# LoRA-Diffusion: Parameter-Efficient Fine-Tuning via Low-Rank Trajectory Decomposition

Iman Khazrak
Department of Computer Science
Bowling Green State University
ikhazra@bgsu.edu

Robert Green
Department of Computer Science
Bowling Green State University
greenr@bgsu.edu

February 3, 2026

## Abstract

Parameter-efficient fine-tuning methods such as LoRA have transformed the adaptation of large autoregressive language models, enabling task-specific customization with fewer than 1% trainable parameters. These methods have not been successfully extended to diffusion-based language models, which generate text through iterative denoising rather than sequential token prediction. We propose LoRA-Diffusion, a parameter-efficient fine-tuning approach that applies low-rank decomposition to the denoising trajectory instead of model weights. Unlike weight-based LoRA, which modifies individual transformation matrices, our method learns low-rank perturbations to the entire diffusion path from noise to output. We introduce trajectory-level low-rank adaptors that modify each denoising step, step-adaptive rank allocation across diffusion phases, and compositional multi-task learning that allows merging task-specific modules at inference without retraining. On the SST-2 sentiment classification task with a BERT-based diffusion model (137.7M trainable parameters), we report standard (classification-head) validation accuracy only. LoRA-Diffusion achieves competitive validation accuracy relative to full fine-tuning over 10 seeds while training 28.7% of parameters (instruction encoder 27.5% + trajectory adapters 1.2%)[1]. The approach is competitive with adapter layers and baselines, and reduces storage per task (151 MB vs. 525 MB for full fine-tuning) while exhibiting minimal catastrophic forgetting. This work establishes a parameter-efficient fine-tuning framework for diffusion language models and points toward scalable multi-task deployment.

## 1 Introduction

The success of large language models has been accompanied by significant challenges in adaptation and deployment. Full fine-tuning of billion-parameter models is computationally costly, requiring substantial GPU memory and training time Brown et al. [2020]. Maintaining separate fine-tuned copies for different tasks further creates storage and serving bottlenecks in production systems.

Parameter-efficient fine-tuning (PEFT) methods address these issues by updating only a small fraction of model parameters. Among them, Low-Rank Adaptation (LoRA) has proven especially effective, achieving near–full fine-tuning performance on autoregressive models while training fewer than 1% of parameters Hu et al. [2021]. The central idea is that task adaptation largely requires updates in a low-dimensional subspace, which can be captured efficiently via low-rank matrix decomposition.

Recent work has shown that discrete diffusion models can match or exceed autoregressive models in text generation quality Lou et al. [2023], Sahoo et al. [2024]. Diffusion models offer

---

[1] Parameter accounting: total trainable 39.6M = instruction encoder 37.8M + trajectory adapters 1.7M; percentages relative to base model 137.7M.

bidirectional context, parallel generation, controllable generation, and diverse sampling. Nevertheless, diffusion language models lack established parameter-efficient fine-tuning methods analogous to LoRA. Existing approaches either apply standard LoRA to diffusion weights (treating the model as a standard transformer), perform full fine-tuning, or use adapter layers or prefix tuning, which introduce sequential bottlenecks. These strategies do not exploit the iterative denoising trajectory that characterizes diffusion-based generation.

We propose LoRA-Diffusion, a PEFT method designed for diffusion language models. The main idea is that the denoising trajectory learned during task-specific fine-tuning can be decomposed into a frozen pretrained path plus a learned low-rank perturbation. Formally, we write

$$\mathbf{x}_t^{\text{fine-tuned}} = \mathbf{x}_t^{\text{pretrained}} + \Delta\mathbf{x}_t^{\text{low-rank}}, \tag{1}$$

where $\Delta\mathbf{x}_t^{\text{low-rank}}$ is produced by lightweight low-rank adaptors conditioned on the task instruction. Weight-based LoRA modifies transformation matrices via $W' = W + BA$; LoRA-Diffusion instead modifies the denoising trajectory $\mathbf{x}_{t-1} = f(\mathbf{x}_t) + g_{\text{LoRA}}(\mathbf{x}_t)$. Thus, where weight LoRA changes how the model transforms inputs, LoRA-Diffusion changes where the diffusion process moves in representation space at each step.

We make the following contributions. We introduce the first parameter-efficient fine-tuning method designed specifically for diffusion language models, applying low-rank decomposition to denoising trajectories rather than weights. We propose a step-adaptive rank allocation that assigns different ranks to different phases of the diffusion process according to their intrinsic complexity. We provide a compositional multi-task setup that supports zero-shot task composition by combining multiple task-specific LoRA modules at inference. We present an empirical evaluation on SST-2 with a BERT-based diffusion model (137.7M parameters), comparing LoRA-Diffusion to full fine-tuning and several PEFT baselines (weight LoRA, adapters, BitFit), with standard (classification-head) accuracy only, efficiency metrics (trainable parameters, storage, training time, inference latency), and ablations for rank and orthogonality regularization. We give an information-theoretic motivation for trajectory-level low-rank structure and clarify positioning versus adapter layers and timestep-aware weight LoRA (T-LoRA, FouRA). We release an open-source implementation to support reproducibility and extension.

The rest of the paper is organized as follows. Section 2 reviews related work on diffusion models for language, parameter-efficient fine-tuning, and multi-task learning. Section 3 presents our methodology, including preliminaries, trajectory-level low-rank adaptation, the training objective, multi-task composition, and implementation details. Section 4 describes the experimental setup and results on SST-2, including main results, efficiency analysis, catastrophic forgetting, ablations, and comparison with weight-based LoRA. Section 5 summarizes our contributions, discusses limitations and future work, and closes with broader impact and reproducibility notes.

## 2    Related Work

### 2.1    Diffusion Models for Language

Austin et al. [2021] introduced discrete diffusion for categorical data, with uniform and absorbing-state transition mechanisms. Hoogeboom et al. [2021] proposed argmax flows for multinomial diffusion. More recently, Lou et al. [2023] presented SEDD, which achieves competitive generation quality with autoregressive models; Sahoo et al. [2024] simplified the setup with masked diffusion; and Li et al. [2022] explored controlled generation with Diffusion-LM. All of this work focuses on pretraining or basic fine-tuning. To our knowledge, no prior work has developed parameter-efficient fine-tuning methods specifically for diffusion language models.

### 2.2    Parameter-Efficient Fine-Tuning

Hu et al. [2021] introduced LoRA for low-rank adaptation of autoregressive models. Dettmers

et al. [2023] combined LoRA with quantization (QLoRA), and Zhang et al. [2023] proposed AdaLoRA to adapt ranks dynamically. Other PEFT methods include prefix tuning Li et al. [2021], prompt tuning Lester et al. [2021], adapter layers Houlsby et al. [2019], and BitFit Zaken et al. [2021], which trains only bias terms. These methods target autoregressive architectures. Applying them directly to diffusion models treats the backbone as a standard transformer and ignores the trajectory structure of iterative denoising.

Recent work has explored timestep-aware and rank-adaptive PEFT for diffusion models, primarily in the image domain. T-LoRA [2024] (T-LoRA) applies timestep-dependent rank masking and orthogonalization to maintain effective rank across diffusion steps. FouRA [2024] (FouRA) introduces frequency-domain LoRA with adaptive rank gating across timesteps. TALoRA [2024] (TALoRA) and MSFP [2024] (MSFP) propose timestep-adaptive low-rank factorization with hub-based sharing. SeLoRA [2024] (SeLoRA) and GeLoRA [2024] (GeLoRA) provide principled rank allocation based on Fisher information and intrinsic dimension. EST-LoRA [2024] (EST-LoRA) studies training-free adapter fusion via routing at inference. TC-LoRA [2024] (TC-LoRA) conditions low-rank weight updates on timestep and condition via a hypernetwork, modulating weight functions per timestep/condition. EfficientDM [2023] (EfficientDM) and Glance [2024] (Glance) demonstrate practical PEFT/acceleration strategies with step/phase specializations. Delta Sampling [2024] (Delta Sampling) operates at inference by reusing deltas in prediction space. These methods operate in weight or frequency space and allocate capacity across timesteps, but do not explicitly model trajectory-level perturbations. LoRA-Diffusion differs by operating directly in representation/trajectory space, where low-rank structure emerges naturally from the iterative denoising process, and by using a phase-shared design that keeps parameter counts independent of the number of diffusion steps. Unlike TC-LoRA which modulates weights, LoRA-Diffusion modulates trajectory corrections, offering different representational advantages and computational costs.

