

PARASOLVER-TURBO: ACCELERATING PARALLEL DIFFUSION INTEGRATOR VIA INTRINSIC PARTIALLY LINEAR STRUCTURE

Anonymous authors

Paper under double-blind review

ABSTRACT

This paper explores the challenge of accelerating the sequential inference process of Diffusion Probabilistic Models (DPMs). We tackle this critical issue from a dynamic system perspective, in which the inherent sequential nature is transformed into a parallel sampling process. Specifically, we first reveal that the sequential integral solver of the diffusion model can be approximated by a full linear solver, enabling efficient computation for parallel integral solvers of DPMs. Based on such a linear formulation, we then introduce a unified framework that reformulates the original nonlinear sequential integral process of diffusion model as a system of partial linear equations. Moreover, we further develop an immediate update strategy to solve the system. In addition, we prove that (1) the system admits a unique root corresponding precisely to the trajectory of the sequential integral solver; (2) solving the system guarantees convergence to the trajectory of sequential integral solvers in equal or fewer iterations. Building on these insights, we present *ParaSolver-Turbo*, a partial linear parallel integral solver to accelerate a broad class of sequential and parallel sampling methods such as DDPM and ParaSolver. Extensive experiments validate that ParaSolver-Turbo achieves $2 \times \sim 50 \times$ speedup in terms of wall-clock time without measurable quality degradation. The source code will be released publicly.

1 INTRODUCTION

In recent years, the rise of diffusion probabilistic models (DPMs) (Ho et al., 2020; Song et al., b) has dramatically transformed the field of generative modeling. These models have become a cornerstone technique for a wide range of applications (Yang et al., 2024; Liu et al., 2023a; Lu et al., 2024; Chung et al., 2023; Lu et al., 2023; You et al., 2025; Liu et al., 2025a), from high-fidelity image and video synthesis (Rombach et al., 2022; Blattmann et al., 2023) to molecular generation (Wu et al., 2024; Nguyen et al., 2023). At their core, diffusion models operate through an iterative denoising process, mathematically formulated as an ordinary or stochastic differential equation (ODE/SDE). This process gradually refines noise from an initial Gaussian distribution, transforming it into a realistic sample that aligns with the target data distribution. However, a key limitation of DPMs is their slow sampling speed. Generating high-quality outputs typically requires numerous sequential denoising steps, each involving computationally intensive evaluations of large neural networks.

To accelerate sampling, researchers have proposed various approaches (Gong et al., 2024; Zheng et al., 2023; Gonzalez et al., 2024; Geng et al., 2024; Liu et al., 2023b; Luo et al., 2023). One line of work focuses on advanced SDE/ODE solvers, such as DDIM (Song et al., a) and DPMSolver (Lu et al., 2022), which leverage mathematical optimizations to reduce sampling steps. However, these methods often trade off sample quality for speed. Another direction involves distilling the ODE trajectory of a pre-trained diffusion model into a faster neural network, as seen in (Salimans and Ho; Song et al., 2023; Lipman et al.; Liu et al.). While effective, such techniques frequently suffer from degraded output quality and diversity. To overcome these limitations, recent work has explored parallel sampling strategies. Shih et al. (2024b) introduced a parallel denoising method based on Picard iteration, while Tang et al. (2024) reformulated the problem as solving triangular nonlinear equations via fixed-point iteration. Building on these advances, ParaSolver (Lu et al., 2025) models the sampling process as a banded nonlinear system, exploiting sparsity for greater efficiency than prior dense-structured approaches like ParaTAA. These methods offer three key advantages: (1) they

are training-free and compatible with existing solvers; (2) they preserve sample quality comparable to sequential sampling; (3) they drastically reduce sampling steps, significantly boosting efficiency.

Despite their promise, current parallel methods face two major bottlenecks. First, each iteration requires substantially more neural function evaluations (NFEs) than sequential sampling, causing parallel efficiency to saturate quickly due to GPU compute limits. Second, while iteration counts are reduced compared to sequential sampling, they remain high to constrain maximum speedup. To push these boundaries further, we propose *ParaSolver-Turbo*, a unified framework that generalizes prior work through the lens of nonlinear equations (NEs). Unlike existing methods that rely solely on fully nonlinear systems, we reinterpret sequential sampling as a mixed system of linear and nonlinear equations. This partial linear structure drastically reduces NFE overhead, unlocking new levels of parallel efficiency. Moreover, our framework is highly general—prior methods like ParaDiGMS, ParaTAA, and ParaSolver emerge as special cases within our formulation.

In summary, this paper makes the following theoretical and practical contributions:

- we introduce ParaSolver-Turbo, a novel parallel sampling algorithm for diffusion models, which reformulates the sequential nonlinear diffusion sampling process into a system of partial linear equations;
- we enhance the updating process through an immediate updating strategy, with theoretical guarantees that convergence to the sequential integral solver’s trajectory is achieved within equal or fewer iterations; and
- we experimentally demonstrate that ParaSolver-Turbo achieves a $2.7\times$ – $50.0\times$ speedup in wall-clock time while maintaining comparable generation quality (measured by FID and CLIP scores), setting a new state-of-the-art in this domain.

2 RELATED WORK

Parallel sampling has become a promising technique to speed up the sequential sampling process in DPMs while maintaining the same sample fidelity. By simultaneously denoising multiple sample points along ODE/SDE trajectories, this method iteratively improves a set of initial trajectories, allowing faster generation.

To the best of our knowledge, only a handful of recent works (Shih et al., 2024a; Tang et al., 2024; Lu et al., 2025) have investigated parallel sampling for DPMs. Among these, ParaDiGMS (Shih et al., 2024a) stands as the first parallel sampling method, considering discretized sample points along ODE trajectories as a sequence of fixed points that are refined iteratively using the fixed-point theorem and Picard-Lindelöf theorem. ParaTAA (Tang et al., 2024) builds upon fixed-point iteration by formulating DPM sampling as solving triangular nonlinear equations, where all discretized sample points serve as unknown variables. ParaSolver (Lu et al., 2025) further progresses the field by modeling the sampling process as a banded nonlinear equation system. Its hierarchical parallel strategy exploits the system’s sparsity, achieving considerably higher computational efficiency compared to ParaTAA’s dense triangular structure.

Although these methods demonstrate promising benefits, their parallel efficiency is constrained by two critical factors. First, they require substantially more NFEs per iteration than the sequential sampling, causing parallel efficiency to plateau quickly as NFEs scale—the single GPU’s floating-point operations per second (FLOPS) become the bottleneck. Second, despite achieving a notable reduction in iterations relative to sequential approaches, the remaining iterations are still high, capping the maximum achievable parallel speedup.

In contrast, we reinterpret the sequential sampling of DPMs as a mixed system of linear and nonlinear equations—a distinct departure from existing methods that require fully nonlinear equations. This partial linear structure substantially cuts down the number of NFEs required by current parallel methods, thereby significantly unleashing the parallel efficiency. In addition, our parallel framework is highly general; existing methods such as ParaDiGMS, ParaTAA, and ParaSolver emerge as special cases within our unified paradigm.

Beyond parallel sampling, we also observe other research directions focused on accelerating diffusion models, including model/sample partitioning (Ma et al., 2024b;a; Wang et al., 2024; Li et al., 2024; Liu et al., 2025b), gradient guidance (Wang et al., 2025), trajectory stitching (Pan et al., 2024), nested diffusion (Elata et al., 2024), minimum denoising step prediction (Yu and Barad, 2024), picard consistency model (So et al., 2025), and theoretical guarantees (Gupta et al., 2024; Chen et al.,

2024; Zhou and Sugiyama, 2024). Our work complements these approaches by parallelizing DPMs’ inherently sequential sampling procedure.

3 PRELIMINARY

Sequential Integral Solvers. The diffusion dynamics are governed by a drift coefficient function $f(t)$ and a diffusion coefficient function $g(t)$. The forward diffusion process, which gradually transforms a clean image \mathbf{X}_T into random noise \mathbf{X}_0 (note the convention where clean images correspond to $t = T$ and noise to $t = 0$), is described by the following stochastic differential equation (SDE):

$$d\mathbf{X}_t = f(t)\mathbf{X}_t dt + g(t)d\mathbf{W}, \quad (1)$$

where $d\mathbf{W}$ indicates the standard Wiener process. The generation process requires reversing this diffusion, leading to the reverse-time SDE:

$$d\mathbf{X}_t = \underbrace{[f(t)\mathbf{X}_t - g^2(t)\nabla_{\mathbf{X}_t} \log p(\mathbf{X}_t)]}_{\varphi(\mathbf{X}_t, t)} dt + \underbrace{g(t)}_{\sigma_t} d\mathbf{W}, \quad (2)$$

where $\nabla_{\mathbf{X}_t} \log p(\mathbf{X}_t)$ is approximated by a neural network-based score function $\mathbf{S}_\theta(\cdot)$ with parameters θ . The terms $\varphi(\mathbf{X}_t, t)$ and σ_t parameterize the deterministic and stochastic processes for the reverse process, respectively. Let $\Phi(t, s, \mathbf{X}_s)$ denote the integral solution of Eq. (2) over the time interval $[s, t]$ with initial condition \mathbf{X}_s :

$$\Phi(t, s, \mathbf{X}_s) = \mathbf{X}_s + \int_s^t \varphi(\mathbf{X}_\tau, \tau) d\tau + \int_s^t \sigma_\tau d\mathbf{W}. \quad (3)$$

Parallel Integral Solvers. The ParaSolver framework (Lu et al., 2025) provides a unified approach that encompasses both parallel and sequential sampling methods. Here we introduce existing parallel algorithms within this framework. Consider an array of ODE/SDE trajectories $\{\mathbf{X}_t\}_{t=0}^T$ for DPMs¹. ParaSolver (Lu et al., 2025) proposes a hierarchical parallel strategy that first divides the complete trajectory into N sub-intervals $[t_0, t_1, \dots, t_N]$ with $0 = t_0 < t_1 < \dots < t_N = T$, then formulates a system of banded nonlinear equations:

$$\begin{cases} \hat{\mathbf{X}}_{t_0} - \mathbf{X}_{t_0} = \mathbf{0}, \\ \hat{\mathbf{X}}_{t_{n+1}} - \Phi(t_{n+1}, t_n, \hat{\mathbf{X}}_{t_n}) = \mathbf{0}, \quad n \in \{0, 1, \dots, N-1\}, \end{cases} \quad (4)$$

where $\Phi(t_{n+1}, t_n, \hat{\mathbf{X}}_{t_n})$ represents the local integration operator between consecutive time steps. The framework generalizes various approaches, becoming equivalent to sequential sampling when $N = 1$ and matching full parallel methods like ParaDiGMS (Shih et al., 2024a) and ParaTAA (Tang et al., 2024) when $N = T$.

The algorithm proceeds through iterative refinement of an initial trajectory estimate $\{\hat{\mathbf{X}}_{t_i}\}_{i=0}^N$, typically initialized with coarse approximations or random noise, denoted as $\{\hat{\mathbf{X}}_{t_i}^{(0)}, i = 0, \dots, N\}$. Let $\hat{\mathbf{X}}_{t_0:t_N} = [\hat{\mathbf{X}}_{t_0}^\top, \dots, \hat{\mathbf{X}}_{t_N}^\top]^\top$ denote the concatenated state vector. Then, at the k -th parallel iteration ($k \in [0, K]$), the estimated trajectory is updated according to:

$$\hat{\mathbf{X}}_{t_0:t_N}^{(k+1)} = \hat{\mathbf{X}}_{t_0:t_N}^{(k)} - \mathcal{G}^{(k)} \mathcal{R}_{t_0:t_N}^{(k)}, \quad (5)$$

where $\mathcal{R}_{t_{n+1}}^{(k)} = \hat{\mathbf{X}}_{t_{n+1}}^{(k)} - \Phi(t_{n+1}, t_n, \hat{\mathbf{X}}_{t_n}^{(k)})$ and $\mathcal{R}_{t_0}^{(k)} = \hat{\mathbf{X}}_{t_0}^{(k)} - \mathbf{X}_{t_0} = \mathbf{0}$ indicate residual terms to be minimized. The matrix $\mathcal{G}^{(k)}$ distinguishes various algorithms for solving the above system of nonlinear equations: $\mathcal{G}^{(k)} = \mathbf{I}$ yields fixed-point iteration (Burden and Faires, 1985), while $\mathcal{G}^{(k)} = (\mathcal{J}^{(k)})^{-1}$ corresponds to the Newton-Raphson iteration (Kelley, 2003), with $\mathcal{J}^{(k)} = \partial \mathcal{R}_{t_0:t_N}^{(k)} / \partial \hat{\mathbf{X}}_{t_0:t_N}$ being the Jacobian matrix.

¹For compatibility with parallel sampling, we discretize the continuous time interval $[0, T]$ into discrete steps $[0, 1, \dots, T]$ with unit step size.

4 PROPOSED METHOD

Parallel solvers reformulate the sequential integration process of the diffusion model into a system of non-linear equations (as outlined in Eq. (4) and Eq. (5)). However, parallel solvers demand significantly more model evaluations per iteration to denoise multiple noisy samples compared to sequential solvers, substantially limiting their acceleration performance. To tackle this problem, we propose a linearization method, which reduces the fully nonlinear computations (i.e., the model evaluation operation) of parallel solvers to partially linear operations.

4.1 ρ -NONLINEAR EQUATION SYSTEM

Existing approaches typically develop parallel solvers based on the sequential integral formulation in Eq. (3), which results in a fully nonlinear system of equations in Eq. (4) when employing the nonlinear score function $\mathbf{S}_\theta(\cdot)$. We reveal that the sequential integral solver in Eq. (3) can be expressed in a more elementary form and reduced to a linear solver, as formally stated in the following Proposition 1.

