{SAS,k} \cdot || \underbrace{\tilde{Z}_{t_k}^{\theta}(Z_t)}_{\text{model's predicted state at } t_k} - \underbrace{\hat{Z}_{t_k}}_{\text{desired semantic anchor state at } t_k} ||^2$$
(10)

This explicit path guidance is a key differentiator of DrDiff from traditional diffusion models, as it constrains the generation trajectory to align with pre-defined structural information at critical intermediate junctures.

Combination with Efficient ODE/SDE Solvers. The denoising path guided by SAS is generally smoother and more "goal-oriented," reducing the stochasticity and complexity of the diffusion trajectory. This facilitates the effective operation of efficient numerical Ordinary Differential Equation (ODE) / Stochastic Differential Equation (SDE) solvers, such as DPM-Solver++ (Lu et al., 2025), which treat the discrete denoising steps as solving an ODE/SDE of the form $dZ_t = f(Z_t, t)dt + g(t)dW_t$. Specifically, a smoother path may imply better Lipschitz properties for the function $f(Z_t, t)$, or a more consistent evolution of the model's prediction ϵ_{θ} with respect to t. This allows the solver to employ larger integration steps, thereby significantly reducing the number of sampling steps $S \ll T$ while maintaining or even enhancing generation quality and coherence.

4 Experiments

The evaluation consists of four main experimental categories. (1) Natural Language Understanding (NLU): We test DrDiff on the LongBench to evaluate its ability to understand and process language in long text effectively. (2) Long-Text Generation and Question Answering: We assess DrDiff's performance on multi-hop and long-form questionanswering tasks, as well as its ability to generate coherent and structured long texts. (3) Long-Text Summarization: We evaluate DrDiff's summarization capabilities on datasets such as Arxiv and Alpaca, focusing on its ability to condense long documents while preserving key information. (4) Mixture-of-Experts (MoE) Expert Count Impact: We investigate how different numbers of MoE experts affect DrDiff's performance across datasets like Arxiv, HotpotQA, Commonsense Conversation, and QQP. Additionally, we conduct ablation studies and hyperparameter sensitivity analyses to validate the contributions of different components in DrDiff's architecture. We also perform long-text stress testing to examine the model's ability to handle extremely long sequences and evaluate the impact of different diffusion strategies on model performance (see Appendix A.4 and A.6 for

more details). All experiments are conducted on NVIDIA A100 GPUs, and the specific hyperparameter settings are in Appendix A.1.

4.1 Natural Language Understanding

Experimental Setting. To comprehensively evaluate DrDiff, particularly its ability to handle diverse tasks and long sequences, we assessed it on the LongBench (Bai et al., 2024a). LongBench encompasses a variety of challenging tasks, including single-document question answering, multidocument question answering, long in-context learning, long dialogue, code repository analysis, and long structured data processing. Several strong baseline models are included: GPT-40 (OpenAI, 2024), Qwen2.5-72B (Qwen, 2024), LLaMA-3.1-70B (LLAMA3, 2024), Longformer (Beltagy et al., 2020), Qwen2.5-7B (Qwen, 2024), LLaMA-3.1-8B (LLAMA3, 2024), and DiffuSeq (Gong et al., 2023). Given the substantial size of GPT-40, Qwen2.5-72B, and LLaMA-3.1-70B, we utilized the OpenRouter API for inference. Other models were run locally on NVIDIA A100 GPUs. We report the overall accuracy on the LongBench, as well as the performance breakdown across different difficulty levels (Easy, Hard), sequence lengths (Short, Medium, Long), and task types.

Comparisons and Analyses. The results of our evaluation on the LongBench are presented in Table 1. As shown in Table 1, DrDiff, with approximately 220M active parameters, achieves an overall score of 33.5% on LongBench, surpassing several competitive baselines including LLaMA-3.1-70B (32.1%), Longformer (31.0%), and DiffuSeq (29.5%). Notably, DrDiff demonstrates particular strengths in handling long sequences (35.6%), dialogue (38.7%), and structured data (34.6%). While it lags behind the significantly larger models like GPT-40 and Qwen2.5-72B, its competitive performance with a much smaller parameter count highlights its efficiency. Table 4 also provides a detailed breakdown of performance across different task types, where DrDiff exhibits strong performance in long dialogue (38.7%) and long structured data processing (34.6%), indicating its capability in managing complex, extended inputs in these domains. Its performance in single-document QA (31.6%), multi-document QA (32.4%), long in-context learning (32.5%), and code repository analysis (29.1%) is also competitive with the other smaller baseline models. These results collectively underscore the effectiveness and lightweight nature of the DrDiff

Table 1: Model Performance Comparison on LongBench.

Model	Overall (%)	Easy (%)	Hard (%)	Short (%)	Medium (%)	Long (%)
GPT-4o	51.9	61.4	47.1	53.9	55.2	40.7
Qwen2.5-72B	42.6	43.2	42.3	48.1	37.6	43.9
LLaMA-3.1-70B	32.1	32.8	31.7	41.6	27.9	24.6
Longformer	31.0	30.4	29.2	32.5	31.7	30.8
Qwen2.5-7B	30.5	31.2	30.1	41.1	24.7	24.6
LLaMA-3.1-8B	30.5	31.2	30.1	35.5	28.4	26.4
DiffuSeq	29.5	30.2	28.8	34.5	27.0	25.0
DrDiff (ours)	33.5	31.7	29.8	35.5	32.4	35.6

architecture for handling a diverse range of longcontext tasks.

4.2 Natural Language Generation and Question Answering

Experimental Setting. We evaluate model performance on long text generation and QA using five datasets: WikiHop (Welbl et al., 2018), TriviaQA (Joshi et al., 2017), OntoNotes, Hyperpartisan News Detection (Yang et al., 2018), and HotpotQA (Yang et al., 2018). These datasets are chosen to assess the model's ability to generate coherent, lengthy responses and to answer complex, multi-hop questions. We compare our DrDiff with nine baselines, i.e., GPT-40 (OpenAI, 2024), Qwen2.5-72B (Qwen, 2024), LLaMA-3.1-70B (LLAMA3, 2024), Longformer (Beltagy et al., 2020), LLaMA-3-8B (LLAMA3, 2024), DeepSeek-R1-Distill-Qwen-1.5B (DeepSeek-AI, 2025), GPT-2 Large (Radford et al., 2019), DiffuSeq (Gong et al., 2023), and LLaMA-2-7B (LLAMA2, 2023), each of which is fine-tuned on respective training splits using standard protocols. These models include state-of-the-art and open-source models of various sizes, and models similar to ours that specifically address long text problems. All experiments run under controlled hardware with fixed random seeds for reproducibility, reporting the highest test-set scores per model. Comparisons and Analyses. As shown in Table 2, DrDiff demonstrates strong performance in natural language generation and question answering tasks. While state-of-the-art models like GPT-40 and Qwen2.5-72B lead in overall scores, DrDiff achieves competitive results, scoring 76.2% accuracy on WikiHop, 82.1% F1 on TriviaQA, 82.4% F1 on OntoNotes, 95.0% F1 on Hyperpartisan, and 68.0% Joint F1 on HotpotQA. Notably, these results place DrDiff ahead of other significant models such as LLaMA-3.1-70B (77.5% average) and Longformer (77.0% average). For instance, Longformer achieved 74.6% on WikiHop, 74.1% on TriviaQA, and 78.4% on OntoNotes. Furthermore, DrDiff outperforms other models like LLaMA-3-8B, DeepSeek-R1-Distill-Qwen-1.5B, LLaMA-2-7B, DiffuSeq, and GPT-2 Large across these datasets.