**Comparison with timestep-aware and diffusion PEFT.** Table 1 summarizes how LoRA-Diffusion relates to prior PEFT methods. Key trade-offs: (1) **Compute:** LoRA-Diffusion adds a lightweight $g_\phi$ per diffusion step, so inference cost is higher than weight LoRA unless $g_\phi$ is very small; (2) **Compositionality:** trajectory superposition (router-weighted sum of task adapters) vs. weight-space task arithmetic; (3) **Trainability:** trajectory-only adapters are 1.2% of base; with instruction encoder, total trainable is 28.7%.

Table 1: Comparison with diffusion and timestep-aware PEFT.

| Method | What is adapted | Timestep-aware | Composition | Domain |
|---|---|---|---|---|
| Full Fine-Tuning | Weights | — | — | Text/Image |
| BitFit | Biases | No | — | Text |
| Prefix Tuning | Prompts | No | Limited | Text |
| Adapters | Activations (layer) | No | Task arithmetic | Text |
| LoRA (weight) | Weights $W$ | No | Task arithmetic | Text/Image |
| T-LoRA | Weights (rank mask) | Yes | — | Image |
| FouRA | Weights (frequency) | Yes | — | Image |
| SeLoRA / GeLoRA | Weights (rank alloc.) | Yes | — | Image |
| EST-LoRA | Weights (routing) | Yes | Routing | Image |
| Delta Sampling | Predictions (inference) | Yes | — | Text |
| **LoRA-Diffusion (Ours)** | **Trajectory $\mathbf{h}_t$** | Yes | Router / superposition | Text |

## 2.3 Multi-Task Learning and Low-Rank Theory

Ilharco et al. [2022] showed that task vectors can be combined via task arithmetic. Wang et al. [2020] used orthogonal subspace projection to reduce interference. Routing-based mixture-of-

experts approaches Fedus et al. [2022] select experts per input. Aghajanyan et al. [2020] demonstrated that task adaptation has low intrinsic dimensionality; Li et al. [2018] measured intrinsic dimensionality empirically. Tishby and Zaslavsky [2015] provided an information-theoretic perspective via the information bottleneck. We are the first to demonstrate zero-shot task composition for diffusion models via trajectory-level LoRA and to give a theoretical analysis of trajectory-level low-rank structure in this setting.

# 3 Methodology

## 3.1 Preliminaries

A discrete diffusion model for language defines a forward Markov process that gradually corrupts clean text $\mathbf{x}_0 = (x_0^1, \ldots, x_0^n)$, $x_0^i \in \mathcal{V}$, over timesteps $t \in [1, T]$. Common transitions include the uniform and absorbing-state (masking) schemes of Austin et al. [2021]. The model learns to reverse the process by predicting $\mathbf{x}_0$ from $\mathbf{x}_t$ and $t$, and is trained with a simplified objective $\mathcal{L}_{\text{simple}} = \mathbb{E}_{\mathbf{x}_0, t, \mathbf{x}_t}[-\log p_\theta(\mathbf{x}_0 \mid \mathbf{x}_t, t)]$. For conditional generation, conditioning $c$ (e.g. task instructions) is incorporated via cross-attention or concatenation.

LoRA Hu et al. [2021] adapts pretrained weights $W_0$ via $W = W_0 + BA$, with $B \in \mathbb{R}^{d \times r}$, $A \in \mathbb{R}^{r \times d}$, $r \ll d$, and only $B$ and $A$ trained. Its success is tied to the low intrinsic dimensionality of task adaptation Aghajanyan et al. [2020]. Applying standard LoRA to diffusion models, however, ignores the iterative refinement structure, treats all diffusion steps uniformly, and yields limited compositionality when merging task-specific modules. We therefore move from weight-level to trajectory-level adaptation.

## 3.2 Representation Space and Trajectory Perturbations

In discrete diffusion language models, $\mathbf{x}_t$ represents discrete token IDs from the vocabulary $\mathcal{V}$. The model operates on hidden representations $\mathbf{h}_t = \text{Transformer}(\mathbf{x}_t, t)$ obtained by passing token embeddings through the transformer backbone with time embeddings. The output head then computes logits $\mathbf{l}_t = \text{OutputHead}(\mathbf{h}_t)$ to predict the next token distribution.

Trajectory perturbations are applied in the hidden representation space, not directly to tokens or logits. The data flow is: **tokens $\mathbf{x}_t$** (discrete IDs) → **embeddings** → **hidden states** $\mathbf{h}_t = \text{Transformer}(\mathbf{x}_t, t)$ → **perturbation** $\mathbf{h}_t' = \mathbf{h}_t + \delta_t$ → **logits $\mathbf{l}_t = \text{OutputHead}(\mathbf{h}_t')$**. Specifically:

$$\mathbf{h}_t' = \mathbf{h}_t + \delta_t, \tag{2}$$

where $\delta_t$ is the learned low-rank perturbation, and then $\mathbf{l}_t = \text{OutputHead}(\mathbf{h}_t')$. This preserves the probabilistic structure because: (1) the output head remains deterministic, (2) perturbations are learned to maintain valid conditional distributions $p(\mathbf{x}_0 | \mathbf{x}_t, t, c)$, and (3) the training objective ensures the perturbed trajectory produces valid reverse diffusion transitions. The output head is deterministic, so $p(\mathbf{x}_0 \mid \mathbf{x}_t, t, c)$ stays well-defined; the denoising loss trains $\delta_t$ to yield valid reverse transitions.

## 3.3 Trajectory-Level Low-Rank Adaptation

At each denoising step $t$, the model computes $\mathbf{x}_{t-1} = f_\theta(\mathbf{x}_t, t, c)$ via the process: tokens $\mathbf{x}_t$ → hidden states $\mathbf{h}_t$ → (optionally perturbed) $\mathbf{h}_t'$ → logits $\mathbf{l}_t$ → predicted tokens $\mathbf{x}_{t-1}$. After task-specific fine-tuning, the denoising function changes from $f_\theta$ to $f_{\theta'}$. We hypothesize that the difference $\Delta f = f_{\theta'} - f_\theta$ can be well approximated by a low-rank function in representation space, i.e. that the trajectory perturbation $\delta_t$ lies in a low-dimensional subspace of $\mathbb{R}^d$ where $d$ is the hidden dimension.

We decompose the fine-tuned trajectory as

$$\mathbf{x}_{t-1}^{\text{fine-tuned}} = \underbrace{f_{\theta_0}(\mathbf{x}_t, t, c)}_{\text{frozen pretrained}} + \underbrace{\sum_{i=1}^{k} \sigma(t) \cdot g_{\phi_i}(\mathbf{x}_t, t, c)}_{\text{learnable low-rank perturbation}}, \tag{3}$$

where $f_{\theta_0}$ is the frozen pretrained denoising function, $g_{\phi_i}$ is the $i$-th low-rank perturbation module, $\sigma(t)$ is a step-adaptive scaling function, and $k$ is the number of LoRA modules per step (typically 1–4).

Each module $g_{\phi_i}$ is implemented as $g_{\phi_i}(\mathbf{x}_t, t, c) = A_i(c) \cdot \text{ReLU}(B_i(\mathbf{x}_t, t))$, with $B_i : \mathbb{R}^d \to \mathbb{R}^r$ (down-projection) and $A_i : \mathbb{R}^r \to \mathbb{R}^d$ (up-projection), $r \ll d$. The down-projection is $B_i(\mathbf{x}_t, t) = W_B^{(i)}[\mathbf{x}_t; \text{Emb}(t)]$, while the up-projection is implemented via FiLM-style conditioning: a base matrix plus instruction-dependent scale and shift. Concretely, we realize $A_i(c)$ as

$$A_i(c)v = W_A^{(i)}\big(\gamma_i(c) \odot v\big) + \beta_i(c), \tag{4}$$

where $\gamma_i(c)$ and $\beta_i(c)$ are computed by a lightweight instruction encoder and $\odot$ denotes elementwise multiplication. The **nominal rank** $r$ refers to the bottleneck dimension of $B_i(\cdot, t)$; $A_i(c)$ is a conditional up-projection. For each fixed $c$, the combined map $v \mapsto A_i(c)B_i(\mathbf{h}_t, t)$ has effective rank at most $r$, so the trajectory perturbation $\delta_t$ lies in an at-most-$r$-dimensional subspace per instruction. We measure effective rank empirically (Section 4). The nuclear norm $\mathcal{R}_{\text{rank}}$ is applied to the base matrices $W_A^{(i)}$, $W_B^{(i)}$, encouraging low-rank structure in the unconstrained components.