Proposition 1 (Linear Approximation of the Sequential Solver). *The sequential integral solver from Eq. (3) can be decomposed into a linear solver $\mathcal{H}(\cdot)$ with a negligible error term \mathcal{Z} :*

$$\Phi(t, s, \mathbf{X}_s) = \mathcal{H}(t, s, \mathbf{X}_s, \mathbf{X}_T(\mathbf{X}_s)) + \mathcal{Z}, \quad (6)$$

where the linear sequential integral solver $\mathcal{H}(\cdot)$ is linear with respect to \mathbf{X}_s and $\mathbf{X}_T(\mathbf{X}_s)$:

$$\mathcal{H}(t, s, \mathbf{X}_s, \mathbf{X}_T(\mathbf{X}_s)) = \frac{\alpha_t}{\alpha_s} \mathbf{X}_s - \alpha_t \int_{\lambda_s}^{\lambda_t} e^{-\lambda} \frac{\mathbf{X}_\lambda - \alpha_\lambda \mathbf{X}_T(\mathbf{X}_s)}{\sigma_\lambda} d\lambda + \int_{\lambda_s}^{\lambda_t} g(\lambda) d\mathbf{W}, \quad (7)$$

and the error term \mathcal{Z} satisfies $\|\mathcal{Z}\| \leq \zeta \sigma_t \left(\frac{\alpha_t \sigma_s}{\alpha_s \sigma_t} - 1 \right) \simeq 0$.

Proof. See Appendix H for the complete derivation. \square

Remark 1. Proposition 1 reveals that the exact sequential solver in Eq. (3) possesses an intrinsically simpler linear structure.

Now, following ParaSolver (Lu et al., 2025), the linear sequential integral solver in Eq. (7) can be decomposed into many integral functions on different time subintervals, written as

$$\mathbf{X}_{t_{n+1}} = \mathcal{H}(t_{n+1}, t_n, \mathbf{X}_{t_n}, \mathbf{X}_{t_N}), n \in \{0, 1, \dots, N-1\}. \quad (8)$$

Motivated by recent parallel algorithms (Lu et al., 2025; Song et al., 2021), this inherently sequential problem can be addressed in a parallel framework. Specifically, such a sequence of cascaded linear functions can be interpreted as a system of fully linear equations. Denote by $\{\mathbf{X}_{t_n}, n = 0, \dots, N\}$ a set of points exactly on the diffusion trajectory, and $\{\hat{\mathbf{X}}_{t_n}, n = 0, \dots, N\}$ the unknown variables to be optimized towards the trajectory. The cascade of integral formulations as Eq. (8) can then be reformulated as the solutions to the following *fully linear equations*.

Definition 1 (Fully Linear Equation System). *We define the system of fully linear equations for the sequential sampling process in Eq. (8) as*

$$\begin{cases} \hat{\mathbf{X}}_{t_0} - \mathbf{X}_{t_0} = \mathbf{0}, \\ \hat{\mathbf{X}}_{t_{n+1}} - \mathcal{H}(t_{n+1}, t_n, \hat{\mathbf{X}}_{t_n}, \hat{\mathbf{X}}_{t_N}) = \mathbf{0}, n \in \{0, 1, \dots, N-1\}. \end{cases} \quad (9)$$

The complete linearity of the above system permits exceptionally efficient solving with negligible computational cost. However, this very full linearity prevents them from achieving the sample quality of the sequential integral solvers in Eq. (3), as they lack information regarding the target data distribution. To bridge this quality gap, we introduce the target distribution information into the linear system by supplying estimated clean sample $\hat{\mathbf{X}}_{t_N}$. Recall that in the forward process of DPMs, we can compute a predicted clean sample $\hat{\mathbf{X}}_{t_N}$ as follows:

$$\hat{\mathbf{X}}_{t_N}(\hat{\mathbf{X}}_{t_n}, \theta) = \frac{\hat{\mathbf{X}}_{t_n} - \sigma_{t_n} \mathbf{S}_\theta(\hat{\mathbf{X}}_{t_n}, t_n)}{\alpha_{t_n}}. \quad (10)$$

where $\mathbf{X}_{t_n} \sim \mathcal{N}(\alpha_{t_n} \mathbf{X}_{t_N}, \sigma_{t_n} \mathbf{I})$ where α and σ are the noise schedule. Note that the diffusion model consumes most computational resources on the score-matching process, i.e., estimating corresponding noise or clean samples by a neural network. In this paper, to lessen such burden of score matching, we explore the consistency of reverse diffusion process and propose to divide the NEs to a partial-linear equation system, namely ϱ -nonlinear equation system. Specifically, due to the diffusion models make a progressive transition from pure noise to clean data distribution, the score function is highly unstable on the high noise region, which is also evidenced by recent research (Karras et al., 2022; Ma et al., 2024b). Based on such an insight, to further balance the score estimation accuracy and computational consumption, we divide the original NEs shown as Eq. 4 into ϱ -Nonlinear Equation System in Definition 2, where for the unstable and stable score region, we apply the nonlinear and linear solvers, respectively.

Definition 2 (ϱ -Nonlinear Equation System). We define a ϱ -nonlinear equation system for the sequential sampling process in Eq. (8) as

$$\begin{cases} \hat{\mathbf{X}}_{t_0} - \mathbf{X}_{t_0} = \mathbf{0}, \\ \hat{\mathbf{X}}_{t_{n+1}} - \mathcal{H}\left(t_{n+1}, t_n, \hat{\mathbf{X}}_{t_n}, \hat{\mathbf{X}}_{t_N}(\hat{\mathbf{X}}_{t_n}, \theta)\right) = \mathbf{0}, n \in \{0, 1, \dots, \varrho - 1\}, \\ \hat{\mathbf{X}}_{t_{n+1}} - \mathcal{H}\left(t_{n+1}, t_n, \hat{\mathbf{X}}_{t_n}, \hat{\mathbf{X}}_{t_N}(\hat{\mathbf{X}}_{t_{\varrho-1}}, \theta)\right) = \mathbf{0}, n \in \{\varrho, \varrho + 1, \dots, N - 1\}. \end{cases} \quad (11)$$

Remark 2. The scalar ϱ serves as a critical transition point between nonlinear (score) and linear (fixed score) regimes in the system. The first ϱ equations evolve into computationally costly nonlinear equations, reducing into exiting nonlinear parallel system in Eq. (4). Subsequent equations (for $n \geq \varrho$) retain computational efficiency as linear equations, as they depend only on the fixed prediction from the $(\varrho - 1)^{\text{th}}$ time step.

Remark 3. The above framework establishes a unified paradigm, subsuming all existing parallel methods such as ParaDiGMS, ParaTAA, and ParaSolver as special cases. In particular, existing methods are fully nonlinear equation systems, corresponding to the limiting case $\varrho = N - 1$.

Remark 4. Our theoretical analysis demonstrates that the root to the ϱ -nonlinear equation system exists uniquely and provides an unbiased estimator for the sampling results from the sequential integral solver, as established in Proposition 2 and Proposition 3, respectively.

To theoretically validate the equivalence, we define the cumulative state error as $\Delta_{t_n} = \hat{\mathbf{X}}_{t_n} - \mathbf{X}_{t_n}$, representing the deviation between the parallel solution and the precise sequential trajectory. Additionally, based on Definition 2, the the clean sample estimator is explicitly defined as:

$$\mathbf{E}_{t_n} = \begin{cases} \hat{\mathbf{X}}_{t_N}(\hat{\mathbf{X}}_{t_n}, \theta) & \text{if } n < \varrho, \\ \hat{\mathbf{X}}_{t_N}(\hat{\mathbf{X}}_{t_{\varrho-1}}, \theta) & \text{if } n \geq \varrho. \end{cases} \quad (12)$$

Moreover, let $\xi_n = \mathbf{E}_{t_n} - \mathbf{X}_T$ denote the clean sample estimation error.

Proposition 2 (Equivalence under Non-Ideal Conditions). For any $\varrho \in \{1, \dots, N\}$, the solution to the ϱ -nonlinear equation system defined in Eq. (11) is equivalent to the true sequential integral trajectory $\{\mathbf{X}_{t_n}\}_{n=0}^N$. Specifically, the discrepancy is governed by the following error recurrence bound:

$$\|\Delta_{t_{n+1}}\| \leq |\mathcal{A}_n| \|\Delta_{t_n}\| + |\mathcal{B}_n| \|\xi_n\| + \|\mathcal{Z}_n\|, \quad (13)$$

where $\mathcal{A}_n = \frac{\alpha_{t_{n+1}}}{\alpha_{t_n}} \log \frac{\alpha_{t_n} \sigma_{t_{n+1}}}{\alpha_{t_{n+1}} \sigma_{t_n}} \approx 0$ and $\mathcal{B}_n = \alpha_{t_{n+1}} \log \frac{\alpha_{t_{n+1}} \sigma_{t_n}}{\alpha_{t_n} \sigma_{t_{n+1}}} \cdot \log \left(\sqrt{\frac{\alpha_{t_n} \sigma_{t_{n+1}}}{\alpha_{t_{n+1}} \sigma_{t_n}}} \right) \approx 0$. Since $\|\Delta_{t_0}\| = 0$, the coefficients \mathcal{A}_n and \mathcal{B}_n approach negligible values (≈ 0), and combined with $\|\mathcal{Z}_n\| \approx 0$, the total error $\|\Delta_{t_{n+1}}\|$ remains negligible throughout the sampling process.

Proof. See Appendix I for the complete derivation. \square

Remark 5. Prop. 2 guarantes that different ϱ does not impact the quality of the generated samples. Our experiments also confirm this, showing that varying ϱ does not degrade the FID and CLIP scores of the generated sample across various diffusion processes like variance preserving (VP) and variance exploding (VE). We note that ϱ is less of a hyperparameter and more of a hardware configuration parameter. For good efficiency, ϱ can simply be set to the number of available processors.

270 5 SOLVING THE ϱ -NONLINEAR EQUATION SYSTEM

271
272 We investigate the update rule in Eq. (5) to solve the prescribed ϱ -nonlinear system. Specializing
273 Eq. (5) to our system yields the following iterative scheme:

$$274 \hat{\mathbf{X}}_{t_0:t_N}^{(k+1)} = \hat{\mathbf{X}}_{t_0:t_N}^{(k)} - \mathcal{G}^{(k)} \mathcal{R}_{t_0:t_N}^{(k)}, \quad (14)$$

275
276 where the residual term $\mathcal{R}_{t_{n+1}}^{(k)}$ with $\mathcal{R}_{t_0}^{(k)} = \hat{\mathbf{X}}_{t_0}^{(k)} - \mathbf{X}_{t_0} = \mathbf{0}$ is defined as:

$$277 \mathcal{R}_{t_{n+1}}^{(k)} = \begin{cases} \hat{\mathbf{X}}_{t_{n+1}}^{(k)} - \mathcal{H}(t_{n+1}, t_n, \hat{\mathbf{X}}_{t_n}^{(k)}, \hat{\mathbf{X}}_{t_N}(\hat{\mathbf{X}}_{t_n}^{(k)}, \theta)), & n \in \{0, 1, \dots, \varrho - 1\}, \\ \hat{\mathbf{X}}_{t_{n+1}}^{(k)} - \mathcal{H}(t_{n+1}, t_n, \hat{\mathbf{X}}_{t_n}^{(k)}, \hat{\mathbf{X}}_{t_N}(\hat{\mathbf{X}}_{t_{\varrho-1}}^{(k)}, \theta)), & n \in \{\varrho, \varrho + 1, \dots, N - 1\}. \end{cases} \quad (15)$$

281 Observe that this update rule applies a Jacobi-type iteration, where all components $\mathbf{X}_{t_1:t_N}^{(k+1)}$ are up-
282 dated exclusively from the old components $\mathbf{X}_{t_1:t_N}^{(k)}$, resulting in delayed information propagation
283 across time steps. Inspired by the classical Gauss-Seidel iteration, we propose a more efficient im-
284 mediate update strategy wherein newly updated components are incorporated immediately into sub-
285 sequent components' updating within the same iteration. Specifically, we first update the nonlinear-
286 region components $\mathbf{X}_{t_n}^{(k)}$ for $n \in \{1, \dots, \varrho\}$. Upon obtaining the refined nonlinear component
287 $\mathbf{X}_{t_\varrho}^{(k+1)}$, we instantly propagate this refinement forward to compute the remaining linear-region
288 components (i.e., $\mathbf{X}_{t_{\varrho+1}}^{(k+1)}, \dots, \mathbf{X}_{t_N}^{(k+1)}$). This approach ensures faster convergence by propagating
289 the latest information forward, leading to the following improved update rule:

$$290 \hat{\mathbf{X}}_{t_0:t_N}^{(k+1)} = \hat{\mathbf{X}}_{t_0:t_N}^{(k)} - \mathcal{G}^{(k)} \mathcal{R}_{t_0:t_N}^{(k)}, \quad (16)$$

291
292 where the residual term $\mathcal{R}_{t_{n+1}}^{(k)}$ now incorporates the latest refined components:

$$293 \mathcal{R}_{t_{n+1}}^{(k)} = \begin{cases} \hat{\mathbf{X}}_{t_{n+1}}^{(k)} - \mathcal{H}(t_{n+1}, t_n, \hat{\mathbf{X}}_{t_n}^{(k)}, \hat{\mathbf{X}}_{t_N}(\hat{\mathbf{X}}_{t_n}^{(k)}, \theta)), & n \in \{0, 1, \dots, \varrho - 1\}, \\ \hat{\mathbf{X}}_{t_{n+1}}^{(k)} - \mathcal{H}(t_{n+1}, t_n, \hat{\mathbf{X}}_{t_n}^{(k+1)}, \hat{\mathbf{X}}_{t_N}(\hat{\mathbf{X}}_{t_{\varrho-1}}^{(k)}, \theta)), & n \in \{\varrho, \varrho + 1, \dots, N - 1\}. \end{cases} \quad (17)$$