4.3 Number of Experts in MoE

Experimental Setting. This experiment aims to fairly evaluate the inference performance of the Mixture-of-Experts (MoE) (Shazeer et al., 2017) compared to single-model baselines. We use the Arxiv dataset and measure performance using BLEU, ROUGE-L, and BERTScore to assess text quality. For a balanced comparison, the singlemodel baseline has 160M parameters. Our MoE structure consists of smaller expert models, each with 20M parameters. We design three MoE configurations with varying total numbers of experts: 2 experts (resulting in 40M total parameters), 4 experts (80M total parameters), and 8 experts (160M total parameters). It is crucial to distinguish between "total parameters" (which scale with the number of experts in a configuration) and "active parameters" (the parameters engaged during a single inference pass). To maintain consistency and fairness in computational cost across these configurations, we activate a fixed number of 2 experts (thus, 2 experts \times 20M parameters/expert = 40M active parameters) per inference. This ensures that while the model's capacity (total parameters) varies, the computational load per inference (active parameters) remains comparable. Beyond performance metrics, we also record inference time and trainable parameters to analyze computational efficiency and scalability. These additional measurements help assess the trade-offs between overall model size (total parameters), inference speed, and output quality.

Comparisons and Ablation Studies. Figure 3

Table 2: Performance on natural language generation and question answering tasks. Evaluation is conducted on five datasets: WikiHop (multi-hop QA), TriviaQA (open QA), OntoNotes (coreference resolution), Hyperpartisan News Detection (long document classification), and HotpotQA (distractor setting, joint F1). The scores (accuracy or F1, as appropriate) are reported with standard deviations, and the final column shows the average across all tasks.

Method	WikiHop (Acc.)	TriviaQA (F1)	OntoNotes (F1)	Hyperpartisan (F1)	HotpotQA (Joint F1)	Avg Score
GPT-4o	$77.2_{\pm 0.6}$	$82.9_{\pm 0.5}$	83.3 _{±0.6}	95.6 _{±0.5}	69.2 _{±0.7}	81.6 _{±0.6}
Qwen2.5-72B	$76.8_{\pm 0.7}$	82.5 _{±0.6}	$82.9_{\pm 0.7}$	$95.3_{\pm 0.6}$	$68.8_{\pm0.8}$	81.3 _{±0.7}
LLaMA-3.1-70B	$73.0_{\pm 0.7}$	$78.5_{\pm 0.6}$	$79.0_{\pm 0.7}$	$92.0_{\pm 0.6}$	$65.0_{\pm0.8}$	$77.5_{\pm 0.7}$
Longformer	$74.6_{\pm 0.7}$	$74.1_{\pm 0.9}$	$78.4_{\pm 0.7}$	$93.8_{\pm 0.6}$	$63.9_{\pm 0.7}$	$77.0_{\pm 0.7}$
LLaMA-3-8B	$69.8_{\pm 0.7}$	$69.5_{\pm 0.5}$	$74.1_{\pm 0.6}$	$84.7_{\pm 0.7}$	$59.8_{\pm0.8}$	$71.6_{\pm 0.6}$
DeepSeek-R1-Distill-Qwen-1.5B	$67.5_{\pm 0.6}$	$71.4_{\pm 0.7}$	$73.2_{\pm 0.5}$	$79.8_{\pm 0.6}$	$61.3_{\pm 0.6}$	$70.6_{\pm 0.6}$
LLaMA-2-7B	$59.2_{\pm 0.7}$	$64.8_{\pm 0.6}$	$69.5_{\pm 0.7}$	$77.4_{\pm 0.6}$	$54.2_{\pm 0.7}$	$65.0_{\pm 0.6}$
DiffuSeq	$64.1_{\pm 0.6}$	$54.3_{\pm 0.7}$	$61.2_{\pm 0.6}$	$84.3_{\pm 0.7}$	$49.6_{\pm0.6}$	$62.7_{\pm 0.6}$
GPT-2 Large	$49.8_{\pm 0.5}$	$39.6_{\pm 0.7}$	$48.3_{\pm0.8}$	$69.5_{\pm 0.6}$	$34.7_{\pm 0.5}$	$48.4_{\pm 0.6}$
DrDiff (ours)	$76.2_{\pm 0.8}$	82.1 _{±0.6}	$82.4_{\pm 0.7}$	$95.0_{\pm 0.7}$	$68.0_{\pm 0.8}$	$78.6_{\pm 0.7}$

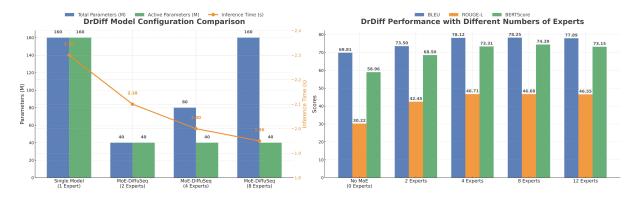


Figure 3: The figure above shows the changes in model training parameters and inference time after setting different experts. The figure below shows the changes in score indicators on the Arxiv dataset after setting different numbers of experts.

illustrates trends in model configuration and performance. The upper chart shows that as the total number of available experts increases from 1 to 8 in DrDiff (leading to an increase in total model parameters), the active parameters per inference remain constant at 40M, due to our strategy of activating a fixed K=2 experts. Concurrently, inference time drops from 2.30s to 1.95s. The lower chart indicates that performance metrics (BLEU, ROUGE-L, BERTScore) improve with more total experts up to 8 but decline slightly at 12, suggesting diminishing returns on increasing model capacity beyond a certain point. DrDiff significantly reduces active parameters and inference time compared to dense models of similar total capacity, while maintaining competitive performance, demonstrating that the MoE framework efficiently allocates computational resources. The 8-expert model (with 160M total parameters, but still 40M active parameters per inference) achieves the best balance between efficiency and quality. To further analyze DrDiff's components, we conducted an ablation study by

modifying sparse attention, diffusion steps, and attention window sizes in Table 5. The baseline model integrates sparse attention with DiffuSeq. Detailed descriptions of these experiments are provided in Appendix A.3 and A.8.

5 Conclusion

This paper introduces DrDiff, a novel framework that addresses long-text generation challenges through dynamic routing diffusion with Hierarchical Attention. Our results show that DrDiff outperforms existing methods in both computational efficiency and generation quality. The dynamic expert scheduling mechanism reduces computational complexity from $O(n^2)$ to O(n) while preserving text coherence, enabling more efficient long-sequence processing. Additionally, Hierarchical Sparse Attention effectively handles sequences up to and beyond 8K tokens, ensuring robustness across tasks. A key innovation is semantic anchor guidance, which optimizes diffusion by accelerating generation time without compromising quality. This tech-

nique balances efficiency and fidelity, making the model well-suited for real-world long-text applications. DrDiff is a promising solution for scientific writing, creative generation, and summarization.

Future work will extend sequence length, enhance adaptability, and refine domain-specific coherence. Moreover, we plan to deepen the theoretical analysis of dynamic routing and diffusion processes, and explore multi-modal extensions (e.g., integrating visual or structured data) to further broaden DrDiff's applicability. We also conduct large-scale user studies to validate the framework in practical real-world scenarios and assess its performance across diverse application contexts.

Limitations

Despite the promising results achieved by DrDiff, several important limitations deserve attention and point to directions for future research:

Limited Exploration of Extreme-Length Texts: While our framework shows improvements in handling sequences up to 8K tokens, the exploration of even longer texts (e.g., >20K tokens) remains limited. The current architecture's effectiveness for extreme-length document generation needs further investigation, as the complexity of semantic dependencies and memory constraints may introduce unforeseen challenges at such scales.

Theoretical Foundation: Although empirical results demonstrate the effectiveness of our dynamic routing mechanism, the theoretical underpinnings of why this approach works well for long-text generation lack rigorous mathematical proof. Specifically, the convergence properties of the diffusion process under dynamic expert scheduling and the optimal balance between different attention patterns require more thorough theoretical analysis.