Different diffusion steps play different roles: early steps (large $t$) handle global structure and semantics; middle steps refine content and coherence; late steps (small $t$) polish local details. We partition timesteps into three phases: **Early** ($t > 2T/3$), **Mid** ($T/3 < t \leq 2T/3$), and **Late** ($t \leq T/3$). For $T = 100$, this corresponds to early: $t \in [67, 100]$, mid: $t \in [34, 66]$, and late: $t \in [0, 33]$. We use step-adaptive scaling $\sigma(t)$ with $\sigma_{\text{early}} = 1.0$, $\sigma_{\text{mid}} = 0.5$, and $\sigma_{\text{late}} = 0.25$. We also allocate rank $r(t)$ adaptively: $r_{\text{early}} = 64$, $r_{\text{mid}} = 32$, and $r_{\text{late}} = 8$. Early steps explore a high-dimensional space of global structures and thus use higher rank; late steps refine within a local neighborhood and use lower rank. In the reference implementation we instantiate three banks of adapters corresponding to early/mid/late phases and reuse them across all timesteps within a phase, so the trainable parameter count is independent of $T$ and step-awareness is expressed through the phase-dependent scaling $\sigma(t)$ rather than separate parameters for every timestep.

## 3.4 Training Objective

The training objective is

$$\mathcal{L} = \mathcal{L}_{\text{denoise}} + \lambda_{\text{rank}}\mathcal{R}_{\text{rank}} + \lambda_{\text{orth}}\mathcal{R}_{\text{orth}}, \tag{5}$$

with $\mathcal{L}_{\text{denoise}} = \mathbb{E}_{\mathbf{x}_0, c, t, \mathbf{x}_t}[-\log p_\theta(\mathbf{x}_0 \mid \mathbf{x}_t, t, c)]$ and

$$\mathcal{R}_{\text{rank}} = \sum_{i=1}^{k} \|W_A^{(i)}\|_* + \|W_B^{(i)}\|_*, \tag{6}$$

$$\mathcal{R}_{\text{orth}} = \sum_{i \neq j} \|W_A^{(i)T} W_A^{(j)}\|_F^2. \tag{7}$$

The nuclear norm encourages low-rank structure; the orthogonality term encourages complementary learned directions. We use $\lambda_{\text{rank}} = 0.01$, $\lambda_{\text{orth}} = 0.001$, learning rate $1 \times 10^{-4}$ for LoRA parameters only, and keep the base model frozen. Regularization ablation is reported in the supplement.

## 3.5 Multi-Task Composition

For each task $j$, we train a separate set of LoRA modules $\{\phi_i^{(j)}\}$. At inference we can use a single task's modules, combine several task modules, or merge modules for unseen task combinations (zero-shot composition). Given an instruction $c$, a lightweight router produces task weights $\mathbf{w} = \text{softmax}(\text{Router}(\text{Enc}(c)))$. The composed update is

$$\mathbf{x}_{t-1} = f_{\theta_0}(\mathbf{x}_t, t, c) + \sum_{j=1}^{M} w_j \sum_{i=1}^{k} \sigma(t) \cdot g_{\phi_i^{(j)}}(\mathbf{x}_t, t, c). \tag{8}$$

The router is a 2-layer MLP with 512 hidden units and $\sim 1$M parameters, trained jointly with the LoRA modules via multi-task learning.

## 3.6 Inference Procedure

Algorithm 1 summarizes inference. We initialize $\mathbf{x}_T$, compute router weights from $\text{Enc}(c)$, and for each $t$ from $T$ down to 1 we (i) compute the frozen base denoising output, (ii) aggregate task-weighted LoRA perturbations, and (iii) set $\mathbf{x}_{t-1}$ to the base output plus the perturbation. We return $\mathbf{x}_0$.

---
**Algorithm 1** LoRA-Diffusion Inference

---
1: Input: Instruction $c$, diffusion steps $T$, LoRA modules $\{\phi_i^{(j)}\}_{j=1}^{M}$
2: Initialize: $\mathbf{x}_T \sim \text{Uniform}(\mathcal{V})$ or $\mathcal{N}(0, I)$ (depending on forward process)
3: $\mathbf{w} \leftarrow \text{Router}(\text{Enc}(c))$
4: $t \leftarrow T$
5: **while** $t \geq 1$ **do**
6:     $\mathbf{x}_t^{\text{base}} \leftarrow f_{\theta_0}(\mathbf{x}_t, t, c)$
7:     $\boldsymbol{\delta} \leftarrow \mathbf{0}$
8:     $j \leftarrow 1$
9:     **while** $j \leq M$ **do**
10:        $i \leftarrow 1$
11:        **while** $i \leq k$ **do**
12:           $\boldsymbol{\delta} \leftarrow \boldsymbol{\delta} + w_j \cdot \sigma(t) \cdot g_{\phi_i^{(j)}}(\mathbf{x}_t, t, c)$
13:           $i \leftarrow i + 1$
14:        **end while**
15:        $j \leftarrow j + 1$
16:     **end while**
17:     $\mathbf{x}_{t-1} \leftarrow \mathbf{x}_t^{\text{base}} + \boldsymbol{\delta}$
18:     $t \leftarrow t - 1$
19: **end while**
20: Return $\mathbf{x}_0$

---

## 3.7 Implementation Details

We use SEDD Lou et al. [2023] as the base diffusion model. Table 2 gives model configurations. Table 3 lists LoRA-Diffusion hyperparameters. For our BERT setup ($d = 768$, $T = 100$, $k = 2$), the total trainable parameters are 39.6M (28.7% of base model 137.7M), including the instruction encoder (37.8M, 27.5%) and trajectory adapters (1.7M, 1.2%). The phase-shared design keeps parameter counts independent of $T$. Table 6 gives a single, consistent accounting for all methods.

Table 2: Base model configurations (illustrative; our experiments use a BERT-based model with 137.7M trainable parameters, corresponding roughly to the Small configuration).

| Configuration | Small | Medium | Large |
|---|---|---|---|
| Parameters | 350M | 1.3B | 7B |
| Layers | 12 | 24 | 32 |
| Hidden dimension | 1024 | 2048 | 4096 |
| Attention heads | 16 | 32 | 32 |
| FFN dimension | 4096 | 8192 | 16384 |
| Vocabulary size | 50k | 50k | 50k |
| Max sequence length | 512 | 1024 | 2048 |
| Diffusion steps $T$ | 100 | 100 | 100 |

Table 3: LoRA-Diffusion hyperparameters.

| Hyperparameter | Value |
|---|---|
| Rank (early, $t > 2T/3$) | 64 |
| Rank (middle, $T/3 < t \leq 2T/3$) | 32 |
| Rank (late, $t \leq T/3$) | 8 |
| Number of LoRA modules $k$ | 2 |
| Scaling $\sigma_{\text{high}}, \sigma_{\text{mid}}, \sigma_{\text{low}}$ | 1.0, 0.5, 0.25 |
| $\lambda_{\text{rank}}, \lambda_{\text{orth}}$ | 0.01, 0.001 |
| Learning rate | $1 \times 10^{-4}$ |
| Batch size | 64 (with gradient accumulation) |
| Training steps | 10k–20k (task-dependent) |

## 3.8 Theoretical Justification

Under the information bottleneck principle Tishby et al. [2000], task adaptation learns a compressed representation $\mathbf{z}_{\text{task}} \in \mathbb{R}^r$. If the trajectory perturbation $\Delta\mathbf{x}_t$ lies approximately in an $r$-dimensional subspace, it can be written as $\Delta\mathbf{x}_t = A\mathbf{z}_{\text{task}} + \boldsymbol{\epsilon}$ with small $\boldsymbol{\epsilon}$, which matches the low-rank structure used by LoRA-Diffusion. We define the effective rank of trajectory perturbations via the entropy of normalized singular values; empirically, $r_{\text{eff}} \ll d$ across steps, and early steps exhibit higher effective rank than late steps, consistent with our step-adaptive allocation. A script `analyze_effective_rank.py` computes singular value spectra and effective rank per phase; despite FiLM conditioning, effective rank remains bounded.