294
295 We observe that once the nonlinear components converge, the predicted clean sample $\hat{\mathbf{X}}_{t_N}(\hat{\mathbf{X}}_{t_{\varrho-1}}^{(k)}, \theta)$
296 becomes unchanged, implying that no further precise information can be propagated for the linear
297 equations in subsequential iterations. To address this, we progressively convert the linear equations
298 into nonlinear ones, such that the number of nonlinear equations remains unchanged as the iteration
299 proceeds. Formally, at iteration k , suppose the unconverged components are $\hat{\mathbf{X}}_{t_n}^{(k)}, \dots, \hat{\mathbf{X}}_{t_N}^{(k)}$, we
300 maintain nonlinearity in the residuals $\{\mathcal{R}_{t_{i+1}}^{(k)}, i = n, n + 1, \dots, n + \varrho - 1\}$ of the first ϱ components
301 through the following formulation:

$$302 \hat{\mathbf{X}}_{t_n:t_N}^{(k+1)} = \hat{\mathbf{X}}_{t_n:t_N}^{(k)} - \mathcal{G}^{(k)} \mathcal{R}_{t_n:t_N}^{(k)}, \quad (18)$$

303
304 where the residual term $\mathcal{R}_{t_{i+1}}^{(k)}$ can adaptively convert the linear equations into non-linear ones:

$$305 \mathcal{R}_{t_{i+1}}^{(k)} = \begin{cases} \hat{\mathbf{X}}_{t_{i+1}}^{(k)} - \mathcal{H}(t_{i+1}, t_i, \hat{\mathbf{X}}_{t_i}^{(k)}, \hat{\mathbf{X}}_{t_N}(\hat{\mathbf{X}}_{t_i}^{(k)}, \theta)), & i \in \{n, n + 1, \dots, n + \varrho - 1\}, \\ \hat{\mathbf{X}}_{t_{i+1}}^{(k)} - \mathcal{H}(t_{i+1}, t_i, \hat{\mathbf{X}}_{t_i}^{(k+1)}, \hat{\mathbf{X}}_{t_N}(\hat{\mathbf{X}}_{t_{n+\varrho-1}}^{(k)}, \theta)), & i \in \{n + \varrho, n + \varrho + 1, \dots, N - 1\}. \end{cases} \quad (19)$$

306
307 where $\mathcal{R}_{t_n}^{(k)} = \hat{\mathbf{X}}_{t_n}^{(k)} - \mathbf{X}_{t_n} = \mathbf{0}$. This approach ensures continuous refinement of $\hat{\mathbf{X}}_{t_N}(\hat{\mathbf{X}}_{t_{n+\varrho-1}}^{(k)}, \theta)$,
308 yielding increasingly accurate solutions for the linear equations.

309
310 **Proposition 3** (Convergence Analysis). *The update rule in Eq. (18) achieves exact convergence to*
311 *the generated sample from the sequential integral solver within $K \leq T$ iterations.*

312
313 *Proof.* See Appendix J for the complete derivation. \square

314
315 **Stopping Criterion.** Following existing works (Shih et al., 2024a; Tang et al., 2024; Lu et al., 2025),
316 We define the stopping criterion as $\frac{1}{D} \left\| \hat{\mathbf{X}}_{t_n}^{(k+1)} - \hat{\mathbf{X}}_{t_n}^{(k)} \right\|_F^2 \leq \delta^2 \sigma_{t_n}^2$, where D is the dimensionality
317 of the sample $\hat{\mathbf{X}}_{t_n}$, $\|\cdot\|_F$ denotes the Frobenius norm, and δ is an tolerance parameter.

Algorithm 1: ParaSolver-Turbo: a partial nonlinear parallel integral solver for diffusion models

Input : Diffusion model S_θ , noise schedule α_t and σ_t , subinterval number N , preconditioning steps M , tolerance δ , batch window size p , sample dimension D , nonlinear degree ϱ , maximum parallel step K .

Output : A sample.

```

329 1 Initialize  $\{\hat{\mathbf{X}}_{t_n}^{(0)}, n = 0, \dots, p\}$  by Eq. (20) // Initialize with a few sampling steps.
330 2  $n, k \leftarrow 0, 0, k \in [0, K]$ , and  $n \in [0, N - 1]$ 
331 3 while  $n < N$  do
332 4    $\hat{\mathbf{X}}_{t_N}(\hat{\mathbf{X}}_{t_i}^{(k)}, \theta) = \frac{\hat{\mathbf{X}}_{t_i}^{(k)} - \sigma_{t_i} \mathbf{S}_\theta(\hat{\mathbf{X}}_{t_i}^{(k)}, t_i)}{\alpha_{t_i}}, \forall i \in \{n, \dots, n + \varrho - 1\}$ . // Predict clean samples in parallel.
333 5    $\mathcal{R}_{t_{i+1}}^{(k)} = \hat{\mathbf{X}}_{t_{i+1}}^{(k)} - \mathcal{H}(t_{i+1}, t_i, \hat{\mathbf{X}}_{t_i}^{(k)}, \hat{\mathbf{X}}_{t_N}(\hat{\mathbf{X}}_{t_i}^{(k)}, \theta)), \forall i \in \{n, \dots, n + \varrho - 1\}$  // Nonlinear residual.
334 6    $\hat{\mathbf{X}}_{t_{n+1}:t_{n+\varrho}}^{(k+1)} = \hat{\mathbf{X}}_{t_{n+1}:t_{n+\varrho}}^{(k)} - \mathcal{G}^{(k)} \mathcal{R}_{t_{n+1}:t_{n+\varrho}}^{(k)}$  // Update the nonlinear components.
335 7   for  $i \in \{n + \varrho, n + \varrho + 1, \dots, n + p - 1\}$  do
336 8      $\mathcal{R}_{t_{i+1}}^{(k)} = \hat{\mathbf{X}}_{t_{i+1}}^{(k)} - \mathcal{H}(t_{i+1}, t_i, \hat{\mathbf{X}}_{t_i}^{(k+1)}, \hat{\mathbf{X}}_{t_N}(\hat{\mathbf{X}}_{t_{n+\varrho-1}}^{(k)}, \theta))$  in Eq. (17) // Linear residual.
337 9      $\hat{\mathbf{X}}_{t_{i+1}}^{(k+1)} = \hat{\mathbf{X}}_{t_{i+1}}^{(k)} - \mathcal{G}^{(k)} \mathcal{R}_{t_{i+1}}^{(k)}$  // Update the linear components using immediate updating
338 // strategy.
339 10 Compute errors  $\mathcal{E}_{t_{i+1}} = \frac{1}{D} \|\mathcal{R}_{t_{i+1}}^{(k)}\|_F^2, \forall i \in \{n, \dots, n + p - 1\}$  // Compute error for each component.
340 11  $s \leftarrow \min_{i+1} (\{i + 1 | \mathcal{E}_{t_{i+1}} > \delta^2 \sigma_{t_{i+1}}^2, \forall i \in \{n, \dots, n + p - 1\}\} \cup \{n + p\})$  // The sliding stride.
341 12  $\hat{\mathbf{X}}_{t_{i+1}}^{(k+1)} = \mathcal{H}(t_{i+1}, t_i, \hat{\mathbf{X}}_{t_i}^{(k+1)}, \hat{\mathbf{X}}_{t_N}(\hat{\mathbf{X}}_{t_{n+\varrho-1}}^{(k)}, \theta)), \forall i \in \{n + p, \dots, n + p + s - 1\}$  // Initialize
342 // new components for the next window.
343 13  $n \leftarrow n + s, k \leftarrow k + 1, p \leftarrow \min(p, N - n)$  // Update the iteration indexes.

```

Return: $\hat{\mathbf{X}}_{t_N}^{(K)}$

Initialization. To establish the diffusion trajectory, we adopt the ParaSolver methodology by performing M steps of a sequential integral solver as preconditioning steps. The initialization process is formally expressed as:

$$\begin{cases} \hat{\mathbf{X}}_{t_n}^{(0)} = \mathbf{X}_{t_n}, & \text{if } 0 \leq n < M, \\ \hat{\mathbf{X}}_{t_n}^{(0)} = \mathcal{H}(t_n, t_{n-1}, \hat{\mathbf{X}}_{t_{n-1}}^{(0)}, \hat{\mathbf{X}}_{t_N}(\hat{\mathbf{X}}_{t_{M-1}}^{(0)}, \theta)), & \text{if } M \leq n < N, \end{cases} \quad (20)$$

where the initial state \mathbf{X}_{t_0} is sampled from a standard Gaussian distribution, $\mathbf{X}_{t_0} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. When $M = 0$, following Shih et al. (2024a), all samples are initialized with random noise.

Sliding Window Strategy. To optimize computational efficiency, ParaSolver utilizes a sliding window of fixed size p , retaining only the most recent p denoised points $\{\hat{\mathbf{X}}_{t_n}^{(k)} : n = 0, \dots, p - 1\}$ for parallel processing. This approach enhances parallel performance while adhering to GPU memory limitations. Adapted to our framework, the update rule becomes:

$$\hat{\mathbf{X}}_{t_n:t_{n+p}}^{(k+1)} = \hat{\mathbf{X}}_{t_n:t_{n+p}}^{(k)} - \mathcal{G}^{(k)} \mathcal{R}_{t_n:t_{n+p}}^{(k)}, \quad (21)$$

where the residual term $\mathcal{R}_{t_{i+1}}^{(k)}$ over a sliding window becomes:

$$\mathcal{R}_{t_{i+1}}^{(k)} = \begin{cases} \hat{\mathbf{X}}_{t_{i+1}}^{(k)} - \mathcal{H}(t_{i+1}, t_i, \hat{\mathbf{X}}_{t_i}^{(k)}, \hat{\mathbf{X}}_{t_N}(\hat{\mathbf{X}}_{t_i}^{(k)}, \theta)), & i \in \{n, \dots, n + \varrho - 1\}, \\ \hat{\mathbf{X}}_{t_{i+1}}^{(k)} - \mathcal{H}(t_{i+1}, t_i, \hat{\mathbf{X}}_{t_i}^{(k+1)}, \hat{\mathbf{X}}_{t_N}(\hat{\mathbf{X}}_{t_{n+\varrho-1}}^{(k)}, \theta)), & i \in \{n + \varrho, \dots, p - 1\}. \end{cases} \quad (22)$$

Complete Algorithm. Algorithm 1 outlines the full procedure of the proposed ParaSolver-Turbo. Initially, an array of starting guesses $\{\hat{\mathbf{X}}_{t_n}^{(0)} : n = 0, \dots, p\}$ is prepared either through random noises or inexpensive preconditioning steps (Line 1). The parallel sampling loop then begins at Line 2. In Lines 4–6, all nonlinear equations are solved simultaneously. Subsequently, Lines 7–9 compute the solutions to the linear equations, leveraging the predicted clean samples obtained from the nonlinear solving phase. This step is computationally efficient since it avoids model evaluations. The discrepancy between the newly computed states and the current states is assessed at Line 10, followed by determining the feasible window sliding stride at Line 11. Line 12 initializes s new points beyond the current window based on the computed stride.

Table 1: Comparisons on Stable Diffusion-v2, with classifier guidance $w = 7.5$. The visual comparisons are shown in *Appendix*.

| Steps | Method | Stable Diffusion-v2 | | | | | |
|-------|------------------------------|---------------------|------------|-------|------|------------|--------------|
| | | Iters ↓ | NFE ↓ | CLIP↑ | FID↓ | Time (s)↓ | Speedup↑ |
| 1000 | DDPM | 1000 | 1000 | 25.6 | 55.9 | 128.0 | 1.0× |
| | DDPM + ParaDiGMS | 65 | 2024 | 25.6 | 55.6 | 32.8 | 3.9× |
| | DDPM + ParaSolver | 32 | 1065 | 25.6 | 55.3 | 15.9 | 8.1× |
| | DDPM + ParaSolver-Turbo | 32 | 195 | 25.6 | 55.7 | 3.2 | 40.0× |
| 50 | DDIM | 50 | 50 | 25.6 | 57.2 | 6.3 | 1.0× |
| | DDIM + ParaDiGMS | 21 | 132 | 25.6 | 56.9 | 3.3 | 1.9× |
| | DDIM + ParaSolver | 13 | 83 | 25.6 | 56.9 | 1.9 | 3.3× |
| | DDIM + ParaSolver-Turbo | 13 | 73 | 25.6 | 57.1 | 1.6 | 3.9× |
| 25 | DDIM | 25 | 25 | 25.4 | 62.9 | 3.4 | 1.0× |
| | DDIM + ParaDiGMS | 18 | 53 | 25.4 | 62.8 | 3.2 | 1.2× |
| | DDIM + ParaSolver | 11 | 49 | 25.4 | 61.9 | 1.2 | 2.8× |
| | DDIM + ParaSolver-Turbo | 11 | 42 | 25.4 | 61.8 | 1.1 | 3.1× |
| 50 | DPMSolver | 50 | 50 | 25.6 | 57.2 | 6.3 | 1.0× |
| | DPMSolver + ParaDiGMS | 25 | 132 | 25.6 | 57.2 | 3.3 | 1.8× |
| | DPMSolver + ParaSolver | 15 | 96 | 25.6 | 57.1 | 2.2 | 3.1× |
| | DPMSolver + ParaSolver-Turbo | 14 | 93 | 25.6 | 57.1 | 1.7 | 3.9× |
| 25 | DPMSolver | 25 | 25 | 25.4 | 62.3 | 3.2 | 1.0× |
| | DPMSolver + ParaDiGMS | 15 | 81 | 25.4 | 62.3 | 2.2 | 1.5× |
| | DPMSolver + ParaSolver | 11 | 58 | 25.4 | 62.1 | 1.5 | 2.1× |
| | DPMSolver + ParaSolver-Turbo | 10 | 55 | 25.4 | 62.1 | 1.2 | 2.7× |

6 EXPERIMENTS

Evaluation Metrics. We assessed ParaSolver-Turbo using five widely adopted metrics: function evaluation count (NFE), iterations (iters), wall-clock time, FID score (Heusel et al., 2017), and CLIP score (Hessel et al., 2021).