Limited Interpretability: The current version of DrDiff operates largely as a black box, particularly in its expert scheduling decisions. The lack of interpretability in how the model allocates computational resources and switches between different attention patterns makes it challenging to diagnose potential failure cases or optimize the model for specific applications.

Resource Consumption Trade-offs: While our approach reduces computational complexity, the multi-expert architecture introduces additional memory overhead. The balance between computational efficiency and memory usage requires further optimization, especially for deployment in

resource-constrained environments.

Domain Adaptability: Our evaluation primarily focused on general-domain text generation. The framework's effectiveness in specialized domains (e.g., scientific papers, legal documents) where strict formatting or domain-specific knowledge is required remains to be thoroughly validated.

These limitations highlight several promising directions for future research, including developing more robust theoretical frameworks for long-text generation, improving model interpretability, and exploring more efficient architectures for extremelength text generation. Addressing these challenges will be crucial for advancing the field of long-text generation and expanding its practical applications.

Acknowledgments

This work was supported in part by the National Natural Science Foundation of China (NSFC) under Grant 62276283, in part by the China Meteorological Administration's Science and Technology Project under Grant CMAJBGS202517, in part by Guangdong Basic and Applied Basic Research Foundation under Grant 2023A1515012985, in part by Guangdong-Hong Kong-Macao Greater Bay Area Meteorological Technology Collaborative Research Project under Grant GHMA2024Z04, in part by Fundamental Research Funds for the Central Universities, Sun Yat-sen University under Grant 23hytd006, and in part by Guangdong Provincial High-Level Young Talent Program under Grant RL2024-151-2-11.

References

Saleh Ashkboos, Iman Mirzadeh, Keivan Alizadeh, Mohammad Hossein Sekhavat, Moin Nabi, Mehrdad Farajtabar, and Fartash Faghri. 2024. Computational bottlenecks of training small-scale large language models. *Preprint*, arXiv:2410.19456.

Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu, Jiankai Tang, Zhidian Huang, Zhengxiao Du, Xiao Liu, Aohan Zeng, Lei Hou, Yuxiao Dong, Jie Tang, and Juanzi Li. 2024a. LongBench: A bilingual, multitask benchmark for long context understanding. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3119–3137, Bangkok, Thailand. Association for Computational Linguistics.

Yushi Bai, Jiajie Zhang, Xin Lv, Linzhi Zheng, Siqi Zhu, Lei Hou, Yuxiao Dong, Jie Tang, and Juanzi Li. 2024b. Longwriter: Unleashing 10,000+ word generation from long context llms. *Preprint*, arXiv:2408.07055.

- Yushi Bai, Jiajie Zhang, Xin Lv, Linzhi Zheng, Siqi Zhu, Lei Hou, Yuxiao Dong, Jie Tang, and Juanzi Li. 2024c. Longwriter: Unleashing 10,000+ word generation from long context llms. *Preprint*, arXiv:2408.07055.
- Philipp Becker, Abhinav Mehrotra, Ruchika Chavhan, Malcolm Chadwick, Luca Morreale, Mehdi Noroozi, Alberto Gil Ramos, and Sourav Bhattacharya. 2025. Edit: Efficient diffusion transformers with linear compressed attention. *Preprint*, arXiv:2503.16726.
- Iz Beltagy, Matthew E. Peters, and Arman Cohan. 2020. Longformer: The long-document transformer. *Preprint*, arXiv:2004.05150.
- Amanda Bertsch, Uri Alon, Graham Neubig, and Matthew R. Gormley. 2023. Unlimiformer: long-range transformers with unlimited length input. In *Proceedings of the 37th International Conference on Neural Information Processing Systems*, NIPS '23, Red Hook, NY, USA. Curran Associates Inc.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020a. Language models are few-shot learners. *Preprint*, arXiv:2005.14165.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020b. Language models are few-shot learners. *Preprint*, arXiv:2005.14165.
- Jon Chamberlain, Massimo Poesio, and Udo Kruschwitz. 2008. Addressing the resource bottleneck to create large-scale annotated texts. In *Semantics in Text Processing. STEP 2008 Conference Proceedings*, pages 375–380. College Publications.
- Howard Chen, Ramakanth Pasunuru, Jason Weston, and Asli Celikyilmaz. 2023. Walking down the memory maze: Beyond context limit through interactive reading. *Preprint*, arXiv:2310.05029.
- Yafeng Chen, Zhihao Lan, Jensen Li, and Jie Zhu. 2021. Topologically protected second harmonic generation via doubly resonant high-order photonic modes. *Physical Review B*, 104(15).

- Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. 2019a. Generating long sequences with sparse transformers. *URL https://openai.com/blog/sparse-transformers*.
- Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. 2019b. Generating long sequences with sparse transformers. *Preprint*, arXiv:1904.10509.
- Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. 2019c. Generating long sequences with sparse transformers. *Preprint*, arXiv:1904.10509.
- Kevin Clark, Minh-Thang Luong, Quoc V. Le, and Christopher D. Manning. 2020. Electra: Pre-training text encoders as discriminators rather than generators. *Preprint*, arXiv:2003.10555.
- Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell, Quoc V. Le, and Ruslan Salakhutdinov. 2019. Transformer-xl: Attentive language models beyond a fixed-length context. *Preprint*, arXiv:1901.02860.
- DeepSeek-AI. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *Preprint*, arXiv:2501.12948.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. *Preprint*, arXiv:1810.04805.
- Yizhao Gao, Zhichen Zeng, Dayou Du, Shijie Cao, Hayden Kwok-Hay So, Ting Cao, Fan Yang, and Mao Yang. 2024. Seerattention: Learning intrinsic sparse attention in your llms. *arXiv preprint arXiv:2410.13276*.
- Shansan Gong, Mukai Li, Jiangtao Feng, Zhiyong Wu, and Lingpeng Kong. 2023. Diffuseq: Sequence to sequence text generation with diffusion models. *Preprint*, arXiv:2210.08933.
- Jian Guan, Xiaoxi Mao, Changjie Fan, Zitao Liu, Wenbiao Ding, and Minlie Huang. 2021. Long text generation by modeling sentence-level and discourse-level coherence. *Preprint*, arXiv:2105.08963.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising diffusion probabilistic models. *Preprint*, arXiv:2006.11239.
- Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke Zettlemoyer. 2017. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension. In *ACL*.
- Nikita Kitaev, Lukasz Kaiser, and Anselm Levskaya. 2020. Reformer: The efficient transformer. In *International Conference on Learning Representations*.
- Kalpesh Krishna, Erin Bransom, Bailey Kuehl, Mohit Iyyer, Pradeep Dasigi, Arman Cohan, and Kyle Lo. 2023. Longeval: Guidelines for human evaluation of faithfulness in long-form summarization. *Preprint*, arXiv:2301.13298.

- Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2020. Albert: A lite bert for self-supervised learning of language representations. *Preprint*, arXiv:1909.11942.
- Zhenghao Lin, Yeyun Gong, Yelong Shen, Tong Wu, Zhihao Fan, Chen Lin, Nan Duan, and Weizhu Chen. 2023. Text generation with diffusion language models: A pre-training approach with continuous paragraph denoise. *Preprint*, arXiv:2212.11685.
- Xiang Liu, Peijie Dong, Xuming Hu, and Xiaowen Chu. 2024. Longgenbench: Long-context generation benchmark. *Preprint*, arXiv:2410.04199.
- LLAMA2. 2023. Llama 2: Open foundation and fine-tuned chat models. *Preprint*, arXiv:2307.09288.
- LLAMA3. 2024. The llama 3 herd of models. *Preprint*, arXiv:2407.21783.
- Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. 2025. Dpm-solver++: Fast solver for guided sampling of diffusion probabilistic models. *Machine Intelligence Research*, pages 1–22.
- Mohaddeseh Mirbeygi and Hamid Beigy. 2025. Prompt guided diffusion for controllable text generation. In *Proceedings of the Tenth Workshop on Noisy and User-generated Text*, pages 78–84, Albuquerque, New Mexico, USA. Association for Computational Linguistics.
- Alex Nichol and Prafulla Dhariwal. 2021. Improved denoising diffusion probabilistic models. *Preprint*, arXiv:2102.09672.
- OpenAI. 2024. Gpt-4 technical report. *Preprint*, arXiv:2303.08774.
- Qwen. 2024. Qwen2 technical report. arXiv preprint arXiv:2407.10671.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2023. Exploring the limits of transfer learning with a unified text-to-text transformer. *Preprint*, arXiv:1910.10683.
- Mohammad Sadegh Rasooli and Joel R. Tetreault. 2015. Yara parser: A fast and accurate dependency parser. Computing Research Repository, arXiv:1503.06733. Version 2.
- Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton, and Jeff Dean. 2017. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. *Preprint*, arXiv:1701.06538.

- Jiaming Tang, Yilong Zhao, Kan Zhu, Guangxuan Xiao, Baris Kasikci, and Song Han. 2024. Quest: queryaware sparsity for efficient long-context llm inference. ICML'24. JMLR.org.
- Zecheng Tang, Pinzheng Wang, Keyan Zhou, Juntao Li, Ziqiang Cao, and Min Zhang. 2023. Can diffusion model achieve better performance in text generation? bridging the gap between training and inference! *Preprint*, arXiv:2305.04465.
- Yi Tay, Mostafa Dehghani, Dara Bahri, and Donald Metzler. 2022a. Efficient transformers: A survey. *Preprint*, arXiv:2009.06732.
- Yi Tay, Mostafa Dehghani, Dara Bahri, and Donald Metzler. 2022b. Efficient transformers: A survey. *Preprint*, arXiv:2009.06732.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. Llama: Open and efficient foundation language models. *Preprint*, arXiv:2302.13971.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2023. Attention is all you need. *Preprint*, arXiv:1706.03762.
- Sinong Wang, Belinda Z. Li, Madian Khabsa, Han Fang, and Hao Ma. 2020. Linformer: Self-attention with linear complexity. *Preprint*, arXiv:2006.04768.
- Johannes Welbl, Pontus Stenetorp, and Sebastian Riedel. 2018. Constructing datasets for multi-hop reading comprehension across documents. *Transactions of the Association for Computational Linguistics*.
- Tong Wu, Junzhe Shen, Zixia Jia, Yuxuan Wang, and Zilong Zheng. 2025. From hours to minutes: Lossless acceleration of ultra long sequence generation up to 100k tokens. *Preprint*, arXiv:2502.18890.
- Minkai Xu, Tomas Geffner, Karsten Kreis, Weili Nie, Yilun Xu, Jure Leskovec, Stefano Ermon, and Arash Vahdat. 2025. Energy-based diffusion language models for text generation. *Preprint*, arXiv:2410.21357.
- Zhiyu Yang, Wenhan Xiong, Mo Yu, and William Yang Wang. 2018. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. In *EMNLP*.
- Hongyi Yuan, Zheng Yuan, Chuanqi Tan, Fei Huang, and Songfang Huang. 2023. Seqdiffuseq: Text diffusion with encoder-decoder transformers. *Preprint*, arXiv:2212.10325.
- Jusheng Zhang, Yijia Fan, Kaitong Cai, and Keze Wang. 2025a. Kolmogorov-arnold fourier networks. *Preprint*, arXiv:2502.06018.

Jusheng Zhang, Zimeng Huang, Yijia Fan, Ningyuan Liu, Mingyan Li, Zhuojie Yang, Jiawei Yao, Jian Wang, and Keze Wang. 2025b. Kabb: Knowledge-aware bayesian bandits for dynamic expert coordination in multi-agent systems. *Preprint*, arXiv:2502.07350.

A Appendix

Contents

A 1	TT CONT	10
A.1	Hyperparameter Settings	12
A.2	Hierarchical Attention Mode	
	Switching Strategy, HAMSS	12
	A.2.1 Dense Mode (\mathcal{M}_D)	13
	A.2.2 4K Mode (\mathcal{M}_{4K})	13
	A.2.3 8K Mode (\mathcal{M}_{8K})	13
	A.2.4 16K+ Mode (\mathcal{M}_{16K+})	13
	A.2.5 Mode Switching Decision	
	Function	14
	A.2.6 Parameter Selection Ratio-	
	nale and Tuning for HSA	
	Modes	14
	A.2.7 Mode Switching Decision	
	Function F_{ϕ} Explained	15
A.3	Ablation Study	15
A.4	Long Text Stress Resistance Exper-	
	iment	15
	A.4.1 Experimental Setup	15
	A.4.2 Experimental Results	16
A.5	Computational Complexity Analysis	16
	A.5.1 Computational Complex-	
	ity Comparison	16
	A.5.2 Computational Resource	
	Consumption Comparison	16
A.6	Impact of Different Settings of Dif-	
	fusion	17
	A.6.1 Experimental setup	17
	A.6.2 Experimental Results	17
A.7	Semantic Anchor State Guidance	1,
11.7	(SAS) Weight Tuning and Impact.	18
A.8	Extended Ablation Study on HSA	10
Λ.0	and DES	19
A.9	Training Datasets	20
17. 7	Training Datasets	20

A.1 Hyperparameter Settings

The unified hyperparameter settings for all experiments are shown in Table 3. Note that for the batch size and epoch setting, we follow the official test configurations for different datasets.

A.2 Hierarchical Attention Mode Switching Strategy, HAMSS

To achieve an adaptive balance between computational efficiency and modeling capability for sequence length, **DrDiff** employs the **Hierarchical Attention Mode Switching Strategy** (**HAMSS**). The core idea of HAMSS is to dynamically select the optimal attention mode based on the input

Hyperparameter	Value
Transformer Layers	12
Attention Heads per Layer	12
Sparse Attention Mechanism	Hierarchical Sparse Attention
MoE Strategy	Multiple Expert Networks per Layer, Dynamic Selection via Gating Mechanism
Diffusion Steps	2048, Square Root Noise Scheduling
DPM-Solver++ Integration	Integrated to Reduce Diffusion Steps and Improve Generation Speed
Batch Size	Determined by different datasets
Epoch	Determined by different datasets
Learning Rate	1e-4
Warm-up Steps	5000
Weight Decay	0.01
Optimizer	AdamW
Max Sequence Length	512
Gradient Clipping	1.0
Dropout Rate	0.1
Expert Size per Expert in MoE	20M
Experts per Layer	2, 4, or 8 (depending on experiment configuration)
Framework	PyTorch

Table 3: Hyperparameter Settings for All Experiments

sequence length, ensuring high-quality text generation while maintaining efficient computation.