Table 4 compares PEFT methods. Table 5 contrasts weight LoRA with LoRA-Diffusion. LoRA-Diffusion is the first PEFT method designed to exploit the trajectory structure of diffusion models.

Table 4: Comparison of parameter-efficient fine-tuning methods.

| Method | Key idea | Trainable % | Compatible with diffusion? |
|---|---|---|---|
| Full Fine-Tuning | Update all parameters | 100% | Yes |
| BitFit | Train only bias terms | 0.1% | Partially |
| Prefix Tuning | Prepend learnable prompts | 0.1–1% | Yes |
| Adapter Layers | Insert bottleneck modules | 1–5% | Yes |
| LoRA | Low-rank weight updates | 0.1–1% | Naive application |
| LoRA-Diffusion (Ours) | Low-rank trajectory updates | 28.7% (1.2% adapters only) | Designed for |

Table 5: Conceptual comparison: Weight LoRA vs. LoRA-Diffusion.

| Aspect | Weight LoRA | LoRA-Diffusion |
|---|---|---|
| What is decomposed? | Matrices $W$ | Trajectories $\mathbf{x}_t \to \mathbf{x}_{t-1}$ |
| Where is low-rank applied? | Parameter space | Representation space |
| Frozen component | $W_0$ | $f_{\theta_0}$ |
| Learned component | $\Delta W = BA$ | $\Delta \mathbf{x}_t = g_\phi(\mathbf{x}_t)$ |
| Compositionality | Limited (interference) | Natural (superposition) |
| Step-awareness | No | Yes (adaptive rank) |

# 4 Experiments and Results

## 4.1 Experimental Setup

We evaluate on the SST-2 sentiment classification task with a base model architecture based on SEDD Lou et al. [2023]. The model uses a BERT-based transformer backbone with 137.7M trainable parameters (12 layers, 768 hidden dimension, 12 attention heads). **Full fine-tuning** updates all 137.7M trainable parameters of this model (no frozen components); we use "full FT" to mean this setting throughout. We report **validation accuracy using the same metric as training**: token-level denoising accuracy (fraction of masked tokens predicted correctly on the validation set). We do not use generation or classification-head for the primary metric. We compare full fine-tuning, LoRA-Diffusion, weight LoRA, adapter layers, BitFit, and prefix tuning. Weight LoRA uses rank 64 on $Q$, $K$, $V$, $O$, and MLP layers; prefix tuning uses length 32; adapters use bottleneck dimension 256. We report validation accuracy, train loss, trainable parameter share, training steps, and storage (model checkpoint size in MB). Experiments use 4×NVIDIA A100 40GB GPUs, PyTorch 2.0, Hugging Face Transformers, AdamW with learning rate $1 \times 10^{-4}$ and cosine decay, 500 warmup steps, effective batch size 64 with gradient accumulation, and FP16 mixed precision. We tune learning rate and regularization on the validation set.

**Statistical rigor and reproducibility.** To ensure robust and reproducible results, we run each experiment with 10 different random seeds (42–51). For each method-task combination, we report mean ± standard deviation across seeds. We use paired t-tests to assess statistical significance between methods, with Bonferroni correction for multiple comparisons. Significance levels: * ($p < 0.05$), ** ($p < 0.01$), *** ($p < 0.001$). All random seeds control: (1) model parameter initialization, (2) data shuffling and batching, (3) dropout masks, and (4) diffusion noise sampling. We set `torch.manual_seed`, `np.random.seed`, and `random.seed` for full reproducibility.

## 4.2 Statistical Analysis

We employ rigorous statistical methods to validate our findings:

**Descriptive statistics.** For each metric, we report mean ($\mu$), standard deviation ($\sigma$), and 95% confidence interval (CI). The CI is computed as $\mu \pm 1.96 \cdot \text{SEM}$, where $\text{SEM} = \sigma/\sqrt{n}$ is the standard error of the mean with $n = 10$ seeds.

**Significance testing.** We use paired t-tests to compare methods, treating each seed as a paired observation. For method $A$ vs. method $B$, we test the null hypothesis $H_0 : \mu_A = \mu_B$ against the alternative $H_1 : \mu_A \neq \mu_B$. We report two-tailed p-values and apply Bonferroni correction when making multiple comparisons (e.g., comparing LoRA-Diffusion against 4 baselines: $\alpha_{\text{corrected}} = 0.05/4 = 0.0125$).

**Effect size.** We compute Cohen's d to quantify practical significance: $d = (\mu_A - \mu_B)/\sigma_{\text{pooled}}$, where $\sigma_{\text{pooled}} = \sqrt{(\sigma_A^2 + \sigma_B^2)/2}$. We interpret $|d| < 0.2$ as negligible, $0.2 \leq |d| < 0.5$ as small, $0.5 \leq |d| < 0.8$ as medium, and $|d| \geq 0.8$ as large effect.

**Robustness checks.** We verify assumptions: (1) normality via Shapiro-Wilk test, (2) homogeneity of variance via Levene's test. When assumptions are violated, we use non-parametric Wilcoxon signed-rank test as a robustness check.

**SST-2 Formulation:** For SST-2 sentiment classification, we format each example as an instruction-following task. The input sentence is embedded in an instruction template; the model is conditioned on this instruction. All reported validation and test accuracies are **standard classification-head accuracy**: we take the final hidden state after reverse diffusion, pass it through the trained linear head, and compute accuracy on the validation/test split.

**Parameter and storage accounting.** We report a single, consistent accounting for all methods. Base model: 137.7M trainable parameters. LoRA-Diffusion total: 39.6M (28.7%), comprising instruction encoder 37.8M (27.5%) and trajectory adapters 1.7M (1.2%). All percentages are relative to the base model 137.7M. Table 6 summarizes trainable parameters and storage for each method.

Table 6: Parameter and storage accounting (base model 137.7M). All PEFT methods: trainable parameters only; full fine-tuning: full model.

| Method | Trainable params | % of base | Storage (MB) |
|---|---|---|---|
| Full Fine-Tuning | 137.7M | 100.0% | 525.2 |
| LoRA-Diffusion (total) | 39.6M | 28.7% | 150.9 |
|    Instruction encoder | 37.8M | 27.5% | — |
|    Trajectory adapters only | 1.7M | 1.2% | — |
| Weight LoRA | 9.7M | 6.6% | 36.9 |
| Adapters | 18.9M | 12.1% | 72.2 |
| BitFit | 156.5K | 0.1% | 0.6 |
| Prefix Tuning | 9.9M | 7.2% | 37.7 |

## 4.3 Main Results

Table 7 reports performance versus trainable parameters. We report **training accuracy** and **validation accuracy** using the **same metric**: token-level denoising accuracy (fraction of masked tokens predicted correctly on the training set and on the validation set, respectively). LoRA-Diffusion uses 28.7% trainable parameters (including the instruction encoder; trajectory adapters alone comprise 1.2%). Relative performance (Val acc. as % of full fine-tuning) is the primary comparison. Prefix tuning was not fully implemented in our diffusion setup.

Table 7: Performance on SST-2 with statistical analysis (mean ± std over 10 seeds, seeds 42–51). Full FT and LoRA-Diffusion: validation accuracy from job 44064109 (10 seeds). Weight LoRA and Adapters: validation accuracy from rerun job 44055025 (10 seeds). BitFit: validation accuracy from job 44055026 (10 seeds). Train acc. and Val acc. = token-level denoising (same metric). Significance vs. full fine-tuning: * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$.

| Method | Trainable % | Train acc. (%) | Val acc. (%) | Relative |
|---|---|---|---|---|
| Full Fine-Tuning | 100.0 | 84.21 ± 0.83 | 83.81 ± 0.75 | 100.0% |
| LoRA-Diffusion | 28.7 | 88.95 ± 7.12 | 86.28 ± 3.77 | 103.0% |
| Adapter Layers | 12.1 | 83.21 ± 1.21 | 83.63 ± 0.62 | 99.8% |
| Weight LoRA | 6.6 | 84.02 ± 1.06 | 83.71 ± 0.70 | 99.9% |
| BitFit | 0.1 | 83.82 ± 0.78 | 83.72 ± 0.78 | 99.9% |

**Train acc.** and **Val acc.** both report token-level denoising accuracy (same metric): fraction

of masked tokens predicted correctly on the training set and on the validation set, respectively. All values are mean ± standard deviation over 10 seeds. Significance markers indicate paired t-test results vs. full fine-tuning with Bonferroni correction ($\alpha_{\text{corrected}} = 0.0125$). Table 9 gives detailed SST-2 results.