Datasets and Models. As per ParaSolver (Lu et al., 2025), we analyze our ParaSolver-Turbo across latent-space model (StableDiffusion-v2) and pixel-space model (LSUN Church) using the same settings.

Algorithms. We apply our approach to accelerate the performance of the state-of-the-art sequential sampling methods: DDPM (Ho et al., 2020), DDIM (Song et al., a), and DPMSolver (Lu et al., 2022). We evaluate the parallel efficiency of our ParaSolver-Turbo against ParaSolver and ParaDiGMS in accelerating the aforementioned three sequential methods as it incorporates all existing parallel methods.

Hyperparameter Settings. Following the previous settings in (Lu et al., 2025; Tang et al., 2024), we apply our ParaSolver-Turbo and ParaDiGMS to DDPM with 1000 sequential sampling steps. For DDIM and DPMSolver, we consider two settings: 25 and 50 sequential sampling steps. Besides, We follow the best tolerances for ParaDiGMS and ParaSolver in their papers. For our ParaSolver-Turbo implementation, we optimized the hyperparameters through grid search. On the StableDiffusion-v2 model, we set $\delta = 0.5$ for DDPM, while using $\delta = 0.03$ for both DDIM and DPMSolver. For the LSUN Church model, we use $\delta = 0.5$ for DDPM, and $\delta = 0.003$ for DDIM. Additionally, we set $\rho = 8$ for ParaSolver-Turbo, as we utilize 8 NVIDIA 3090 GPUs. More details are provided in the *Appendix D*.

Acceleration Comparison. Table 1 shows the experimental results on DDPM, DDIM, DPMSolver, and their parallel variants when combined with ParaSolver, ParaDiGMS, and our ParaSolver-Turbo, respectively. ParaSolver-Turbo demonstrates significant acceleration for DDPM, DDIM, and DPMSolver while maintaining comparable FID and CLIP scores. Notably, it achieves a remarkable wall-clock time speedup of up to 40.0×, setting a new benchmark for parallel diffusion methods. Although parallelization inherently trades computational efficiency for speed, leading to increased NFEs relative to sequential approaches, our method capitalizes on partial linearity to substantially

Table 2: The effect of the number of GPUs using LSUN Church. We report Batch Size for feeding into the model per GPU (BS), GPU Utilization (Util), and Speedup (Spd) for 1000 and 50 diffusion steps. We set $\varrho = 8$ for our method and $p = 24$ for all methods.

| Step | Method | 1 GPU | | | 4 GPUs | | | 8 GPUs | | |
|------|------------|-------|------|-------------------------------|--------|------|--------------------------------|--------|------|--------------------------------|
| | | BS | Util | Spd | BS | Util | Spd | BS | Util | Spd |
| 1000 | DDPM | 1 | 19% | 1.0 \times | 1 | 19% | 1.0 \times | 1 | 19% | 1.0 \times |
| | ParaDiGMS | 24 | 100% | 1.4 \times | 6 | 48% | 3.3 \times | 3 | 27% | 4.6 \times |
| | ParaSolver | 24 | 100% | 3.1 \times | 6 | 46% | 6.9 \times | 3 | 25% | 12.1 \times |
| | Ours | 8 | 59% | 8.3\times | 2 | 25% | 16.3\times | 1 | 19% | 43.7\times |
| 50 | DDPM | 1 | 19% | 1.0 \times | 1 | 19% | 1.0 \times | 1 | 19% | 1.0 \times |
| | ParaDiGMS | 24 | 98% | 0.4 \times | 6 | 41% | 0.7 \times | 3 | 23% | 1.6 \times |
| | ParaSolver | 24 | 98% | 0.8 \times | 6 | 40% | 1.1 \times | 3 | 24% | 2.9 \times |
| | Ours | 8 | 43% | 1.7\times | 2 | 24% | 2.6\times | 1 | 19% | 3.8\times |

reduce NFEs, thereby unleashing the parallel efficiency. Similarly, ParaSolver-Turbo consistently surpasses existing methods, delivering speed improvements of 2.7 \times to **50.0 \times** on LSUN Church model. The experimental results on the LSUN Church model are summarized in the *Appendix*.

Generalization on VP, VE, and EDM Process. We investigated the impact of the VP, VE, and EDM diffusion processes by converting the sequential sampler in Alg. 1 of work (Karras et al., 2022) into our parallel version. Adopting the Tab. 1 configurations specified in work (Karras et al., 2022), we then analyzed the effect of the nonlinear degree ϱ on FID performance across these three processes using the provided ImageNet dataset. As shown in Tab. 3, our parallel sampling method demonstrates its applicability and effectiveness across all three diffusion processes. For each process, our parallel approach consistently reduces the number of sampling iterations while achieving comparable or improved FID scores compared to the sequential baseline. This highlights the robustness and efficiency of our method, making it a viable accelerator for various diffusion model formulations. We notice that our method yields a better FID score in the VP process. This improvement likely stems from the inherently stiff nature of the VP-ODE, combined with the step-parallel updates of our parallel solver. The step-parallel updates sever the temporal dependence chain across steps, thus effectively curtail the accumulation of truncation error across steps (Karras et al., 2022). As a result, our solver acts as an implicit numerical stabilizer for the stiff VP-ODE, which typically exhibits significant truncation error that degrades the performance of conventional VP samplers.

Table 3: The impact of nonlinear degree ϱ on FID across EDM, VE, and VP processes over the ImageNet dataset. We set $\delta = 0.01$ for the EDM and VP processes, and $\delta = 0.1$ for the VE process.

| (a) EDM Process ($\delta = 0.01$) | | | (b) VE Process ($\delta = 0.1$) | | | (c) VP Process ($\delta = 0.01$) | | |
|-------------------------------------|---------|-------|-----------------------------------|----------|-------|------------------------------------|---------|-------|
| Method | FID | Iters | Method | FID | Iters | Method | FID | Iters |
| Baseline | 11.1649 | 25 | Baseline | 423.0628 | 25 | Baseline | 104.15 | 25 |
| $\varrho = 4$ | 10.8811 | 22 | $\varrho = 4$ | 422.9794 | 14 | $\varrho = 4$ | 35.8963 | 12 |
| $\varrho = 6$ | 10.8718 | 21 | $\varrho = 6$ | 422.9293 | 14 | $\varrho = 6$ | 41.1196 | 10 |
| $\varrho = 8$ | 10.8621 | 20 | $\varrho = 8$ | 422.8247 | 13 | $\varrho = 8$ | 58.5588 | 8 |
| $\varrho = 10$ | 10.8626 | 15 | $\varrho = 10$ | 422.7672 | 10 | $\varrho = 10$ | 58.5215 | 7 |
| $\varrho = 12$ | 10.8618 | 12 | $\varrho = 12$ | 422.7651 | 9 | $\varrho = 12$ | 58.8340 | 7 |

GPU Scaling Analysis. Table 2 examines GPU scaling performance. Our method achieves 43.7 \times speedup, significantly outperforming existing approaches. Notably, our approach delivers outstanding performance even on a single GPU, achieving 8.3 \times acceleration. This single-GPU efficiency offers significant practical advantages, enabling high-speed generation for resource-constrained sce-

486 narios. The efficiency stems from our linearized sampling process that solves only a partial nonlinear
 487 system.

488 **Nonlinear Degree ϱ Analysis:** Tab. 4 conducted an ablation
 489 study by accelerating a 50-step DPMSolver on the Stable
 490 Diffusion v2 model using 8 GPUs. As shown in the
 491 table, the FID and CLIP scores remain remarkably stable
 492 across a wide range of ϱ values. This stability confirms
 493 our theoretical finding in Prop. 2 that ϱ does not compro-
 494 mise the fidelity of the generated samples. Furthermore,
 495 the results reveal a clear trade-off between ϱ and efficiency.
 496 While increasing ϱ initially yields significant speedups, the
 497 gains diminish and eventually reverse. This is because a
 498 larger ϱ requires more NFEs, as seen in the table. The in-
 499 creased NFE count results in a greater computational over-
 500 head for each GPU (proportional to NFEs / GPUs), which
 501 leads to a rapid increase in runtime after an optimal point
 502 (around $\varrho = 15$).

Table 4: Impact of varying ϱ on generation quality (FID, CLIP) and performance (NFEs, Speedup).

| Method | FID | CLIP | Iters | NFEs | Speedup |
|----------------|-------|-------|-------|------|---------|
| Baseline | 57.23 | 25.64 | 50 | 50 | 1× |
| $\varrho = 5$ | 57.18 | 25.55 | 18 | 82 | 2.9× |
| $\varrho = 10$ | 57.25 | 25.59 | 12 | 101 | 3.0× |
| $\varrho = 15$ | 57.14 | 25.61 | 11 | 123 | 3.1× |
| $\varrho = 20$ | 57.22 | 25.60 | 12 | 160 | 2.6× |
| $\varrho = 25$ | 57.14 | 25.53 | 12 | 202 | 2.1× |
| $\varrho = 30$ | 57.06 | 25.58 | 12 | 223 | 2.1× |
| $\varrho = 35$ | 57.18 | 25.62 | 12 | 245 | 1.8× |
| $\varrho = 40$ | 57.15 | 25.61 | 12 | 264 | 1.8× |
| $\varrho = 45$ | 57.16 | 25.59 | 12 | 275 | 1.6× |
| $\varrho = 50$ | 57.14 | 25.57 | 12 | 284 | 1.6× |

503 7 CONCLUSION

504 In this work, we have presented ParaSolver-Turbo, a uni-
 505 fied framework that reinterprets diffusion model sampling as a mixed system of linear and nonlinear
 506 equations. By exploiting the partial linear structure inherent in sequential sampling, ParaSolver-
 507 Turbo significantly reduces the computational overhead of prior parallel approaches while preserv-
 508 ing sample quality. Theoretical analysis guarantees convergence to the sequential solver’s trajectory
 509 with equal or fewer iterations, and extensive experiments demonstrate 2×–50× wall-clock time
 510 speedups without compromising fidelity.

486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

ETHICS & REPRODUCIBILITY STATEMENTS

Our work focuses on a fundamental algorithmic improvement for sampling from generative models and does not introduce new ethical concerns beyond those already associated with large-scale text-to-image models. The models used in our experiments are developed by third parties, and we use them as is. Our method could be used to accelerate the generation of harmful content, but it does not inherently make such generation easier or more likely than with standard samplers. We believe the primary positive impact is making high-fidelity generative AI more accessible for research and creative applications by lowering the inference time barrier. For reproducibility, we have detailed our methodology, algorithm, and experimental setup in the paper. We will release our source code, built upon standard open-source libraries, upon publication to allow for full verification of our results.

REFERENCES

- Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion: Scaling latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*, 2023.
- William E Boyce, Richard C DiPrima, and Douglas B Meade. *Elementary differential equations and boundary value problems*. John Wiley & Sons, 2021.
- Richard L Burden and J Douglas Faires. Fixed-point iteration. *Numerical Analysis (3rd ed.)*. PWS Publishers, page 64, 1985.
- Chi-Tsong Chen. *Linear system theory and design*. Saunders college publishing, 1984.
- Haoxuan Chen, Yinuo Ren, Lexing Ying, and Grant M Rotskoff. Accelerating diffusion models with parallel sampling: Inference at sub-linear time complexity. *arXiv preprint arXiv:2405.15986*, 2024.
- Hyungjin Chung, Jeongsol Kim, Sehui Kim, and Jong Chul Ye. Parallel diffusion models of operator and image for blind inverse problems. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6059–6069, 2023.
- Noam Elata, Bahjat Kawar, Tomer Michaeli, and Michael Elad. Nested diffusion processes for anytime image generation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 5018–5027, 2024.
- Zhengyang Geng, Ashwini Pokle, and J Zico Kolter. One-step diffusion distillation via deep equilibrium models. *Advances in Neural Information Processing Systems*, 36, 2024.
- Mingming Gong, Shaoan Xie, Wei Wei, Matthias Grundmann, Kayhan Batmanghelich, Tingbo Hou, et al. Semi-implicit denoising diffusion models (siddms). *Advances in Neural Information Processing Systems*, 36, 2024.
- Martin Gonzalez, Nelson Fernandez Pinto, Thuy Tran, Hatem Hajri, Nader Masmoudi, et al. Seeds: Exponential sde solvers for fast high-quality sampling from diffusion models. *Advances in Neural Information Processing Systems*, 36, 2024.
- Shivam Gupta, Linda Cai, and Sitan Chen. Faster diffusion-based sampling with randomized mid-points: Sequential and parallel. *arXiv preprint arXiv:2406.00924*, 2024.
- Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A reference-free evaluation metric for image captioning. *arXiv preprint arXiv:2104.08718*, 2021.
- Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. In *Proc. NeurIPS*, 2022.