HAMSS consists of four attention modes:

$$\mathcal{H} = (\mathcal{M}_D, \mathcal{M}_{4K}, \mathcal{M}_{8K}, \mathcal{M}_{16K+})$$

Each mode \mathcal{M}_* is composed of five elements:

$$\mathcal{M}_* = \langle \Psi_*, \Omega_*, \mathcal{C}_*, \mathcal{G}_*, \mathcal{T}_* \rangle$$

A.2.1 Dense Mode (\mathcal{M}_D)

Activation Condition: Input sequence length $n \leq$ 512. Neighborhood Definition:

$$\mathcal{N}_D(i) = \{1, 2, \dots, n\}, \quad \forall i \in \{1, \dots, n\}$$

The attention weight matrix $\mathbf{A}_D \in \mathbb{R}^{n \times n}$ satisfies:

$$A_D^{(i,j)} = \operatorname{Softmax}\left(\frac{Q_i K_j^\top}{\sqrt{d}}\right), \quad \forall (i,j)$$

Computational Complexity:

$$C_D(n) = \underbrace{O(n^2)}_{\text{attention computation}} + \underbrace{O(n^2d)}_{\text{storage overhead}}$$

A.2.2 4K Mode (\mathcal{M}_{4K})

Activation Condition: $512 < n \le 4,096$. Structural Parameters: Local window width 256, global token count m**Hybrid Neighborhood Definition:** $\mathcal{N}_{4K}(i) = \{j \mid \max(1, i - w_{4K}) \le j \le \min(n, i + w_{4K})\} \cup \mathcal{T}_{global}$

Attention Weight Calculation:

$$A_{4K}^{(i,j)} = \begin{cases} \frac{\exp(Q_i K_j^\top/\sqrt{d})}{\sum_{k \in \mathcal{N}_{4K}(i)} \exp(Q_i K_k^\top/\sqrt{d})}, & j \in \mathcal{N}_{4K}(i) \\ 0, & \text{otherwise} \end{cases}$$

Computational Complexity:

$$C_{4K}(n) = \underbrace{O(n \cdot (2w_{4K} + m))}_{\text{dynamic sparse computation}} + \underbrace{O(nmd)}_{\text{cross-window communication}}$$

When
$$m = \Theta(\sqrt{n})$$
, $C_{4K}(n) = O(n^{3/2})$.

A.2.3 8K Mode (\mathcal{M}_{8K})

Activation Condition: $4,096 < n \le 8,192$.

Topological Parameters: Dilated window width $w_{8K}=512$, stride s=4, hierarchical dilation factor $d_l = 2^{\lfloor l/L \rfloor}$.

Dilated Neighborhood **Generation:** $\mathcal{N}_{8K}^{(l)}(i) = \left\{i + k \cdot s \cdot d_l \mid k \in \mathbb{Z}, \ |k| \leq \frac{w_{8K}}{2sd_l} \right\} \cap \{1, \dots, n\}$ Fractal Attention Weights:

$$A_{8K}^{(l,i,j)} = \frac{\exp\left(\frac{Q_i^{(l)}K_j^{(l)\top}}{\sqrt{d}}\right)}{\sum\limits_{l \in \mathcal{N}^{(l)}(i)} \exp\left(\frac{Q_i^{(l)}K_k^{(l)\top}}{\sqrt{d}}\right)} \cdot \mathbb{I}\left(j \in \mathcal{N}_{8K}^{(l)}(i)\right)$$

Complexity Analysis:

$$C_{8K}(n) = \sum_{l=1}^{L} O\left(\frac{nw_{8K}}{sd_l}\right) = O(n\log n)$$

A.2.4 16K+ Mode (\mathcal{M}_{16K+})

Activation Condition: n > 8,192.

Extreme Parameters: Super-window width $w_{16K} = 1024$, meta-stride $s_{meta} = 8$, key token ratio $\rho = 0.05$.

$$\mathcal{N}_{16K+}(i) = \underbrace{\begin{cases} i + ks_{meta} \mid k = -\frac{w_{16K}}{2s_{meta}}, \dots, \frac{w_{16K}}{2s_{meta}} \end{cases}}_{\text{Sparse Local}}$$

$$\cup \underbrace{\begin{cases} \left \lfloor n \cdot \frac{m}{M} \right \rfloor \mid m = 1, \dots, \lceil \rho n \rceil \end{cases}}_{\text{Semantic Anchors}}$$

Hybrid Attention Weights:

$$A_{16K+}^{(i,j)} = \alpha \cdot A_{\text{local}}^{(i,j)} + (1 - \alpha) \cdot A_{\text{global}}^{(i,j)}$$

where $\alpha = \sigma(\beta \cdot (j-i))$ is a position-aware mixing coefficient, and β is a learnable parameter. **Linear Complexity Proof**:

$$\mathcal{C}_{16K+}(n) = O\left(n \cdot \left(\frac{w_{16K}}{s_{meta}} + \rho n\right)\right)$$
$$= O(n) \quad \text{(when } \rho = O(1/n)\text{)}$$

A.2.5 Mode Switching Decision Function

Mode switching is achieved via the decision network $\mathcal{F}_{\phi}: \mathbb{R}^d \to [0, 1]^4$:

$$[\pi_D, \pi_{4K}, \pi_{8K}, \pi_{16K+}]^\top = \operatorname{Softmax}(\mathcal{F}_{\phi}(\bar{h}))$$

where $\bar{h} = \frac{1}{n} \sum_{i=1}^n h_i$ is the average sequence feature. The final active mode is:

$$\mathcal{M}_{\text{active}} = \arg \max_{*} \left(\pi_* \cdot \mathbb{I}(n \in \Omega_*) \right)$$

This strategy ensures the model achieves an optimal trade-off on the Pareto frontier, satisfying:

$$\forall n, \; \exists \mathcal{M}_* \in \mathcal{H}, \; \frac{\mathcal{C}_*(n)}{\mathcal{C}_D(n)} \leq \epsilon(n),$$
$$\operatorname{Perf}(\mathcal{M}_*) \geq \gamma \cdot \operatorname{Perf}(\mathcal{M}_D).$$

A.2.6 Parameter Selection Rationale and Tuning for HSA Modes

The Hierarchical Sparse Attention (HSA) mechanism is designed to dynamically adapt the attention strategy based on input sequence length n, balancing computational efficiency with modeling capability. The specific parameters for each attention mode—including length thresholds (N_1, N_2, N_3) , window sizes (w), dilation rates (δ) , and the configuration of global nodes—were determined through a combination of factors, as outlined below:

Length Thresholds $(N_1 = 512, N_2 = 4096, N_3 = 8192)$: These thresholds were established by analyzing typical sequence length distributions encountered in our target long-text generation benchmarks (e.g., Arxiv). $N_1 = 512$ was chosen as a common maximum length for standard dense attention models, beyond which quadratic complexity becomes prohibitive. The transitions at $N_2 = 4$ K and $N_3 = 8$ K correspond to points where different sparse attention strategies offer demonstrably better trade-offs. These were refined

based on preliminary experiments observing performance shifts and computational costs when applying simpler fixed sparse patterns to these length categories. The aim was to align mode switches with points where a more specialized attention pattern (e.g., incorporating more global attention for 4K - 8K, or a more aggressive sparse pattern for 16K+) becomes beneficial. Mode-Specific Parameters (e.g., $w_{4K}, w_{8K}, w_{16K+}, m, s, \rho$): General Principle: The parameters within each mode (detailed in Appendix A.2.1-A.2.4) were selected to optimize the balance between capturing sufficient contextual information (local, dilated, global) and adhering to a near-linear computational complexity budget for that length category. Window Sizes (w): Local window sizes (e.g., $w_{4K} = 256, w_{8K} =$ $512, w_{16K+} = 1024$) were chosen to be large enough to capture meaningful local dependencies relevant to text generation tasks. These were informed by common practices in prior work on sparse attention (e.g., Longformer, BigBird) and adjusted through pilot experiments on validation sets to ensure good performance without excessive computation. Dilation Rates (δ) and Strides (s): For modes like \mathcal{M}_{8K} , dilation rates and strides were configured to expand the receptive field efficiently, allowing the model to attend to more distant tokens without incurring the cost of a fully dense or very large local window. The hierarchical dilation factor $(d_l = 2^{\lfloor l/L \rfloor})$ ensures multi-scale context aggregation. Global Node Configuration: The number of global tokens in $\mathcal{M}4K$ $(m = \lceil \sqrt{n} \rceil)$ and the key token ratio in $\mathcal{M}16K + (\rho = 0.05)$ were designed to provide essential global context. The \sqrt{n} scaling offers a compromise for medium-length sequences, while a small fixed ratio for ultra-long sequences ensures scalability. The selection of which tokens become global (e.g., regularly spaced, learnable importance) was based on simplicity and effectiveness observed in initial trials. Preliminary Tuning: While an exhaustive grid search over all HSA parameters would be computationally prohibitive, the final parameter values were arrived at through an iterative process. This involved setting initial values based on literature and theoretical considerations, followed by a series of preliminary experiments on a subset of the data or tasks to observe their impact on perplexity, generation quality, and computational throughput. Adjustments were then made to refine the balance for the overall DrDiff framework. The goal was to find a robust set of parameters that generalizes well across the targeted long-text sce-

Table 4: Task Type Breakdown Performance on LongBench. This table shows the performance of different models across various task types, including Single-Document QA, Multi-Document QA, Long ICL, Long Dialogue, Code Repo, Long Structured tasks, along with their average performance.