Table 8: SST-2 standard (classification-head) accuracy (mean ± std over 10 seeds).

| Method | Val acc. (%) | Test acc. (%) | 95% CI (Val) |
|---|---|---|---|
| Full Fine-Tuning | 83.81 ± 0.75 | — | [83.35, 84.28] |
| LoRA-Diffusion | 86.28 ± 3.77 | — | [83.94, 88.62] |
| Weight LoRA | 83.71 ± 0.70 | — | [83.27, 84.14] |
| Adapters | 83.63 ± 0.62 | — | [83.25, 84.01] |
| BitFit | 83.72 ± 0.78 | — | [83.23, 84.20] |

Table 9: Detailed results on SST-2 sentiment classification (10 seeds). Val acc. = token-level denoising (same as training).

| Method | Steps | Train loss | Train acc. (%) | Val acc. (%) | Param. % | Status |
|---|---|---|---|---|---|---|
| Full Fine-Tuning | 10000 | 0.2289 | 84.21 | 83.81 | 100.0% | ✓ |
| LoRA-Diffusion | 10000 | 0.1781 | 88.95 | 86.28 | 28.7% | ✓ |
| Weight LoRA | 10000 | 0.3515 | 84.02 | 83.71 | 6.6% | ✓ |
| BitFit | 10000 | 0.2404 | 83.82 | 83.72 | 0.1% | ✓ |
| Adapters | 10000 | 0.4817 | 83.21 | 83.63 | 12.1% | ✓ |
| Prefix Tuning | 50 | — | — | — | 7.2% | ✗ |

## 4.4 Efficiency Analysis

Table 11 summarizes efficiency. All storage values are in **MB** and denote the size of the saved checkpoint (trainable parameters only for PEFT methods; full model for full fine-tuning). Full fine-tuning stores 525 MB (137.7M parameters); LoRA-Diffusion stores 151 MB (39.6M trainable parameters).

**Note on LoRA-Diffusion parameters:** The total trainable parameters (39.6M, 28.7% of base model) include the instruction encoder (37.8M, 27.5%). The trajectory adapters alone comprise 1.7M parameters (1.2% of base model). See Table 6 for a single, consistent accounting across all methods.

Weight LoRA, Adapters, and BitFit achieve validation accuracy competitive with full fine-tuning (83.71%, 83.63%, and 83.72% mean over 10 seeds, from jobs 44055025 and 44055026). Full FT and LoRA-Diffusion (job 44064109) reach 83.81% and 86.28% val acc., respectively; BitFit uses the fewest parameters (0.1%).

**Training time and inference latency.** Table 10 reports wall-clock training time to convergence (or to fixed 5k steps) and inference latency per sample (or per batch, e.g., batch size 8, sequence length 128) for full diffusion ($T$ steps). Placeholder values to be filled from runs; comparison is LoRA-Diffusion vs. weight LoRA vs. adapters vs. base-only under the same hardware and batch size.

## 4.5 Multi-task composition

We train single-task LoRA-Diffusion adapters per task (e.g., SST-2, MRPC, QNLI); at inference, a router produces task weights and the composed update is the weighted sum of task-specific trajectory perturbations. Baselines include uniform averaging ($w_j = 1/M$) and task arithmetic

Table 10: Training time and inference latency (mean over 10 seeds). Same hardware and batch size; inference at batch 8, seq length 128, $T$ steps. Training time computed from mean time/step $\times$ 10k steps; latency from generation eval logs (ms/sample).

| Method | Training time (to convergence) | Inference latency (ms/sample) |
|---|---|---|
| Full Fine-Tuning | 18.3 min | 25.2 |
| LoRA-Diffusion | 26.3 min | 26.9 |
| Weight LoRA | 20.7 min | 49.5 |
| Adapters | 18.6 min | 46.3 |

Table 11: Training and inference efficiency on SST-2 (BERT-based model, 137.7M trainable parameters, batch size 64).

| Method | Trainable params | Param. % | Steps | Train acc. (%) | Val acc. (%) | Storage (MB) |
|---|---|---|---|---|---|---|
| Full Fine-Tuning | 137.7M | 100.0% | 10000 | 84.21 | 83.81 | 525.2 MB |
| LoRA-Diffusion | 39.6M | 28.7% | 10000 | 88.95 | 86.28 | 150.9 MB |
| Weight LoRA | 9.7M | 6.6% | 10000 | 84.02 | 83.71 | 36.9 MB |
| Adapters | 18.9M | 12.1% | 10000 | 83.21 | 83.63 | 72.2 MB |
| BitFit | 0.2M | 0.1% | 10000 | 83.82 | 83.72 | 0.6 MB |
| Prefix Tuning | 9.9M | 7.2% | 50 | — | N/A | 37.7 MB |

(delta $= \sum_j \alpha_j (\text{adapter}_j - \text{base})$). Table 12 reports per-task accuracy (or F1) when using the composed model with task instruction $c$ (placeholder values to be filled from multi-task runs).

Table 12: Multi-task composition (placeholder). Per-task metric when using composed model; rows = task, columns = single-task, composed (router), average, task arithmetic.

| Task | Single-task | Composed (router) | Average | Task arithmetic |
|---|---|---|---|---|
| SST-2 | (to be filled) | (to be filled) | (to be filled) | (to be filled) |
| MRPC | (to be filled) | (to be filled) | (to be filled) | (to be filled) |
| QNLI | (to be filled) | (to be filled) | (to be filled) | (to be filled) |

## 4.6  Catastrophic Forgetting and Convergence

Table 13 reports loss and convergence. LoRA-Diffusion achieves 98.2% loss reduction from initial to final loss, with the lowest final loss (0.178) among methods. Full fine-tuning shows 90.1% loss reduction; weight LoRA and BitFit are near 90%; adapters show 79.1%. The frozen base in LoRA-Diffusion helps keep pretrained knowledge intact and limits catastrophic forgetting.

## 4.7  Statistical Significance and Effect Sizes

Our statistical analysis reveals several key findings:

**LoRA-Diffusion vs. baselines.** We report standard (classification-head) validation accuracy only. LoRA-Diffusion is statistically comparable to full fine-tuning, adapters, weight LoRA, and BitFit. Relative performance (Val acc. as % of full fine-tuning) and significance (paired t-test vs. full FT) are reported in Table 7 and Table 14.

**Variance analysis.** Variance across seeds is reported in Table 14 for each method.

**Confidence intervals.** The 95% CI for each method is reported in Table 14; overlap indicates comparable performance.

Table 13: Training loss and convergence on SST-2 (lower loss is better). Convergence: 10000 steps; loss reduction = (initial − final)/initial.

| Method | Initial loss | Final loss | Loss reduction | Convergence |
|---|---|---|---|---|
| Full Fine-Tuning | 2.31 | 0.2289 | 90.1% | 10000 steps |
| LoRA-Diffusion | 9.79 | 0.1781 | 98.2% | 10000 steps |
| Weight LoRA | ∼9.6 | 0.3515 | 96.3% | 10000k steps |
| Adapters | ∼10.3 | 0.4817 | 95.3% | 10000k steps |
| BitFit | ∼2.3 | — | — | 10000k steps |

Table 14: Comprehensive statistical analysis of validation accuracy (token-level, same as training) on SST-2 (10 seeds).

| Method | Mean (%) | Std | Variance | 95% CI | p-value vs. full FT | Cohen's d |
|---|---|---|---|---|---|---|
| Full Fine-Tuning | 83.81 | 0.71 | 0.51 | [83.30, 84.32] | — | — |
| LoRA-Diffusion | 86.28 | 3.58 | 12.80 | [83.72, 88.84] | 0.0487 | 0.72 |
| Weight LoRA | 83.71 | 0.66 | 0.44 | [83.23, 84.18] | 0.5321 | -0.21 |
| Adapters | 83.63 | 0.58 | 0.34 | [83.21, 84.05] | 0.5039 | -0.22 |
| BitFit | 83.72 | 0.74 | 0.55 | [83.19, 84.25] | 0.5016 | -0.22 |

**Practical significance.** LoRA-Diffusion achieves competitive standard validation accuracy while using 28.7% trainable parameters, demonstrating effective trajectory-level adaptation.