- 594 Carl T Kelley. *Solving nonlinear equations with Newton's method*. SIAM, 2003.
595
- 596 Muyang Li, Tianle Cai, Jiaxin Cao, Qinsheng Zhang, Han Cai, Junjie Bai, Yangqing Jia, Kai Li,
597 and Song Han. Distrifusion: Distributed parallel inference for high-resolution diffusion models.
598 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages
599 7183–7193, 2024.
- 600 Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matthew Le. Flow match-
601 ing for generative modeling. In *The Eleventh International Conference on Learning Representa-*
602 *tions*.
603
- 604 Kendong Liu, Zhiyu Zhu, Chuanhao Li, Hui Liu, Huanqiang Zeng, and Junhui Hou. Prefpaint:
605 Aligning image inpainting diffusion model with human preference. *Advances in Neural Informa-*
606 *tion Processing Systems*, 37:30554–30589, 2025a.
- 607 Xingchao Liu, Chengyue Gong, et al. Flow straight and fast: Learning to generate and transfer data
608 with rectified flow. In *The Eleventh International Conference on Learning Representations*.
609
- 610 Xingchao Liu, Xiwen Zhang, Jianzhu Ma, Jian Peng, et al. InstafLOW: One step is enough for
611 high-quality diffusion-based text-to-image generation. In *The Twelfth International Conference*
612 *on Learning Representations*, 2023a.
- 613 Xingchao Liu, Xiwen Zhang, Jianzhu Ma, Jian Peng, et al. InstafLOW: One step is enough for
614 high-quality diffusion-based text-to-image generation. In *The Twelfth International Conference*
615 *on Learning Representations*, 2023b.
616
- 617 Xuewen Liu, Zhikai Li, and Qingyi Gu. Cachequant: Comprehensively accelerated diffusion mod-
618 els. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,
619 2025b.
- 620 Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver: A fast
621 ode solver for diffusion probabilistic model sampling in around 10 steps. *Advances in Neural*
622 *Information Processing Systems*, 35:5775–5787, 2022.
623
- 624 Jianrong Lu, Lulu Xue, Wei Wan, Minghui Li, Leo Yu Zhang, and Shengqing Hu.
625 Preserving privacy of input features across all stages of collaborative learning. In
626 *ISPA/BDCLOUD/SocialCom/SustainCom*, 2023.
- 627 Jianrong Lu, Shengshan Hu, Wei Wan, Minghui Li, Leo Yu Zhang, Lulu Xue, Haohan Wang, and
628 Hai Jin. Depriving the survival space of adversaries against poisoned gradients in federated learn-
629 ing. *IEEE Transactions on Information Forensics and Security*, 2024.
- 630 Jianrong Lu, Zhiyu Zhu, and Junhui Hou. Parasolver: A hierarchical parallel integral solver for
631 diffusion models. In *The Thirteenth International Conference on Learning Representations*, 2025.
632
- 633 Simian Luo, Yiqin Tan, Longbo Huang, Jian Li, and Hang Zhao. Latent consistency models: Synthe-
634 sizing high-resolution images with few-step inference. *arXiv preprint arXiv:2310.04378*, 2023.
635
- 636 Xinyin Ma, Gongfan Fang, Michael Bi Mi, and Xinchao Wang. Learning-to-cache: Accelerating
637 diffusion transformer via layer caching. *Advances in Neural Information Processing Systems*, 37:
638 133282–133304, 2024a.
- 639 Xinyin Ma, Gongfan Fang, and Xinchao Wang. Deepcache: Accelerating diffusion models for free.
640 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages
641 15762–15772, 2024b.
- 642 Loc X Nguyen, Pyae Sone Aung, Huy Q Le, Seong-Bae Park, and Choong Seon Hong. A new chap-
643 ter for medical image generation: the stable diffusion method. In *2023 International Conference*
644 *on Information Networking (ICOIN)*, pages 483–486. IEEE, 2023.
645
- 646 Zizheng Pan, Bohan Zhuang, De-An Huang, Weili Nie, Zhiding Yu, Chaowei Xiao, Jianfei Cai,
647 and Anima Anandkumar. T-stitch: Accelerating sampling in pre-trained diffusion models with
trajectory stitching. *arXiv preprint arXiv:2402.14167*, 2024.

- 648 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor
649 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-
650 performance deep learning library. *Advances in neural information processing systems*, 32, 2019.
651
- 652 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
653 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer-
654 ence on computer vision and pattern recognition*, pages 10684–10695, 2022.
- 655 Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng
656 Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei.
657 ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision
658 (IJCV)*, 115(3):211–252, 2015. doi: 10.1007/s11263-015-0816-y.
659
- 660 Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. In
661 *International Conference on Learning Representations*.
- 662 Andy Shih, Suneel Belkhale, Stefano Ermon, Dorsa Sadigh, and Nima Anari. Parallel sampling of
663 diffusion models. *Advances in Neural Information Processing Systems*, 36, 2024a.
664
- 665 Andy Shih, Suneel Belkhale, Stefano Ermon, Dorsa Sadigh, and Nima Anari. Parallel sampling of
666 diffusion models. *Advances in Neural Information Processing Systems*, 36, 2024b.
- 667 Junhyuk So, Jiwoong Shin, Chaeyeon Jang, and Eunhyeok Park. Pcm: Picard consistency model
668 for fast parallel sampling of diffusion models. In *Proceedings of the IEEE/CVF Conference on
669 Computer Vision and Pattern Recognition*, 2025.
670
- 671 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *Internat-
672 ional Conference on Learning Representations*, a.
- 673 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben
674 Poole. Score-based generative modeling through stochastic differential equations. In *Internat-
675 ional Conference on Learning Representations*, b.
676
- 677 Yang Song, Chenlin Meng, Renjie Liao, and Stefano Ermon. Accelerating feedforward computation
678 via parallel nonlinear equation solving. In *International Conference on Machine Learning*, pages
679 9791–9800. PMLR, 2021.
- 680 Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency models. *International
681 Conference on Machine Learning*, 2023.
682
- 683 Zhiwei Tang, Jiasheng Tang, Hao Luo, Fan Wang, and Tsung-Hui Chang. Accelerating parallel
684 sampling of diffusion models. In *Forty-first International Conference on Machine Learning*, 2024.
685
- 686 Patrick von Platen, Suraj Patil, Anton Lozhkov, Pedro Cuenca, Nathan Lambert, Kashif Ras-
687 sul, Mishig Davaadorj, Dhruv Nair, Sayak Paul, William Berman, Yiyi Xu, Steven Liu, and
688 Thomas Wolf. Diffusers: State-of-the-art diffusion models. [https://github.com/
689 huggingface/diffusers](https://github.com/huggingface/diffusers), 2022.
- 690 Guangyi Wang, Yuren Cai, Lijiang Li, Wei Peng, and Songzhi Su. Pfdiff: Training-free acceler-
691 ation of diffusion models through the gradient guidance of past and future. In *The Thirteenth
692 International Conference on Learning Representations*, 2025.
693
- 694 Jiannan Wang, Jiarui Fang, Aoyu Li, and PengCheng Yang. Pipefusion: Displaced patch pipeline
695 parallelism for inference of diffusion transformer models. *arXiv preprint arXiv:2405.14430*,
696 2024.
- 697 Junde Wu, Rao Fu, Huihui Fang, Yu Zhang, Yehui Yang, Haoyi Xiong, Huiying Liu, and Yanwu
698 Xu. Medsegdiff: Medical image segmentation with diffusion probabilistic model. In *Medical
699 Imaging with Deep Learning*, pages 1623–1639. PMLR, 2024.
- 700 Run Yang, Yuling Yang, Fan Zhou, and Qiang Sun. Directional diffusion models for graph repre-
701 sentation learning. *Advances in Neural Information Processing Systems*, 36, 2024.

702 Meng You, Zhiyu Zhu, Hui Liu, and Junhui Hou. Nvs-solver: Video diffusion model as zero-shot
703 novel view synthesizer. *International Conference on Learning Representations*, 2025.
704

705 Fisher Yu, Ari Seff, Yinda Zhang, Shuran Song, Thomas Funkhouser, and Jianxiong Xiao. Lsun:
706 Construction of a large-scale image dataset using deep learning with humans in the loop. *arXiv
707 preprint arXiv:1506.03365*, 2015.

708 Jean Yu and Haim Barad. Step saver: Predicting minimum denoising steps for diffusion model
709 image generation. *arXiv preprint arXiv:2408.02054*, 2024.
710

711 Hongkai Zheng, Weili Nie, Arash Vahdat, Kamyar Azizzadenesheli, and Anima Anandkumar. Fast
712 sampling of diffusion models via operator learning. In *International conference on machine learn-
713 ing*, pages 42390–42402. PMLR, 2023.

714 Huanjian Zhou and Masashi Sugiyama. Parallel simulation for sampling under isoperimetry and
715 score-based diffusion models. *arXiv preprint arXiv:2412.07435*, 2024.
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

| | | |
|-----|---|-----------|
| 756 | CONTENTS | |
| 757 | | |
| 758 | 1 Introduction | 1 |
| 759 | | |
| 760 | 2 Related Work | 2 |
| 761 | | |
| 762 | 3 Preliminary | 3 |
| 763 | | |
| 764 | 4 Proposed Method | 4 |
| 765 | | |
| 766 | 4.1 ρ -Nonlinear Equation System | 4 |
| 767 | | |
| 768 | 5 Solving the ρ-Nonlinear Equation System | 6 |
| 769 | | |
| 770 | 6 Experiments | 8 |
| 771 | | |
| 772 | 7 Conclusion | 10 |
| 773 | | |
| 774 | A Use of LLM | 15 |
| 775 | | |
| 776 | B Limitations, Future Directions, and Broader Impacts | 15 |
| 777 | | |
| 778 | C The Framework of ParaSolver-Turbo | 16 |
| 779 | | |
| 780 | D Experimental Details | 16 |
| 781 | | |
| 782 | E Acceleration for LSUN Church model | 17 |
| 783 | | |
| 784 | F Vision Comparision | 17 |
| 785 | | |
| 786 | G Vision Results | 17 |
| 787 | | |
| 788 | H Proof of Proposition 1 | 22 |
| 789 | | |
| 790 | I Proof of Proposition 2 | 23 |
| 791 | | |
| 792 | J Proof of Convergence in Proposition 3 | 25 |
| 793 | | |
| 794 | A USE OF LLM | |
| 795 | | |
| 796 | | |
| 797 | During the preparation of this work, we used Large Language Models (LLMs) to assist with the | |
| 798 | writing process. The primary uses included polishing and improving the fluency of the text, gener- | |
| 799 | ating preliminary drafts of proofs, and assisting in the creation and formatting of tables. After using | |
| 800 | these tools, the author(s) reviewed and edited the content extensively. We take full responsibility | |
| 801 | for the entire content of this publication, including the ideas, proofs, and presentations ultimately | |
| 802 | contained in the final manuscript. | |
| 803 | | |
| 804 | B LIMITATIONS, FUTURE DIRECTIONS, AND BROADER IMPACTS | |
| 805 | | |
| 806 | Despite the promising results, our work has several limitations that open up avenues for future | |
| 807 | research. | |
| 808 | | |
| 809 | • Higher NFE than Sequential Methods: Although our linearization method significantly | |
| | reduces the Number of Function Evaluations (NFE) required for parallel sampling, the NFE | |

is still higher than that of sequential methods. Further reduction in NFE would substantially enhance parallel efficiency.

- **Computational Redundancy:** Compared to sequential sampling, our method requires more computational power. Considerable redundancy exists in the current parallel computation. Designing mechanisms, such as a cache, to minimize this redundant computation could further improve parallel efficiency.
- **Iterative Nature and Single-Step Generation:** Parallel sampling still remains an iterative correction process. A fascinating direction would be to investigate the distillation of parallel iterations or the development of a consistency model. This could potentially lead to few-step or even single-step parallel generation, representing a significant leap forward.

The broader impacts. The acceleration of diffusion models, while beneficial for creative applications and scientific research (e.g., molecular generation), also carries the dual-use risks inherent in all powerful generative technologies. We acknowledge that faster model inference could potentially be misused for creating disinformation or other harmful content. We believe that by advancing the technical understanding of these models, our work can also contribute to the development of more robust detection and mitigation strategies.

C THE FRAMEWORK OF PARASOLVER-TURBO

Figure 1 shows the process of our ParaSolver-Turbo. The framework of ParaSolver-Turbo involves an iterative process for parallel noise prediction and state update. It starts by defining a window of size p with N time subintervals. Within this window, the algorithm identifies linear and nonlinear regions.

The process begins with an initial window of size p and iteratively updates the state based on the latest samples. For each iteration k , the algorithm calculates residuals in parallel. These residuals are then used to update the state of the system.

The framework also includes an error check to determine convergence. If the error is within an acceptable range, the state is updated, and the window shifts to the next set of samples. If the error is not converged, the process continues with the next iteration. The window size is dynamically adjusted based on the remaining subintervals, ensuring efficient computation and accurate noise prediction.

D EXPERIMENTAL DETAILS

We evaluate ParaSolver-Turbo across various high-dimensional image generation models, such as the latent-space diffusion model StableDiffusion-v2 (Rombach et al., 2022) and the pixel-space diffusion model LSUN Church (Yu et al., 2015). The results of these experiments demonstrate that ParaSolver-Turbo enhances the efficiency of sequential sampling methods all while maintaining consistent sample quality as measured by metrics like FID score or CLIP score.

Evaluation Metrics. We assessed ParaSolver-Turbo using five widely adopted metrics: function evaluation count (NFE), iterations (iters), wall-clock time, FID score (Heusel et al., 2017), and CLIP score (Hessel et al., 2021).