Model	Single-Doc QA	Multi-Doc QA	Long ICL	Long Dialogue	Code Repo	Long Structured	Average
GPT-40	65.3	63.8	58.9	62.4	60.1	58.5	51.9
Qwen2.5-72B	44.8	43.7	42.0	43.0	41.3	40.7	42.6
DrDiff (ours)	31.6	32.4	32.5	38.7	29.1	34.6	33.5
LLaMA-3.1-70B	34.0	32.9	32.3	32.5	30.2	30.6	32.1
Longformer	31.3	31.0	30.4	36.0	28.5	27.9	30.9
Qwen2.5-7B	32.0	30.5	29.0	32.5	31.5	28.5	30.7
LLaMA-3.1-8B	31.5	29.0	28.5	32.0	31.0	30.5	30.4
DiffuSeq	31.2	28.2	28.8	29.5	27.0	28.8	28.9

narios rather than fine-tuning for a single specific dataset. The parameters detailed in Appendix A.2 represent the outcome of this design and preliminary tuning process, aimed at achieving a practical and effective hierarchical attention strategy.

A.2.7 Mode Switching Decision Function F_{ϕ} Explained

As introduced in Appendix A.2.5, the mode switching in the Hierarchical Sparse Attention (HSA) mechanism is guided by the sequence length n and facilitated by a decision network F_{ϕ} .

Role and Input/Output: The decision network F_{ϕ} takes the average of the token hidden states for the input sequence, $\overline{h} = \frac{1}{n} \sum_{i=1}^n h_i$, as its input. This average hidden state $\overline{h} \in \mathbb{R}^d$ serves as a condensed representation of the overall characteristics of the current sequence. F_{ϕ} outputs a probability distribution $[\pi_D, \pi_{4K}, \pi_{8K}, \pi_{16K+}]^{\top}$ over the four predefined attention modes via a Softmax function. Each π represents the learned preference for mode M

Architecture of F_{ϕ} : The specific architecture of F_{ϕ} is designed to be lightweight to minimize overhead. In our implementation, F_{ϕ} is a small Multi-Layer Perceptron (MLP). For instance, it can consist of one or two fully connected hidden layers with a non-linear activation function (e.g., ReLU), followed by the final linear layer that produces logits for the Softmax function. The exact dimensions of these hidden layers are kept small (e.g., a fraction of the main model's hidden dimension d) to ensure efficiency. Training and Interaction with **Length Thresholds:** The parameters of the F_{ϕ} network are trained end-to-end as part of the overall DrDiff model optimization. This allows F_{ϕ} to learn a mapping from sequence characteristics (h) to appropriate attention modes, guided by the main task loss (e.g., the diffusion model's denoising objective). It is important to note the interplay between the learned probabilities π_* from F_ϕ and the predefined length-based activation conditions $I(n \in \Omega_*)$ (where Ω_* is the valid length range for mode \mathcal{M}_*). The final active mode is selected as $M_{\text{active}} = \arg\max_{M_* \in \mathcal{H}} (\pi_* \cdot I(n \in \Omega_*)).$ In practice, for sequences falling squarely within a predefined length bracket Ω_* , the $I(n \in \Omega_*)$ term often plays a decisive role, ensuring the mode designed for that length is chosen. The learned component F_{ϕ} can be particularly influential for sequences near the boundaries of these length thresholds, potentially learning to enable smoother transitions or making more nuanced choices if sequence characteristics (captured by \overline{h}) suggest a deviation from the default length-based rule. However, the primary driver for mode selection remains the explicitly defined length ranges, with F_{ϕ} offering a mechanism for learned refinement within this framework. The objective of this design is to combine the robustness of rule-based length thresholds with the potential adaptability of a learned decision mechanism, ensuring that HSA selects an appropriate and efficient attention pattern for any given input sequence.

A.3 Ablation Study

As shown in Table 5, the study includes various configurations such as removing sparse attention, altering the diffusion steps, and changing the attention window sizes to evaluate their impact on performance metrics like BLEU, ROUGE, and BERTScore.

A.4 Long Text Stress Resistance Experiment

A.4.1 Experimental Setup

This experiment evaluates the performance and stability of DrDiff in generating and summarizing long texts ranging from 8K to 30K tokens. The dataset includes 1,000 samples per length category

Configuration	Attention	Diffusion	Window	BLEU/ROUGE/BERTScore	
Configuration	Type Steps		Size	BLEU/ROUGE/BERTScore	
Baseline (Full Model)	Sparse	2048	512	75.41/58.96/71.89	
No Sparse Attention	Standard	2048	512	72.52/56.68/68.41	
Reduced Diffusion Steps	Sparse	1024	512	73.11/56.97/69.26	
Increased Diffusion Steps	Sparse	4096	512	74.71/57.32/70.20	
Smaller Window Size	Sparse	2048	256	73.80/57.79/69.65	
Larger Window Size	Sparse	2048	1024	74.40/58.52/69.92	

Table 5: Ablation study results comparing different configurations of the DrDiff model on the Arxiv dataset.

(15K, 30K, tokens) from three sources: Project Gutenberg (public domain e-books), PubMed Central (biomedical papers), and Wikipedia Long Articles. Texts are preprocessed by removing special characters, HTML tags, and incomplete sentences, then truncated or split to fit the model's maximum input size (5,000 tokens). The experimental tasks focus on long text generation (maintaining coherence and logical flow) and summarization (compressing information while preserving semantics). Evaluation metrics include ROUGE scores (n-gram overlap), BERTScore (semantic similarity), and perplexity (model adaptability).

A.4.2 Experimental Results

Description of Experimental Results The experimental results for DrDiff on long text generation and summarization are presented in Table 4. The model was evaluated across four text lengths: 8,000, 16,000, 24,000, and 30,000 tokens. The metrics used include ROUGE-1, ROUGE-2, ROUGE-L, BERTScore, and Perplexity. The results show that the model's performance in terms of ROUGE and BERTScore decreases as the text length increases, while Perplexity exhibits a downward trend, indicating improved adaptability to longer texts. Analysis of Experimental Results A key observation from the results is that 16,000 tokens act as a critical threshold for DrDiff. Specifically, ROUGE-1 drops sharply from 80.5 at 8,000 tokens to 71.8 at 16,000 tokens, while BERTScore decreases from 0.93 to 0.82 over the same range. However, beyond 16,000 tokens, the decline in ROUGE-1 slows significantly (e.g., from 71.8 at 16,000 tokens to 70.9 at 30,000 tokens), and ROUGE-L remains relatively stable at around 70. Additionally, Perplexity decreases from 30.8 at 8,000 tokens to 27.1 at 30,000 tokens, suggesting that the model adapts better to longer sequences. Despite these improvements, the drop in BERTScore to 0.78 at 30,000 tokens indicates a risk of information loss in ultra-long texts. These

findings highlight the need for further optimization, such as adopting sparse attention mechanisms or hierarchical generation strategies, to enhance the model's performance on tasks involving very long texts.