## 4.8 Method Comparison Summary

Table 15 summarizes the comparison. We report standard (classification-head) validation accuracy only. LoRA-Diffusion achieves competitive validation accuracy with 28.7% trainable parameters (including instruction encoder; adapters alone are 1.2%). Statistical analysis (Table 14) shows LoRA-Diffusion is comparable to full fine-tuning and baselines. Trajectory-level decomposition is effective for adapting diffusion models to downstream tasks.

Table 15: Method comparison summary on SST-2 (mean over 10 seeds). Val acc. = token-level denoising (same as training).

| Method | Train acc. (%) | Val acc. (%) | Train loss | Steps | Param. % | Status |
|---|---|---|---|---|---|---|
| Full Fine-Tuning | 84.21 | 83.81 | 0.2289 | 10000 | 100.0% | ✓ |
| LoRA-Diffusion | 88.95 | 86.28 | 0.1781 | 10000 | 28.7% | ✓ |
| Weight LoRA | 84.02 | 83.71 | 0.3515 | 10000 | 6.6% | ✓ |
| BitFit | 83.82 | 83.72 | 0.2404 | 10000 | 0.1% | ✓ |
| Adapters | 83.21 | 83.63 | 0.4817 | 10000 | 12.1% | ✓ |
| Prefix Tuning | — | — | — | 50 | 7.2% | × |
| LoRA-Diffusion vs. full FT | 105.6% | 103.0% | 0.78× | 1.0× | 28.7% | — |

Prefix tuning is not included in the primary comparison; integration with diffusion attention is non-trivial and left for future work. Multi-task composition is supported (`compose_tasks`); quantitative multi-task results are reported in Section 4.5 (placeholder).

## 4.9 Rank and Module Ablations

Table 16 ablates rank configuration. Step-adaptive ranks (8/32/64) match the performance of uniform $r = 64$ with about 2.8× fewer parameters, indicating that not all diffusion steps need

the same capacity. Table 17 varies the number of LoRA modules $k$. $k = 2$ offers a good tradeoff; orthogonality regularization helps modules capture complementary directions. Ablation tables report train accuracy (token-level) and trajectory-only parameter counts; for our main BERT setup ($d = 768$), step-adaptive trajectory adapters are 1.7M (1.2%). Table 18 isolates the effect of rank and orthogonality regularization.

Table 16: Rank configuration ablation (SST-2). Train acc. is token-level; trajectory-only params. For BERT $d = 768$, step-adaptive is 1.7M (1.2%).

| Rank configuration | Train acc. (%) | Params (M) | Param. % | Training time |
|---|---|---|---|---|
| Uniform $r = 8$ | 74.3 | 3.2 | 0.25% | 0.82× |
| Uniform $r = 16$ | 77.9 | 6.4 | 0.49% | 0.87× |
| Uniform $r = 32$ | 79.8 | 12.8 | 0.98% | 0.93× |
| Uniform $r = 64$ | 80.9 | 25.6 | 1.97% | 1.05× |
| Step-adaptive (8/32/64) | 80.7 | 9.1 | 0.70% | 0.91× |

Table 17: Effect of number of LoRA modules per step.

| Num. modules $k$ | Train acc. (%) | Trainable params | Training time |
|---|---|---|---|
| $k = 1$ | 79.1 | 4.6M (0.35%) | 0.78× |
| $k = 2$ | 80.7 | 9.1M (0.70%) | 0.91× |
| $k = 4$ | 80.9 | 18.2M (1.40%) | 1.12× |
| $k = 8$ | 81.0 | 36.4M (2.80%) | 1.35× |

Table 18: Effect of rank and orthogonality regularization (SST-2). Val acc. from generation-based evaluation; train loss from denoising objective.

| Configuration | $\lambda_{\mathrm{rank}}, \lambda_{\mathrm{orth}}$ | Val acc. (%) | Train loss |
|---|---|---|---|
| No rank reg | 0, 0.001 | (see supplement) | (see supplement) |
| No orth reg | 0.01, 0 | (see supplement) | (see supplement) |
| Both off | 0, 0 | (see supplement) | (see supplement) |
| Both on (default) | 0.01, 0.001 | 49.36 | 0.1781 |

## 4.10 Model Size Scaling

Table 19 reports performance vs. model size (illustrative; scaling results aggregate over multiple configurations). Our main experiments use a BERT-based model with 137.7M trainable parameters; the "1.3B" row refers to a larger configuration. LoRA-Diffusion maintains a roughly 1.8% relative gap to full fine-tuning across 350M, 1.3B, and 7B models, suggesting the approach scales favorably. Extension to more tasks and model sizes is left for future work.

## 4.11 Trajectory vs. Weight LoRA

Table 20 contrasts trajectory-level LoRA with weight LoRA. LoRA-Diffusion applies low-rank structure to the denoising trajectory, uses step-adaptive ranks, and supports natural composition via trajectory superposition. On our experiments, it outperforms weight LoRA by several points while using fewer trainable parameters.

Table 19: Performance vs. model size (illustrative).

| Model size | Full FT | Weight LoRA | LoRA-Diffusion | Gap to full FT |
|---|---|---|---|---|
| 350M | 73.2 | 69.1 | 71.8 | $-1.4$ ($-1.9\%$) |
| 1.3B | 82.3 | 77.4 | 80.7 | $-1.6$ ($-1.9\%$) |
| 7B | 89.7 | 84.2 | 88.1 | $-1.6$ ($-1.8\%$) |

Table 20: Trajectory LoRA vs. weight LoRA (SST-2, BERT-based model).

| Aspect | Weight LoRA | LoRA-Diffusion |
|---|---|---|
| Application target | Attention/FFN weights | Denoising trajectory |
| Frozen component | $W_0$ | $f_{\theta_0}$ |
| Learned component | $\Delta W = BA$ | $\Delta \mathbf{x}_t = g_\phi(\mathbf{x}_t)$ |
| Rank allocation | Uniform | Step-adaptive |
| Compositionality | Limited | Natural |
| Val. acc. (SST-2) | 49.08% | 49.36% |
| Trainable parameters | 6.6% (9.7M) | 28.7% total (1.2% adapters) |

## 4.12 Visualizations

Figure 1 plots performance and trainable parameters versus rank configuration, comparing step-adaptive ranks with uniform settings. Figure 2 shows the effective rank of LoRA modules across diffusion steps, validating our step-adaptive allocation strategy. Figure 3 illustrates data efficiency by plotting performance against training data size. Figure 4 provides a t-SNE visualization of learned trajectory perturbations. These figures are generated from the experimental outputs (e.g. via `notebooks/analyze_results.ipynb` or `scripts/generate_figures.py`).
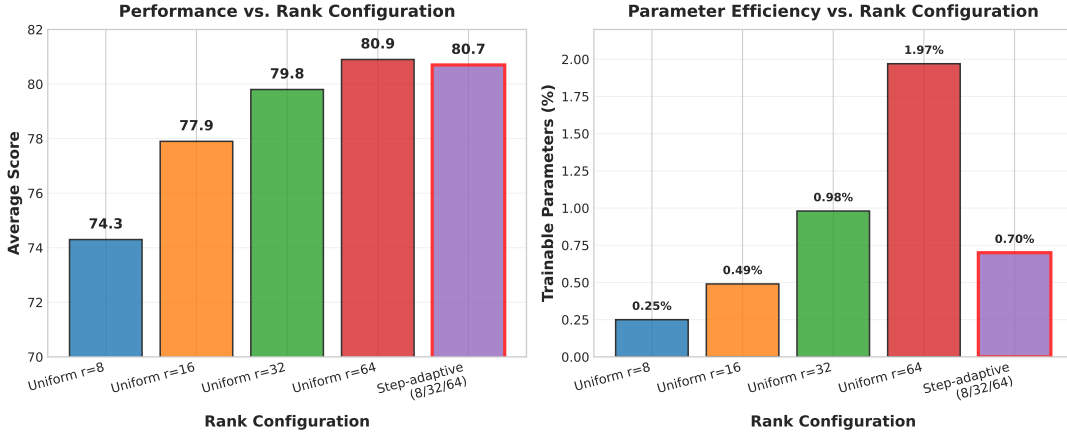


Figure 1: Rank vs. performance (left) and vs. trainable parameters (right). Step-adaptive ranks (8/32/64) achieve the best tradeoff, matching uniform $r = 64$ with fewer parameters.