Datasets and Models. As per ParaSolver (Lu et al., 2025), we analyze our ParaSolver-Turbo across latent-space model (StableDiffusion-v2) and pixel-space model (LSUN Church) using the same settings. For latent-space models, we leverage the StableDiffusion-v2 model.² (Russakovsky et al., 2015). The StableDiffusion-v2 model generates images at a resolution of 768 by 768 pixels, utilizing a diffusion model that operates within a latent space dimension of 4 by 96 by 96. For pixel-space models, we employ models pretrained on the LSUN Church dataset³ available on Huggingface (Ho et al., 2020; von Platen et al., 2022). This LSUN Church pre-trained model operates in the pixel space with a resolution of 256 by 256. For all the computation of the FID score and CLIP score, we use the standard evaluation API⁴ provided by Huggingface. For evaluation of the wall-clock

²<https://huggingface.co/datasets/ILSVRC/imagenet-1k>

³<https://huggingface.co/google/ddpm-ema-church-256>

⁴<https://huggingface.co/docs/diffusers/conceptual/evaluation>

time, we use the `torch.cuda.Event`⁵ method provided by Pytorch 2.0 [Paszke et al. \(2019\)](#). We assess the performance of all methods on 8 NVIDIA RTX 3090 GPUs, each equipped with 24268 MB of memory.

Algorithms. We apply our approach to accelerate the performance of the state-of-the-art sequential sampling methods: DDPM ([Ho et al., 2020](#)), DDIM ([Song et al., a](#)), and DPMSolver ([Lu et al., 2022](#)). We evaluate the parallel efficiency of our ParaSolver-Turbo against ParaSolver and ParaDiGMS in accelerating the aforementioned three sequential methods as it incorporates all existing parallel methods. ParaDiGMS⁶ is implemented based on the the Diffusers library, and the other algorithms are all accessible in the widely-used library Diffusers ([von Platen et al., 2022](#)); hence we utilize them directly. However, we exclude the comparison with another parallel method ([Tang et al., 2024](#)) as it has not been integrated into the Diffusers library yet.

Hyperparameter Settings. Following the previous settings in ([Lu et al., 2025](#); [Tang et al., 2024](#)), we apply our ParaSolver-Turbo and ParaDiGMS to DDPM with 1000 sequential sampling steps. For DDIM and DPMSolver, we consider two settings: 25 and 50 sequential sampling steps. as they are commonly utilized and capable of producing samples of similar quality to DDPM. Besides, We follow the best tolerances for ParaDiGMS and ParaSolver in their papers. For StableDiffusion-v2, ParaDiGMS for DDPM, DDIM, and DPMSolver need to set the tolerance as 0.5, 0.01, and 0.01 to achieve a similar sample quality as the corresponding sequential method, respectively; ParaSolver for DDPM, DDIM, and DPMSolver need to set the tolerance as 0.55, 0.01, and 0.01 to achieve a similar sample quality. For LSUN Church model, the tolerance of ParaDiGMS should be 0.5 and 0.001 for DDPM and DDIM; the tolerances of ParaSolver are set as 0.55 and 0.005 for DDPM and DDIM, respectively.

For our ParaSolver-Turbo implementation, we optimized the hyperparameters through grid search. On the StableDiffusion-v2 model, for DDPM, we search the best values on grid $\delta \in \{0.5, 0.6\}$; for DDIM and DPMSolver, we search the best values on grid $\delta \in \{0.05, 0.04, 0.03, 0.02, 0.01\}$. On the LSUN Church model, for DDPM, we search the best values on grid $\delta \in \{0.5, 0.6\}$; for DDIM and DPMSolver, we search the best values on grid $\delta \in \{0.005, 0.004, 0.003, 0.002, 0.001\}$.

E ACCELERATION FOR LSUN CHURCH MODEL

Similarly, ParaSolver-Turbo consistently surpasses existing methods, delivering speed improvements of $2.8\times$ to $50.0\times$ on the LSUN Church model across varying step configurations. At 1000 steps, ParaSolver-Turbo achieves a remarkable $50.0\times$ speedup while maintaining identical FID (12.7), demonstrating its efficiency in high-fidelity generation. For 50-step and 25-step settings, it outperforms DDIM by $3.0\times$ and $2.8\times$, respectively. Notably, ParaSolver-Turbo’s parallelism significantly reduces iterations compared to baselines and minimizes NFE, underscoring its computational efficiency. This acceleration is critical for real-time applications, such as interactive design tools, where latency directly impacts user experience.

F VISION COMPARISION

This section shows the visual comparisons of when our ParaSolver-Turbo is applied to speed up DDPM, DDIM, and DPMSolver on Stable Diffusion v2. The results are shown in [Figure 2](#), [Figure 3](#), and [Figure 4](#). We can see that our ParaSolver-Turbo significantly outperforms the competitors, with a faster speed to generate an image.

G VISION RESULTS

This section shows the visual results of when our ParaSolver-Turbo is applied to speed up DDPM, DDIM, and DPMSolver on Stable Diffusion v2. We generated 16 prompts with different styles by DeepSeek-v3 for each method. The results for DDIM plus ParaSolver-Turbo are shown in [Figure 5](#). The results for DDPM plus ParaSolver-Turbo and the results for DPMSolver plus ParaSolver-Turbo are shown in supplementary.

⁵<https://pytorch.org/docs/stable/generated/torch.cuda.Event.html>

⁶<https://github.com/AndyShih12/paradigms>

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

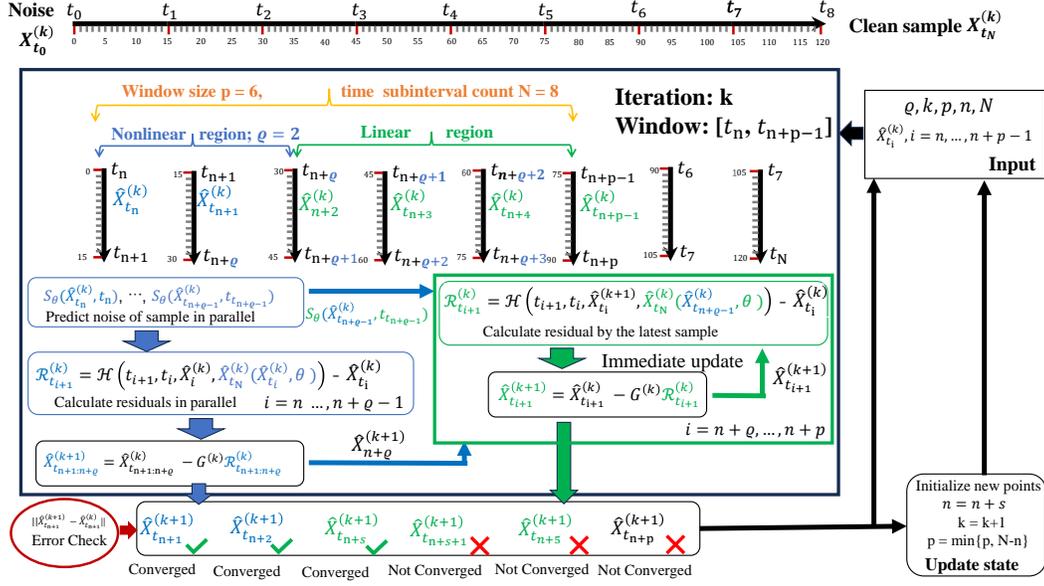


Figure 1: The Framework of ParaSolver-Turbo

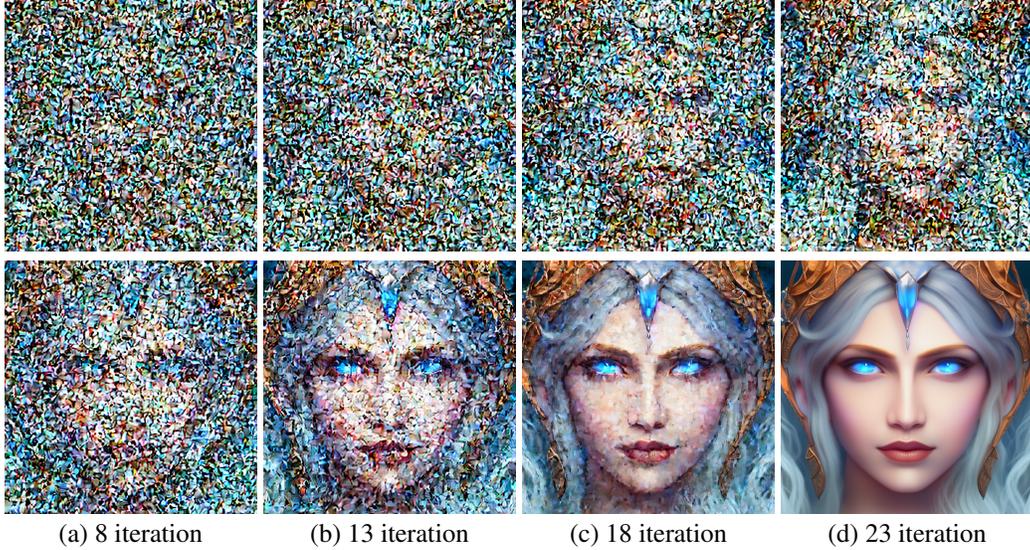
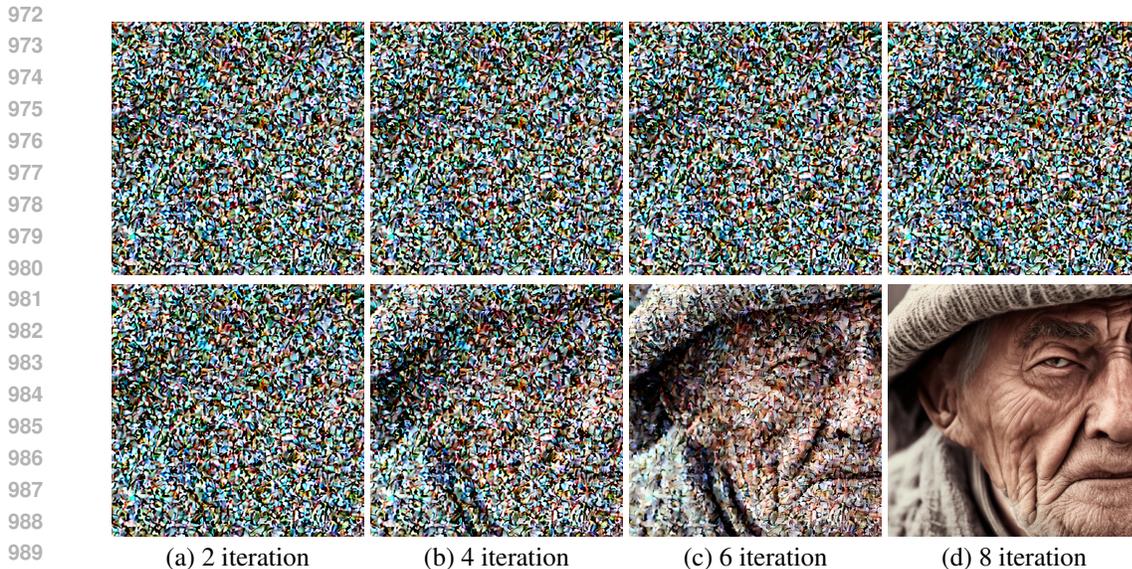
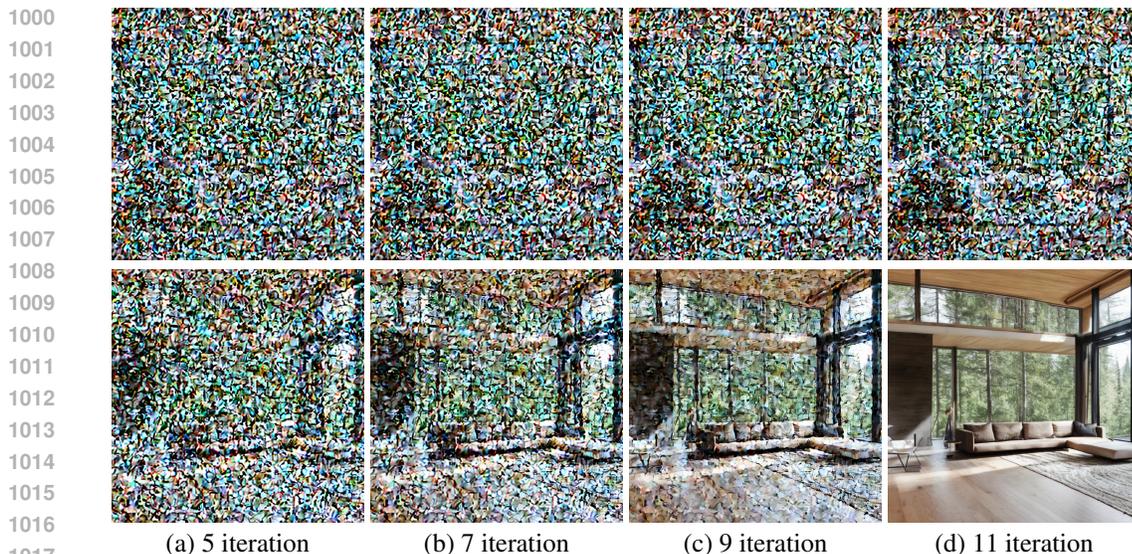


Figure 2: The intermediate generated images for expediting DDPM with 1000 sequential steps on Stable Diffusion Model v2. The images in the first row are produced by DDPM. The images in the second row are generated by our ParaSolver-Turbo. For our ParaSolver-Turbo, we set the number of subintervals N as 1000, the preconditioning steps M as 0, the tolerance δ as 0.5, the maximum parallel steps K as 23, the parallel window size p as 100, the nonlinear degree q as 8. The prompt is "A stunning portrait of an ethereal elf queen with intricate silver jewelry, glowing blue eyes, flowing white hair, soft cinematic lighting, highly detailed, 8K."