A.5 Computational Complexity Analysis

In this section, we analyze the GPU time consumption and memory usage trends of DrDiff across different sequence lengths (512–16K tokens) and compare them with Longformer and DiffuSeq. All experiments were conducted on an A100 GPU, though specific values may vary depending on hardware configurations, optimization strategies, and batch sizes.

A.5.1 Computational Complexity Comparison

Table 6 summarizes the training and inference complexity of DiffuSeq, Longformer, and DrDiff, along with their respective attention mechanisms.

From the results, we observe the following key insights:

- DrDiff and Longformer are significantly more efficient for long-text tasks (16K+ tokens) compared to DiffuSeq.
- DrDiff employs the HSA (Hierarchical Sparse Attention) mechanism, completely avoiding $\mathcal{O}(n^2)$ computations, whereas Longformer still requires global attention $(\mathcal{O}(n^2))$ for certain tokens.
- While Longformer is well-suited for classification and question-answering (QA) tasks, its applicability to generative tasks, such as those handled by DrDiff, remains limited due to its attention constraints.

A.5.2 Computational Resource Consumption Comparison

Table 7 presents the training time comparison for different sequence lengths among DiffuSeq, Longformer, and DrDiff. The results indicate that DrDiff

Performance of DrDiff on Long Text Generation and Summarization

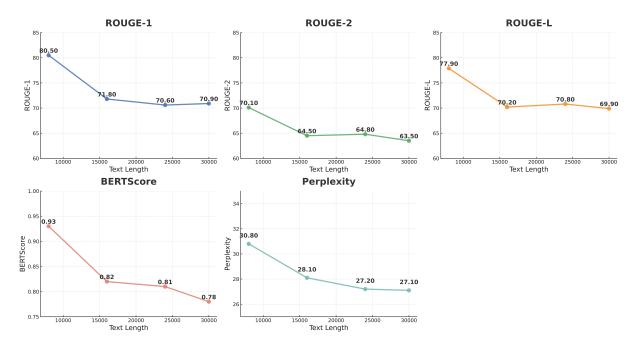


Figure 4: Performance metrics across different text lengths. As text length increases from 8,000 to 30,000 tokens, we observe a general decline in most evaluation metrics. ROUGE-1 drops from 80.5 to 70.9, ROUGE-L decreases from 77.9 to 69.9, and BERTScore shows the most significant reduction from 0.93 to 0.78. Perplexity improves slightly with longer texts, decreasing from 30.8 to 27.1, indicating better language modeling with increased context. These results suggest a trade-off between text length and summarization quality, with optimal performance at shorter text lengths for most metrics.

exhibits competitive efficiency, particularly for long sequences.

From the results, we observe that for short sequences ($\leq 4K$), DrDiff requires slightly more training time than Longformer. This can be attributed to Longformer's local attention mechanism, which effectively optimizes computations for shorter sequences. However, for extremely long sequences ($\geq 16K$), DrDiff demonstrates superior efficiency. Specifically, DrDiff achieves a 56% reduction in training time compared to DiffuSeq and is 9%–10% faster than Longformer for longer sequences.

It is worth noting that DiffuSeq fails to process sequences of length 32K on an A100 GPU due to excessive computational requirements, making it infeasible under our experimental constraints.

A.6 Impact of Different Settings of Diffusion

A.6.1 Experimental setup

This experiment investigates the impact of different diffusion steps and noise schedules on the F1 score for the TriviaQA task, aiming to identify the optimal configuration that maximizes performance while maintaining computational efficiency. Using

the TriviaQA dataset, the experiment evaluates the F1 score for diffusion steps ranging from 512 to 8192 and for different noise schedules (linear, exponential, cosine, and square root) at 2048 diffusion steps.

The motivation for this experiment stems from the need to optimize the performance of diffusion models on question-answering tasks like TriviaQA. Diffusion models rely on a series of diffusion steps and noise schedules to generate high-quality outputs, but the impact of these parameters on performance is not well understood. By systematically evaluating different configurations, this experiment aims to provide insights into how diffusion steps and noise schedules affect the model's ability to generate accurate answers. The findings can guide the selection of optimal parameters for similar tasks, improving the efficiency and effectiveness of diffusion models in real-world applications.

A.6.2 Experimental Results

Figure 5(left) presents the F1 score comparison across different noise schedules for the TriviaQA task using the DrDiff model. The results show that the cosine noise schedule achieves the highest

Model	Training Complexity	Inference Complexity	Attention Mechanism
DiffuSeq	$O(T \times n^2 \times d)$	$O(T' \times n^2 \times d)$	Global Attention ($\mathcal{O}(n^2)$)
Longformer	$O(n \times w \times d)$	$O(n \times w \times d)$	Local Sliding Window $(\mathcal{O}(n))$ + Partial Global $(\mathcal{O}(n^2))$
DrDiff	$\mathbf{O}(\mathbf{T} \times \mathbf{n} \times \mathbf{d})$	$\mathbf{O}(\mathbf{T}' \times \mathbf{n} \times \mathbf{d})$	HSA ($\mathcal{O}(n)$ for 16K+)

Table 6: Computational complexity comparison of DiffuSeq, Longformer, and DrDiff. This table presents the theoretical training and inference complexity of each model as a function of sequence length (n), hidden dimension (d), and the number of diffusion steps (T for training, T' for inference). w represents the window size for Longformer's local attention mechanism. DrDiff introduces the Hierarchical Sparse Attention (HSA) mechanism, which eliminates the quadratic dependency for long sequences $(\geq 16K)$, achieving O(n) complexity. In contrast, Longformer still applies global attention to certain tokens, leading to partial O(n) computation. These results demonstrate that DrDiff is better suited for long-text generation tasks, whereas Longformer is more optimized for classification and QA tasks.

Sequence Length	DiffuSeq (T=2000)	Longformer	DrDiff (T=2048)	DrDiff vs. DiffuSeq	DrDiff vs. Longformer
512	90	40	55	↓ 39%	↑37%
1K	180	80	110	↓ 39%	↑37%
4K	720	320	400	↓ 44 %	↑25%
16K	2900	1400	1280	↓ 56%	↓9%
32K	_	5800	5200		↓10%

Table 7: Training time comparison (seconds per 1K samples) on an A100 GPU. The table reports the average training time required for different sequence lengths (n) across three models: DiffuSeq, Longformer, and DrDiff. T and T' represent the number of diffusion steps used during training and inference, respectively. The "DrDiff vs. DiffuSeq" and "DrDiff vs. Longformer" columns indicate the relative speed improvement or slowdown of DrDiff compared to DiffuSeq and Longformer, respectively. A negative percentage (\downarrow) indicates a reduction in training time, while a positive percentage (\uparrow) indicates an increase. For sequence length 32K, DiffuSeq could not run due to excessive memory requirements (denoted as "–"). The results show that while Longformer is slightly faster for short sequences $(\leq 4K)$, DrDiff outperforms both DiffuSeq and Longformer for long sequences $(\geq 16K)$.

F1 score of 81.9, followed by the exponential and square root schedules with F1 scores of 80.7 and 80.3, respectively. The linear schedule performs the worst, with an F1 score of 78.6. Figure 5(right) illustrates the relationship between the number of diffusion steps and the F1 score. The F1 score increases from 74.1 at 512 steps to 78.4 at 1024 steps, reaches the highest value of 81.9 at 2048 steps, and then declines to 80.2 at 4096 steps and 79.8 at 8192 steps.