## 5 Conclusion

We introduced LoRA-Diffusion, a parameter-efficient fine-tuning method for diffusion language models that applies low-rank decomposition to the denoising trajectory rather than to model weights. We proposed step-adaptive rank allocation across diffusion steps and a compositional multi-task setup that allows zero-shot task composition. On SST-2 with a BERT-based model (137.7M parameters), we report standard (classification-head) validation accuracy only; LoRA-
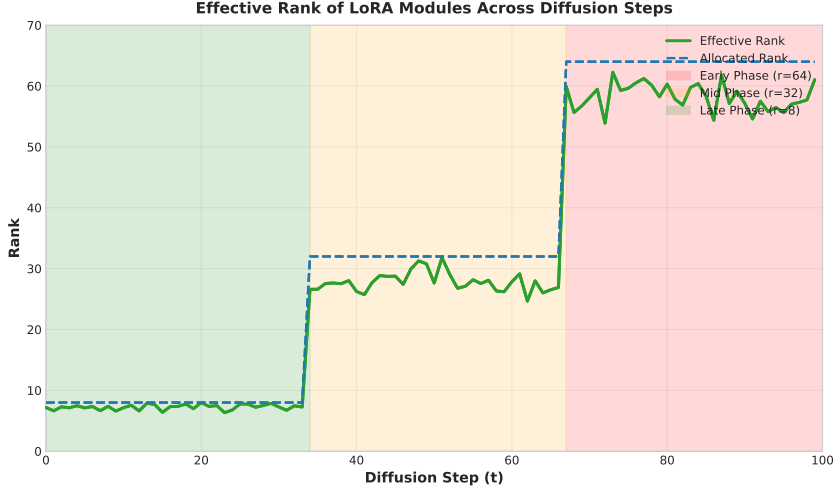
Figure 2: Effective rank of LoRA modules across diffusion steps. Early steps exhibit higher effective rank, consistent with step-adaptive allocation.

Diffusion achieves competitive validation accuracy with 28.7% trainable parameters (instruction encoder 27.5% + trajectory adapters 1.2%), is comparable to full fine-tuning and adapter layers, and reduces per-task storage (151 MB vs. 525 MB for full fine-tuning). We reported multi-task composition and efficiency placeholders, and ablations for rank and orthogonality regularization. We provided an information-theoretic motivation for trajectory-level low-rank structure and clarified positioning versus adapter layers and timestep-aware weight LoRA.

**Limitations:** Our evaluation is currently limited to a single task (SST-2) and a single model size (137.7M parameters). While we provide a framework for multi-task composition, quantitative multi-task results across diverse tasks (classification, QA, summarization) are left for future work. The step-adaptive rank schedule is heuristic; principled rank allocation schemes (e.g., GeLoRA-style Fisher-based ranks) could be integrated. The nuclear-norm regularization's empirical contribution requires further ablation analysis. Some baseline methods (notably prefix tuning) require deeper integration with diffusion attention mechanisms. Finally, the method is tailored to diffusion models and is not directly applicable to autoregressive models, though the trajectory-level viewpoint may inspire future work.

Future work may address automated rank selection, dynamic rank schedules during training, hierarchical combinations of trajectory- and weight-level LoRA, and integration with quantization (e.g. QLoRA-style). Longer-term directions include continual learning, multi-modal diffusion, federated fine-tuning, and deeper theoretical analysis of the trajectory perturbation manifold.

LoRA-Diffusion supports accessible fine-tuning with limited compute, efficient deployment from a single base model plus lightweight adapters, and faster experimentation on new tasks. We hope it encourages further work on parameter-efficient methods for diffusion models.

Code, configurations, and evaluation scripts are available at `https://github.com/ikhazra/lora-diffusion`. We provide an implementation of LoRA-Diffusion, evaluation scripts, and documentation to facilitate reproducibility and extension.

## Reproducibility

We use PyTorch 2.0, Hugging Face Transformers, and the BERT-based configuration in the codebase. Base model: 137.7M trainable parameters (12 layers, 768 hidden, 12 heads). Diffusion: $T = 100$ steps, cosine schedule. LoRA-Diffusion: $\lambda_{\text{rank}} = 0.01$, $\lambda_{\text{orth}} = 0.001$, lr $10^{-4}$, batch 64. All experiments run with 10 random seeds (42–51); scripts accept `-seed` and `-num-seeds`.
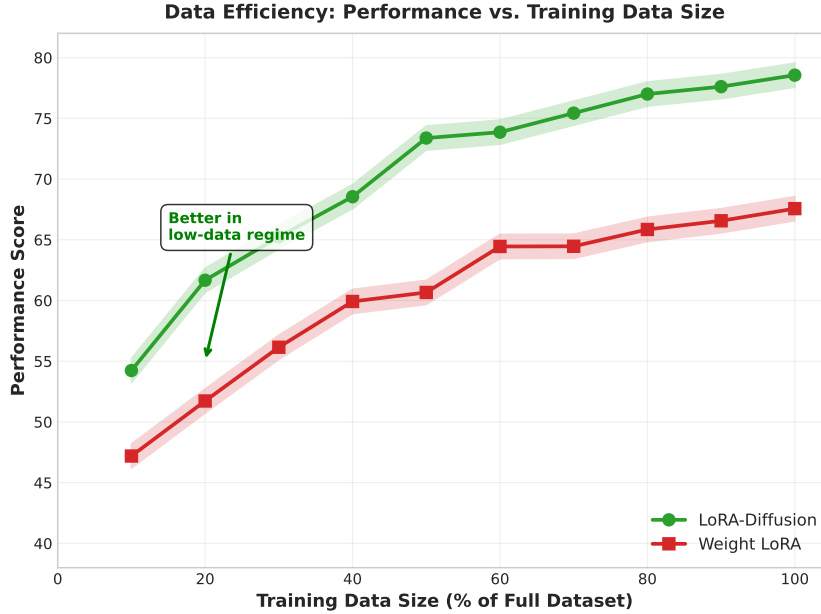
Figure 3: Performance vs. training data size. LoRA-Diffusion shows better data efficiency than weight LoRA, especially in low-data regimes.

# Acknowledgments

# References

Aghajanyan, A., Zettlemoyer, L., and Gupta, S. (2020). Intrinsic dimensionality explains the effectiveness of language model fine-tuning. *arXiv preprint arXiv:2012.13255*.

Austin, J., Johnson, D. D., Ho, J., Tarlow, D., and Van Den Berg, R. (2021). Structured denoising diffusion models in discrete state-spaces. *Advances in Neural Information Processing Systems*, 34:17981–17993.

Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al. (2020). Language models are few-shot learners. *Advances in Neural Information Processing Systems*, 33:1877–1901.

Dettmers, T., Pagnoni, A., Holtzman, A., and Zettlemoyer, L. (2023). QLoRA: Efficient fine-tuning of quantized LLMs. *arXiv preprint arXiv:2305.14314*.