992 Figure 3: The intermediate generated images for expediting DDIM with 50 sequential steps on
993 Stable Diffusion Model v2. The images in the first row are produced by DDIM. The images in
994 the second row are generated by ParaSolver-Turbo. For our ParaSolver-Turbo, we set the number
995 of subintervals N as 50, the preconditioning steps M as 0, the tolerance δ as 0.03, the maximum
996 parallel steps K as 11, the parallel window size p as 50, the nonlinear degree ϱ as 8. The prompt is
997 "A highly detailed portrait of an elderly man with wrinkles, wearing a traditional woolen hat,
998 cinematic lighting, 8K, ultra-realistic, photorealistic, depth of field, soft shadows, film grain".
999



1019 Figure 4: The intermediate generated images for expediting DPMSolver with 50 sequential steps on
1020 Stable Diffusion Model v2. The images in the first row are produced by DPMSolver. The images
1021 in the second row are generated by ParaSolver-Turbo. For our ParaSolver-Turbo, we set the number
1022 of subintervals N as 50, the preconditioning steps M as 0, the tolerance δ as 0.03, the maximum
1023 parallel steps K as 11, the parallel window size p as 50, the nonlinear degree ϱ as 8. The prompt is
1024 "A modern minimalist living room with floor-to-ceiling windows overlooking a forest, Scandinavian
1025 design, natural wood and neutral tones, soft daylight".

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079



(a) A bustling farmer's market at sunset with vibrant fruits and vegetables on display, people chatting, and warm golden light
(b) A misty morning in a Japanese bamboo forest with sunlight filtering through the tall green stalks
(c) An elderly craftsman hand-making leather shoes in a small workshop filled with tools and natural light
(d) A crowded subway station during rush hour with commuters wearing winter coats and carrying briefcases



(e) A Mediterranean coastal village with whitewashed houses, blue shutters, and fishing boats in the harbor
(f) A wildlife photographer crouching in tall grass to capture a lioness with her cubs on the African savanna
(g) A rainy city street at night with colorful neon signs reflecting on wet pavement and people with umbrellas
(h) A traditional Italian kitchen with fresh pasta drying on racks, tomatoes, basil, and copper pots hanging



(i) A high-altitude mountain base camp with tents, climbers in heavy gear, and snow-capped peaks in background
(j) A vintage bookstore with floor-to-ceiling wooden shelves, ladder, and sunlight streaming through windows
(k) A busy construction site with workers in hard hats, cranes, and a half-built steel structure against blue sky
(l) A peaceful lakeside at dawn with mist rising off still water and a single wooden rowboat tied up



(m) A street food vendor in Bangkok cooking pad thai on a sizzling wok with smoke and aromatic steam rising
(n) An autumn forest path covered in fallen orange leaves with sunlight creating dappled patterns through trees
(o) A professional chef plating an elegant dish in a high-end restaurant with precise attention to detail
(p) A traditional Moroccan riad courtyard with intricate tile work, fountain, and potted orange trees

1080
 1081
 1082
 1083
 1084
 1085
 1086
 1087
 1088
 1089
 1090
 1091
 1092
 1093
 1094
 1095
 1096
 1097
 1098
 1099
 1100
 1101
 1102
 1103
 1104
 1105
 1106
 1107
 1108
 1109
 1110
 1111
 1112
 1113
 1114
 1115
 1116
 1117
 1118
 1119
 1120
 1121
 1122
 1123
 1124
 1125
 1126
 1127
 1128
 1129
 1130
 1131
 1132
 1133

Table 5: Quantitative comparisons of different methods on LSUN Church over 5000 samples. DPM-Solver is not integrated with the LSUN model in the Diffusers library; thus we exclude it.

| Steps | Method | LSUN Church | | | | |
|-------|-------------------------|-------------|-------|-------|------------|---|
| | | Iters ↓ | NFE ↓ | FID ↓ | Time (s) ↓ | Speedup ↑ |
| 1000 | DDPM | 1000 | 1000 | 12.7 | 50.0 |  1.0× |
| | DDPM + ParaDiGMS | 65 | 2082 | 12.7 | 10.8 |  4.6× |
| | DDPM + ParaSolver | 42 | 1079 | 12.7 | 4.1 |  12.1× |
| | DDPM + ParaSolver-Turbo | 42 | 196 | 12.7 | 1.0 |  50.0 × |
| 50 | DDIM | 50 | 50 | 15.5 | 1.8 |  1.0× |
| | DDIM + ParaDiGMS | 23 | 202 | 15.8 | 1.5 |  1.2× |
| | DDIM + ParaSolver | 14 | 73 | 15.7 | 0.8 |  2.2× |
| | DDIM + ParaSolver-Turbo | 14 | 67 | 15.7 | 0.6 |  3.0 × |
| 25 | DDIM | 25 | 25 | 15.6 | 1.1 |  1.0× |
| | DDIM + ParaDiGMS | 15 | 96 | 15.7 | 0.8 |  1.4× |
| | DDIM + ParaSolver | 10 | 44 | 15.7 | 0.5 |  2.2× |
| | DDIM + ParaSolver-Turbo | 10 | 39 | 15.7 | 0.4 |  2.8 × |

1134 H PROOF OF PROPOSITION 1

1135
1136 Before proceeding with the proof, we first state the following practical assumption regarding the
1137 denoising model.

1138 **Assumption 1.** *The denoising network $\mathbf{S}_{\theta^*}(\mathbf{X}_t, t)$ is a well-trained noise predictor. Specifically,*
1139 *for any noisy sample $\mathbf{X}_t = \alpha_t \mathbf{X}_T + \sigma_t \epsilon_t$ (where \mathbf{X}_T is the clean sample and $\epsilon_t \sim \mathcal{N}(0, \mathbf{I})$ is the*
1140 *specific noise instance), there exists an upper bound $\zeta > 0$ such that the model’s prediction error is*
1141 *bounded:*

$$1142 \|\mathbf{S}_{\theta^*}(\mathbf{X}_t, t) - \epsilon_t\| < \zeta$$

1143
1144 We recall the forward process $p(\mathbf{X}_t | \mathbf{X}_T) = \mathcal{N}(\alpha_t \mathbf{X}_T, \sigma_t^2 \mathbf{I})$, where α_t and σ_t are the noise schedule
1145 coefficients. We define the log-SNR as $\lambda_t = \log(\alpha_t / \sigma_t)$. It follows that $e^{-\lambda_t} = \sigma_t / \alpha_t$ and λ_t is a
1146 strictly decreasing function of t .

1147 Following Proposition 3.1 in DPM-Solver (Lu et al., 2022), the exact solution of the corresponding
1148 SDE from time s to t is given by:

$$1149 \Phi(t, s, \mathbf{X}_s) = \frac{\alpha_t}{\alpha_s} \mathbf{X}_s + \alpha_t \int_{\lambda_s}^{\lambda_t} e^{-\lambda} \sigma_\lambda \nabla_{\mathbf{X}_\lambda} \log p(\mathbf{X}_\lambda) d\lambda$$

$$1150 + \int_{\lambda_s}^{\lambda_t} g(\lambda) d\mathbf{W}, \quad (23)$$

1151
1152 where $p(\mathbf{X}_\lambda)$ is the marginal distribution. In a practical solver, the true score $\nabla_{\mathbf{X}_\lambda} \log p(\mathbf{X}_\lambda) =$
1153 $-\mathbb{E}[\epsilon | \mathbf{X}_\lambda] / \sigma_\lambda$ is unavailable and is replaced by the model’s prediction. Specifically, the solver
1154 approximates the score using the noise prediction model \mathbf{S}_{θ^*} :

$$1155 \hat{\nabla}_{\mathbf{X}_\lambda} \log p(\mathbf{X}_\lambda) = -\frac{\mathbf{S}_{\theta^*}(\mathbf{X}_\lambda, \lambda)}{\sigma_\lambda}$$

1156
1157 By substituting this approximation into Eq. (23), we define the **practical sequential solver** $\hat{\Phi}$ as:

$$1158 \hat{\Phi}(t, s, \mathbf{X}_s) = \frac{\alpha_t}{\alpha_s} \mathbf{X}_s + \alpha_t \int_{\lambda_s}^{\lambda_t} e^{-\lambda} \sigma_\lambda \left(-\frac{\mathbf{S}_{\theta^*}(\mathbf{X}_\lambda, \lambda)}{\sigma_\lambda} \right) d\lambda$$

$$1159 + \int_{\lambda_s}^{\lambda_t} g(\lambda) d\mathbf{W} \quad (24)$$

$$1160 = \frac{\alpha_t}{\alpha_s} \mathbf{X}_s - \alpha_t \int_{\lambda_s}^{\lambda_t} e^{-\lambda} \mathbf{S}_{\theta^*}(\mathbf{X}_\lambda, \lambda) d\lambda + \int_{\lambda_s}^{\lambda_t} g(\lambda) d\mathbf{W}.$$

1161
1162 Our goal is to show that this practical solver $\hat{\Phi}$ can be precisely approximated by a linear solver.
1163 We use Assumption 1 to decompose the model’s prediction \mathbf{S}_{θ^*} into the ground-truth noise ϵ_λ and
1164 a bounded error term $\delta(\mathbf{X}_\lambda, \lambda)$:

$$1165 \mathbf{S}_{\theta^*}(\mathbf{X}_\lambda, \lambda) = \epsilon_\lambda + \delta(\mathbf{X}_\lambda, \lambda),$$

1166 where $\|\delta(\mathbf{X}_\lambda, \lambda)\| < \zeta$.

1167 Substituting this decomposition back into our practical solver Eq. (24):

$$1168 \hat{\Phi}(t, s, \mathbf{X}_s) = \frac{\alpha_t}{\alpha_s} \mathbf{X}_s - \alpha_t \int_{\lambda_s}^{\lambda_t} e^{-\lambda} (\epsilon_\lambda + \delta(\mathbf{X}_\lambda, \lambda)) d\lambda$$

$$1169 + \int_{\lambda_s}^{\lambda_t} g(\lambda) d\mathbf{W}$$

$$1170 = \underbrace{\frac{\alpha_t}{\alpha_s} \mathbf{X}_s - \alpha_t \int_{\lambda_s}^{\lambda_t} e^{-\lambda} \epsilon_\lambda d\lambda + \int_{\lambda_s}^{\lambda_t} g(\lambda) d\mathbf{W}}_{\mathcal{H}(t, s, \mathbf{X}_s, \mathbf{X}_T)} \quad (25)$$

$$1171 - \underbrace{\alpha_t \int_{\lambda_s}^{\lambda_t} e^{-\lambda} \delta(\mathbf{X}_\lambda, \lambda) d\lambda}_{\mathcal{Z}}.$$

We can identify the first part, \mathcal{H} , as the ideal linear solver. By substituting the definition of the ground-truth noise, $\epsilon_\lambda = (\mathbf{X}_\lambda - \alpha_\lambda \mathbf{X}_T(\mathbf{X}_\lambda))/\sigma_\lambda$, this term becomes:

$$\begin{aligned} \mathcal{H}(t, s, \mathbf{X}_s, \mathbf{X}_T) &= \frac{\alpha_t}{\alpha_s} \mathbf{X}_s - \alpha_t \int_{\lambda_s}^{\lambda_t} e^{-\lambda} \frac{\mathbf{X}_\lambda - \alpha_\lambda \mathbf{X}_T(\mathbf{X}_\lambda)}{\sigma_\lambda} d\lambda \\ &\quad + \int_{\lambda_s}^{\lambda_t} g(\lambda) d\mathbf{W}, \end{aligned} \quad (26)$$

which is the linear solver structure presented in the proposition.

The second part, \mathcal{Z} , represents the cumulative error introduced by the non-ideal denoising model:

$$\mathcal{Z} = -\alpha_t \int_{\lambda_s}^{\lambda_t} e^{-\lambda} \delta(\mathbf{X}_\lambda, \lambda) d\lambda$$

We can bound the norm of this error term using the triangle inequality for integrals and Assumption 1:

$$\begin{aligned} \|\mathcal{Z}\| &= \left\| -\alpha_t \int_{\lambda_s}^{\lambda_t} e^{-\lambda} \delta(\mathbf{X}_\lambda, \lambda) d\lambda \right\| \\ &\leq \alpha_t \int_{\lambda_s}^{\lambda_t} e^{-\lambda} \|\delta(\mathbf{X}_\lambda, \lambda)\| d\lambda \\ &\leq \alpha_t \zeta \int_{\lambda_s}^{\lambda_t} e^{-\lambda} d\lambda \\ &= \alpha_t \zeta \left[-e^{-\lambda} \right]_{\lambda_s}^{\lambda_t} \\ &= \alpha_t \zeta \left(-e^{-\lambda_t} - (-e^{-\lambda_s}) \right) \\ &= \alpha_t \zeta \left(e^{-\lambda_s} - e^{-\lambda_t} \right). \end{aligned} \quad (27)$$

By substituting the definition $e^{-\lambda} = \sigma/\alpha$, we get the final bound:

$$\|\mathcal{Z}\| \leq \alpha_t \zeta \left(\frac{\sigma_s}{\alpha_s} - \frac{\sigma_t}{\alpha_t} \right) = \zeta \sigma_t \left(\frac{\alpha_t \sigma_s}{\alpha_s \sigma_t} - 1 \right)$$

Given that s and t represent consecutive time steps (i.e., the step size $s - t$ is small), we have $t \approx s$, $\alpha_t \approx \alpha_s$, and $\sigma_t \approx \sigma_s$. As the step size approaches zero, the bound approaches:

$$\lim_{t \rightarrow s} \int_{t_0}^{t_1} \zeta \sigma_t \left(\frac{\alpha_t \sigma_s}{\alpha_s \sigma_t} - 1 \right) dt = \zeta (t_1 - t_0).$$

Since ζ is assumed to be small (for a well-trained model); $t_1 - t_0$ is also a limited value denoting the integral time range; and the term $\zeta \sigma_t \left(\frac{\alpha_t \sigma_s}{\alpha_s \sigma_t} - 1 \right)$ is small for typical step sizes, the error $\|\mathcal{Z}\|$ is negligible. This demonstrates that the practical sequential solver $\hat{\Phi}$ is well-approximated by the linear solver \mathcal{H} , completing the proof.