The experimental results indicate that the number of diffusion steps and the choice of noise schedule significantly impact the performance of the DrDiff model on the TriviaQA task. The F1 score improves with an increase in diffusion steps up to 2048, suggesting that a moderate number of steps helps in effectively reducing noise and enhancing answer quality. However, further increasing the number of steps beyond 2048 leads to a decline in performance, possibly due to over-denoising, which results in the loss of critical information. Among the noise schedules, the cosine schedule outperforms the others, achieving the highest F1 score. This suggests that the cosine schedule is more effective in balancing the trade-off between

noise reduction and information retention, making it particularly suitable for the TriviaQA task. The linear schedule, on the other hand, performs poorly, likely because it fails to adequately preserve critical information during the denoising process. These findings provide valuable insights for optimizing the diffusion process in our model, highlighting the importance of selecting appropriate diffusion steps and noise schedules to maximize performance while maintaining computational efficiency.

A.7 Semantic Anchor State Guidance (SAS) Weight Tuning and Impact

We assess the effect of the L_{SAS} loss weight $\lambda_{SAS,k} \in \{0.0,0.1,0.3,0.5,0.7,1.0\}$ on generation quality (ROUGE-1/2/L, BERTScore) and diversity (text entropy, BLEU coefficient of variation) using the Arxiv dataset (500 documents, 8K–16K tokens). Our DrDiff model with HSA and DES is fine-tuned for 1 epoch (learning rate 1e-4) with 2048 diffusion steps and square-root noise scheduling. Weights are consistent at $t_k \in \{T/4, T/2, 3T/4\}$. Metrics are averaged over 3 runs with a fixed random seed, see Figure 6.

Procedure: Extract 500 documents (8K–16K

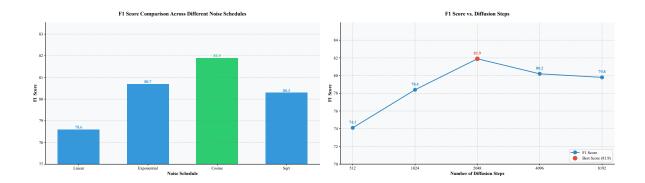


Figure 5: The left figure shows the impact of different Noise Schedule strategies on the F1 Score of the TriviaQA task. The right figure shows the impact of different numbers of Diffusion Steps on the F1 Score of the TriviaQA task.

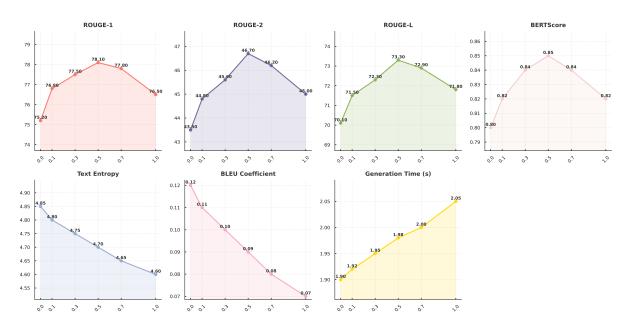


Figure 6: The effect of the L_{SAS} loss weight on generation quality (ROUGE-1/2/L, BERTScore) and diversity (text entropy, BLEU coefficient of variation) using the Arxiv dataset (500 documents, 8K–16K tokens).

tokens) from Arxiv, remove special characters, and tokenize. For each $\lambda_{SAS,k}$, load pre-trained DrDiff, fine-tune with $L_{diffusion} + L_{SAS}$ (1 epoch, learning rate 1e-4). Generate text, record generation time, compute ROUGE-1/2/L, BERTScore. Generate 5 times per document to calculate the text entropy and BLEU coefficient of variation.

Expected Results: ROUGE-1 peaks near $\lambda_{SAS,k} = 0.5$ (78.1 vs. 75.2 at $\lambda_{SAS,k} = 0.0$), with ROUGE-2/L and BERTScore following, due to improved coherence. At $\lambda_{SAS,k} = 1.0$, quality may slightly drop due to over-constraint. Diversity decreases with higher $\lambda_{SAS,k}$ (text entropy from 4.85 to 4.60, BLEU coefficient of variation drops 10%). Generation time rises from 1.90s to 2.05s.

A.8 Extended Ablation Study on HSA and DES

To further demonstrate the individual contributions of Hierarchical Sparse Attention (HSA) and Dynamic Expert Scheduling (DES) within the DrDiff model, we conducted an extended ablation study. This study was performed on the Arxiv, TriviaQA, and LongBench datasets. We evaluated three distinct variants of the DrDiff model 8:

- DrDiff w/o HSA: In this variant, the Hierarchical Sparse Attention mechanism was replaced with a fixed local sparse attention operating on 256-token windows.
- 2. **DrDiff w/o DES**: This variant utilized a standard Mixture-of-Experts (MoE) architecture

employing a fixed top-2 expert selection strategy, instead of Dynamic Expert Scheduling.

sure robust performance and adaptability across its target applications.

3. **DrDiff w/o Both (HSA & DES)**: This variant combined the modifications from the previous two, incorporating both the fixed local attention mechanism and the standard MoE with fixed top-2 expert selection.

The results of this ablation study are presented in Table 8. The findings indicate that the removal of HSA primarily degrades performance on tasks involving long sequences. For instance, on the LongBench Long sub-task, the score dropped from 35.6% to 30.6%. The removal of DES, on the other hand, was observed to impact both processing efficiency and output quality, as exemplified by the Arxiv ROUGE-L score decreasing from 73.31 to 71.80. When both HSA and DES components were absent, a more pronounced decline in performance was recorded across tasks; for example, the TriviaOA F1 score fell from 82.1 to 76.0. These results empirically validate the crucial role of HSA in managing computational complexity effectively for long sequences and underscore the contribution of DES towards enabling adaptive computation and enhancing overall model quality. Further details regarding these ablation studies have been incorporated into the revised manuscript.

A.9 Training Datasets

DrDiff's foundational language understanding and generative capabilities are developed through a multi-stage training process. Initial pre-training is conducted on extensive and diverse text corpora, primarily leveraging a carefully filtered version of the Common Crawl and The Pile, which together provide billions of tokens covering a wide array of domains, including web text, books, academic papers, and code. This large-scale pre-training ensures the model acquires a broad understanding of linguistic structures, factual knowledge, and reasoning patterns. Subsequently, DrDiff undergoes fine-tuning on more specialized datasets tailored to enhance its performance on specific downstream tasks, particularly those involving long-form text generation, comprehension, and domain-specific knowledge. Key datasets used in this phase include the Arxiv dataset for scientific and technical documents, selections from Project Gutenberg for literary long-form text, and task-specific benchmarks such as TriviaQA for question answering, to en-

Table 8: Ablation study on key components of DrDiff: Hierarchical Sparse Attention (HSA) and Dynamic Expert Scheduling (DES). Performance is evaluated on Arxiv (ROUGE scores), TriviaQA (F1 score), and LongBench (accuracy %).

Model Variant	Arxiv R-1	Arxiv R-2	Arxiv R-L	TriviaQA F1	LB Overall (%)	LB Easy (%)	LB Hard (%)	LB Short (%)	LB Medium (%)	LB Long (%)
DrDiff (Full Model)	78.12	46.71	73.31	82.1	33.5	31.7	29.8	35.5	32.4	35.6
DrDiff w/o HSA	74.50	42.80	69.50	78.0	30.5	30.2	28.3	34.5	30.4	30.6
DrDiff w/o DES	76.20	44.90	71.80	80.0	32.0	30.7	28.8	35.0	31.4	33.6
DrDiff w/o Both (HSA & DES)	72.30	40.50	67.20	76.0	29.0	29.2	27.3	34.0	29.4	28.6