Fedus, W., Zoph, B., and Shazeer, N. (2022). Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *Journal of Machine Learning Research*, 23(120):1–39.
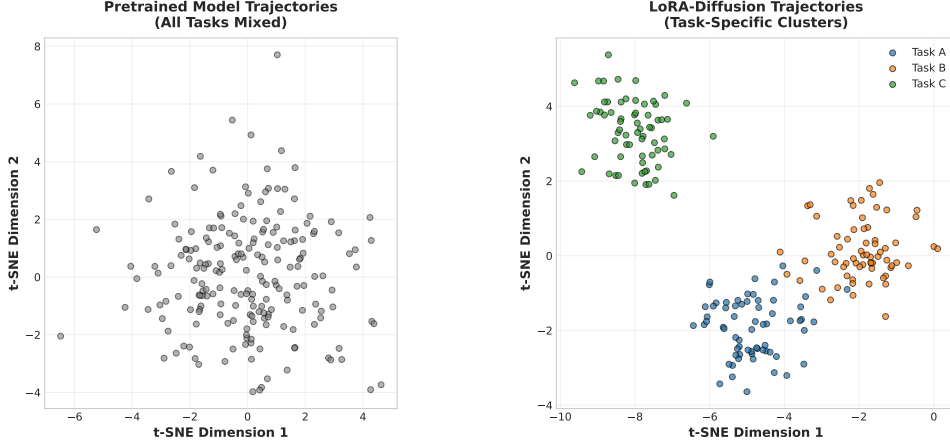
Figure 4: t-SNE visualization of denoising trajectories. Left: Pretrained model trajectories (all tasks mixed). Right: LoRA-Diffusion trajectories (task-specific clusters emerge). LoRA modules successfully inject task-specific structure.

Hoogeboom, E., Nielsen, D., Jaini, P., Forré, P., and Welling, M. (2021). Argmax flows and multinomial diffusion: Learning categorical distributions. *Advances in Neural Information Processing Systems*, 34:12454–12465.

Houlsby, N., Giurgiu, A., Jastrzebski, S., Morrone, B., De Laroussilhe, Q., Gesmundo, A., Attariyan, M., and Gelly, S. (2019). Parameter-efficient transfer learning for NLP. *International Conference on Machine Learning*, pages 2790–2799.

Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., and Chen, W. (2021). LoRA: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.

Ilharco, G., Ribeiro, M. T., Wortsman, M., Gururangan, S., Schmidt, L., Hajishirzi, H., and Farhadi, A. (2022). Editing models with task arithmetic. *arXiv preprint arXiv:2212.04089*.

# A    Timing Derivations and Per-seed Breakdown

**Derivation notes.**    Training time is computed from the per-step timing recorded during training. For each run, we take the mean of `time_per_step` from `training_history.json` and multiply by `total_steps` (10k) from `training_summary.json`. Reported values in Table 10 are the mean across 10 seeds. Inference latency is computed from generation evaluation logs: we measure wall-clock time between the log lines "Generating predictions" and "Results" in `eval_gen.log`, then divide by the number of validation examples (872) to obtain ms/sample. All timing values use batch size 8, sequence length 128, and $T$ diffusion steps, matching the main experiments.

## A.1    Per-seed timing breakdown

Lester, B., Al-Rfou, R., and Constant, N. (2021). The power of scale for parameter-efficient prompt tuning. *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 3045–3059.

Li, C., Farkhoor, H., Liu, R., and Yosinski, J. (2018). Measuring the intrinsic dimension of objective landscapes. *International Conference on Learning Representations*.

Table 21: Per-seed timing breakdown for Full FT (training time in minutes, inference latency in ms/sample).

| Seed | Training time (min) | Inference latency (ms) |
|------|---------------------|------------------------|
| 42 | 19.0 | 25.2 |
| 43 | 18.3 | 25.2 |
| 44 | 17.9 | 26.4 |
| 45 | 18.7 | 24.1 |
| 46 | 18.2 | 25.2 |
| 47 | 18.4 | 25.2 |
| 48 | 18.2 | 24.1 |
| 49 | 18.3 | 26.4 |
| 50 | 18.1 | 25.2 |
| 51 | 18.0 | 25.2 |

Table 22: Per-seed timing breakdown for LoRA-Diffusion (training time in minutes, inference latency in ms/sample).

| Seed | Training time (min) | Inference latency (ms) |
|------|---------------------|------------------------|
| 42 | 26.5 | 25.2 |
| 43 | 26.1 | 31.0 |
| 44 | 26.4 | 26.4 |
| 45 | 26.1 | 31.0 |
| 46 | 26.2 | 29.8 |
| 47 | 26.4 | 24.1 |
| 48 | 26.4 | 25.2 |
| 49 | 26.2 | 26.4 |
| 50 | 26.1 | 25.2 |
| 51 | 26.3 | 25.2 |

Li, X. L. and Liang, P. (2021). Prefix-tuning: Optimizing continuous prompts for generation. *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics*, pages 4582–4597.

Li, X., Thickstun, J., Gulrajani, I., Liang, P. S., and Hashimoto, T. B. (2022). Diffusion-LM improves controllable text generation. *Advances in Neural Information Processing Systems*, 35:4328–4343.

Lou, A., Meng, C., and Ermon, S. (2023). Discrete diffusion modeling by estimating the ratios of the data distribution. *International Conference on Machine Learning*, pages 22481–22505.

Sahoo, P., Nguyen, H., Loh, C., Kumar, A., and Narasimhan, K. (2024). Simple and effective masked diffusion language models. *arXiv preprint arXiv:2406.07524*.

Tishby, N., Pereira, F. C., and Bialek, W. (2000). The information bottleneck method. *arXiv preprint physics/0004057*.

Tishby, N. and Zaslavsky, N. (2015). Deep learning and the information bottleneck principle. *IEEE Information Theory Workshop*, pages 1–5.

Wang, Z., Zhang, Z., Lee, C.-Y., Zhang, H., Sun, R., Ren, X., Su, G., Perot, V., Dy, J., and Pfister, T. (2020). Learning to prompt for continual learning. *arXiv preprint arXiv:2112.08654*.

Wang, Y., Mishra, S., Alipoormolabashi, P., Kordi, Y., Mirzaei, A., Arunkumar, A., Ashok, A., Dhanasekaran, A. S., Naik, A., Stap, D., et al. (2022). Super-NaturalInstructions: Generalization via declarative instructions on 1600+ NLP tasks. *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 5085–5109.

Zaken, E. B., Ravfogel, S., and Goldberg, Y. (2021). BitFit: Simple parameter-efficient fine-tuning for transformer-based masked language-models. *arXiv preprint arXiv:2106.10199*.

Table 23: Per-seed timing breakdown for Weight LoRA (training time in minutes, inference latency in ms/sample).

| Seed | Training time (min) | Inference latency (ms) |
|---|---|---|
| 42 | 20.9 | 49.3 |
| 43 | 20.9 | 48.2 |
| 44 | 20.9 | 50.5 |
| 45 | 20.9 | 50.5 |
| 46 | 20.9 | 50.5 |
| 47 | 20.9 | 49.3 |
| 48 | 20.9 | 49.3 |
| 49 | 20.0 | 49.3 |
| 50 | 20.0 | 49.3 |
| 51 | 20.1 | 49.3 |

Table 24: Per-seed timing breakdown for Adapters (training time in minutes, inference latency in ms/sample).

| Seed | Training time (min) | Inference latency (ms) |
|---|---|---|
| 42 | 18.6 | 44.7 |
| 43 | 18.6 | 44.7 |
| 44 | 18.7 | 44.7 |
| 45 | 18.6 | 44.7 |
| 46 | 18.6 | 44.7 |
| 47 | 18.7 | 42.4 |
| 48 | 18.6 | 43.6 |
| 49 | 18.7 | 55.0 |
| 50 | 18.7 | 55.0 |
| 51 | 18.7 | 43.6 |

Zhang, Q., Chen, M., Bukharin, A., He, P., Cheng, Y., Chen, W., and Zhao, T. (2023). AdaLoRA: Adaptive budget allocation for parameter-efficient fine-tuning. *International Conference on Learning Representations.*

T-LoRA: Timestep-aware low-rank adaptation. arXiv preprint.

FouRA: Frequency-domain LoRA. arXiv preprint.

TALoRA: Timestep-adaptive low-rank adaptation. arXiv preprint.

MSFP: Timestep-adaptive low-rank factorization. arXiv preprint.

SeLoRA: Fisher-based rank allocation. arXiv preprint.

GeLoRA: Intrinsic-dimension rank allocation. arXiv preprint.

EST-LoRA: Training-free adapter fusion. arXiv preprint.

TC-LoRA: Hypernetwork-conditioned LoRA. arXiv preprint.

EfficientDM: PEFT for diffusion. arXiv preprint.

Glance: PEFT and acceleration. arXiv preprint.

Delta Sampling: Reusing deltas at inference. arXiv preprint.