I PROOF OF PROPOSITION 2

Proof. Let $\{\mathbf{X}_{t_n}\}_{n=0}^N$ be the exact trajectory generated by the sequential integral solver. From Proposition 1, the step transition is given by:

$$\mathbf{X}_{t_{n+1}} = \mathcal{H}(t_{n+1}, t_n, \mathbf{X}_{t_n}, \mathbf{X}_T) + \mathcal{Z}_n, \quad (28)$$

where \mathbf{X}_T is the ground-truth data and $\|\mathcal{Z}_n\| \approx 0$ is the linearization error. Let $\{\hat{\mathbf{X}}_{t_n}\}_{n=0}^N$ be the solution to the ϱ -nonlinear system with the update rule:

$$\hat{\mathbf{X}}_{t_{n+1}} = \mathcal{H}(t_{n+1}, t_n, \hat{\mathbf{X}}_{t_n}, \mathbf{E}_{t_n}). \quad (29)$$

Here, \mathbf{E}_{t_n} denotes the estimated clean sample. Based on Definition 2, this estimator is explicitly defined as:

$$\mathbf{E}_{t_n} = \begin{cases} \hat{\mathbf{X}}_{t_N}(\hat{\mathbf{X}}_{t_n}, \theta) & \text{if } n < \varrho, \\ \hat{\mathbf{X}}_{t_N}(\hat{\mathbf{X}}_{t_{\varrho-1}}, \theta) & \text{if } n \geq \varrho. \end{cases} \quad (30)$$

We define the cumulative state error as $\Delta_{t_n} = \hat{\mathbf{X}}_{t_n} - \mathbf{X}_{t_n}$ and the clean sample estimation error as $\xi_n = \mathbf{E}_{t_n} - \mathbf{X}_T$. To derive the error recurrence, we subtract Eq. (28) from Eq. (29). Using the integral definition of \mathcal{H} from Eq. (7), we obtain:

$$\Delta_{t_{n+1}} = \frac{\alpha_{t_{n+1}}}{\alpha_{t_n}} \Delta_{t_n} - \alpha_{t_{n+1}} \int_{\lambda_{t_n}}^{\lambda_{t_{n+1}}} e^{-\lambda} \left(\frac{\hat{\mathbf{X}}_{\lambda} - \alpha_{\lambda} \mathbf{E}_{t_n}}{\sigma_{\lambda}} - \frac{\mathbf{X}_{\lambda} - \alpha_{\lambda} \mathbf{X}_T}{\sigma_{\lambda}} \right) d\lambda - \mathcal{Z}_n. \quad (31)$$

Using the Log-SNR property that $e^{-\lambda} = \sigma_{\lambda}/\alpha_{\lambda}$, the integrand simplifies as follows:

$$\frac{\sigma_{\lambda}}{\alpha_{\lambda}} \left(\frac{\hat{\mathbf{X}}_{\lambda} - \mathbf{X}_{\lambda}}{\sigma_{\lambda}} - \frac{\alpha_{\lambda}(\mathbf{E}_{t_n} - \mathbf{X}_T)}{\sigma_{\lambda}} \right) = \frac{\Delta_{\lambda}}{\alpha_{\lambda}} - \xi_n. \quad (32)$$

Substituting this back into the error equation yields:

$$\Delta_{t_{n+1}} = \frac{\alpha_{t_{n+1}}}{\alpha_{t_n}} \Delta_{t_n} - \alpha_{t_{n+1}} \int_{\lambda_{t_n}}^{\lambda_{t_{n+1}}} \left(\frac{\Delta_{\lambda}}{\alpha_{\lambda}} - \xi_n \right) d\lambda - \mathcal{Z}_n. \quad (33)$$

To evaluate the intermediate error Δ_{λ} inside the integral, we invoke the **Superposition Principle** for linear differential equations (Boyce et al., 2021; Chen, 1984). The total error Δ_{λ} can be rigorously decomposed into the sum of the *Zero-Input Response* (Homogeneous) and the *Zero-State Response* (Forced):

- **Homogeneous Response:** This component represents the natural evolution of the initial state error Δ_{t_n} through the system dynamics, independent of the external estimation error. Due to the signal scaling inherent in the diffusion process \mathcal{H} , this is given by:

$$\Delta_{\lambda}^{\text{homo}} = \frac{\alpha_{\lambda}}{\alpha_{t_n}} \Delta_{t_n}. \quad (34)$$

- **Forced Response:** This component captures the error accumulated over the interval $[\lambda_{t_n}, \lambda]$ solely due to the input source ξ_n . Based on the integral structure of the linear solver in Eq. (7), the contribution accumulates as:

$$\Delta_{\lambda}^{\text{force}} = \alpha_{\lambda} \left[\int_{\lambda_{t_n}}^{\lambda} e^{-\tau} \frac{\alpha_{\tau}}{\sigma_{\tau}} d\tau \right] \xi_n. \quad (35)$$

Substituting $e^{-\tau} = \sigma_{\tau}/\alpha_{\tau}$, the integrand simplifies perfectly to unity:

$$\Delta_{\lambda}^{\text{force}} = \alpha_{\lambda} \left[\int_{\lambda_{t_n}}^{\lambda} 1 d\tau \right] \xi_n = \alpha_{\lambda}(\lambda - \lambda_{t_n}) \xi_n. \quad (36)$$

Combining these components via superposition, we obtain the explicit expression for the intermediate error:

$$\Delta_{\lambda} = \underbrace{\frac{\alpha_{\lambda}}{\alpha_{t_n}} \Delta_{t_n}}_{\text{Homogeneous Response}} + \underbrace{\alpha_{\lambda}(\lambda - \lambda_{t_n}) \xi_n}_{\text{Forced Response}}. \quad (37)$$

We substitute this expression for Δ_{λ} into the term $\frac{\Delta_{\lambda}}{\alpha_{\lambda}}$. This leads to an algebraic cancellation of the time-dependent coefficient α_{λ} in the homogeneous term:

$$\frac{\Delta_{\lambda}}{\alpha_{\lambda}} = \frac{1}{\alpha_{\lambda}} \left(\frac{\alpha_{\lambda}}{\alpha_{t_n}} \Delta_{t_n} + \alpha_{\lambda}(\lambda - \lambda_{t_n}) \xi_n \right) = \frac{1}{\alpha_{t_n}} \Delta_{t_n} + (\lambda - \lambda_{t_n}) \xi_n. \quad (38)$$

Now we can analytically evaluate the integral in Eq. (33):

$$\begin{aligned} \int_{\lambda_{t_n}}^{\lambda_{t_{n+1}}} \left(\frac{\Delta_{\lambda}}{\alpha_{\lambda}} - \xi_n \right) d\lambda &= \int_{\lambda_{t_n}}^{\lambda_{t_{n+1}}} \left(\frac{\Delta_{t_n}}{\alpha_{t_n}} + (\lambda - \lambda_{t_n}) \xi_n - \xi_n \right) d\lambda \\ &= \frac{\Delta_{t_n}}{\alpha_{t_n}} \Delta_{t_n} + \left[\frac{(\lambda - \lambda_{t_n})^2}{2} - (\lambda - \lambda_{t_n}) \right]_{\lambda_{t_n}}^{\lambda_{t_{n+1}}} \xi_n \\ &= \frac{\Delta_{t_n}}{\alpha_{t_n}} \Delta_{t_n} + \left(\frac{(\Delta_{t_n})^2}{2} - \Delta_{t_n} \right) \xi_n, \end{aligned} \quad (39)$$

where $\Delta\lambda_n = \lambda_{t_{n+1}} - \lambda_{t_n}$.

Plugging this integral result back into the expression for $\Delta_{t_{n+1}}$ and grouping the terms by error source:

$$\Delta_{t_{n+1}} = \underbrace{\frac{\alpha_{t_{n+1}}}{\alpha_{t_n}} \left(1 - \Delta\lambda_n\right)}_{\mathcal{A}_n} \Delta_{t_n} + \underbrace{\alpha_{t_{n+1}} \Delta\lambda_n \left(1 - \frac{\Delta\lambda_n}{2}\right)}_{\mathcal{B}_n} \xi_n - \mathcal{Z}_n. \quad (40)$$

Analysis of Coefficients: We analyze the magnitude of the coefficients \mathcal{A}_n and \mathcal{B}_n to establish convergence.

First, for the stability coefficient \mathcal{A}_n , we utilize the identity $\Delta\lambda_n = \log \frac{\alpha_{t_{n+1}} \sigma_{t_n}}{\alpha_{t_n} \sigma_{t_{n+1}}}$ (since we define the log-SNR as $\lambda_t = \log(\alpha_t/\sigma_t)$). The term $(1 - \Delta\lambda_n)$ can be expressed as a logarithmic difference:

$$1 - \Delta\lambda_n \equiv \log \frac{\alpha_{t_{n+1}} \sigma_{t_n}}{\alpha_{t_{n+1}} \sigma_{t_n}} - \log \frac{\alpha_{t_{n+1}} \sigma_{t_n}}{\alpha_{t_n} \sigma_{t_{n+1}}} = \log \frac{\alpha_{t_n} \sigma_{t_{n+1}}}{\alpha_{t_{n+1}} \sigma_{t_n}}. \quad (41)$$

Since the time step size is generally small, the ratio inside the logarithm is close to 1, making this term approach 0. Thus, $\mathcal{A}_n \approx 0$, confirming that the historical error is rapidly dampened.

Second, for the estimation error coefficient \mathcal{B}_n , we apply a similar transformation to the term $(1 - \frac{1}{2}\Delta\lambda_n)$:

$$1 - \frac{1}{2}\Delta\lambda_n \equiv \log \frac{\alpha_{t_{n+1}} \sigma_{t_n}}{\alpha_{t_{n+1}} \sigma_{t_n}} - \log \left(\frac{\alpha_{t_{n+1}} \sigma_{t_n}}{\alpha_{t_n} \sigma_{t_{n+1}}} \right)^{1/2} = \log \left(\sqrt{\frac{\alpha_{t_n} \sigma_{t_{n+1}}}{\alpha_{t_{n+1}} \sigma_{t_n}}} \right). \quad (42)$$

In a similar vein, this logarithmic term approaches 0. Consequently, the entire coefficient \mathcal{B}_n is the product of bounded terms and this vanishing log-difference:

$$\mathcal{B}_n = \alpha_{t_{n+1}} \Delta\lambda_n \left(1 - \frac{\Delta\lambda_n}{2}\right) \approx 0. \quad (43)$$

Conclusion: The error recurrence bound is established as:

$$\|\Delta_{t_{n+1}}\| \leq |\mathcal{A}_n| \|\Delta_{t_n}\| + |\mathcal{B}_n| \|\xi_n\| + \|\mathcal{Z}_n\|. \quad (44)$$

Regardless of the value of ϱ (which determines the source of \mathbf{E}_{t_n}), the non-ideal estimator \mathbf{E}_{t_n} introduces an error ξ_n . However, this estimation error is effectively suppressed by the coefficient \mathcal{B}_n , which approaches zero. Given that the initial error $\|\Delta_{t_0}\| = 0$, and the terms $|\mathcal{A}_n|$, $|\mathcal{B}_n|$, and $\|\mathcal{Z}_n\|$ are all negligible, the total trajectory deviation $\Delta_{t_{n+1}}$ remains negligible throughout the sampling process. Thus, all ϱ -nonlinear systems are statistically equivalent. \square

J PROOF OF CONVERGENCE IN PROPOSITION 3

Assume by induction that for iteration k , the first n components $\{\hat{\mathbf{X}}_{t_i}^{(k)}\}_{i=0}^{n-1}$ have converged to the sequential solution $\{\mathbf{X}_{t_i}\}_{i=0}^{n-1}$ from Proposition 1. Then for the next component $\hat{\mathbf{X}}_{t_n}^{(k+1)}$, we have:

$$\begin{aligned} \hat{\mathbf{X}}_{t_n}^{(k+1)} &= \hat{\mathbf{X}}_{t_n}^{(k)} - \mathcal{G}^{(k)} \mathcal{R}_{t_n}^{(k)} \\ &= \hat{\mathbf{X}}_{t_n}^{(k)} - \mathcal{G}^{(k)} \left(\hat{\mathbf{X}}_{t_n}^{(k)} - \mathcal{H} \left(t_n, t_{n-1}, \hat{\mathbf{X}}_{t_{n-1}}^{(k)}, \hat{\mathbf{X}}_{t_N}(\hat{\mathbf{X}}_{t_{n-1}}^{(k)}, \theta) \right) \right) \\ &\stackrel{(a)}{=} \hat{\mathbf{X}}_{t_n}^{(k)} - \mathcal{G}^{(k)} \left(\hat{\mathbf{X}}_{t_n}^{(k)} - \mathcal{H} \left(t_n, t_{n-1}, \mathbf{X}_{t_{n-1}}, \hat{\mathbf{X}}_{t_N}(\mathbf{X}_{t_{n-1}}, \theta) \right) \right) \\ &\stackrel{(b)}{=} \mathcal{H} \left(t_n, t_{n-1}, \mathbf{X}_{t_{n-1}}, \hat{\mathbf{X}}_{t_N}(\mathbf{X}_{t_{n-1}}, \theta) \right) \\ &= \mathbf{X}_{t_n} \end{aligned} \quad (45)$$

where:

- Eq. (a) follows from the induction hypothesis that $\hat{\mathbf{X}}_{t_{n-1}}^{(k)} = \mathbf{X}_{t_{n-1}}$

1350 • Eq. (b) holds since $\mathcal{G}^{(k)} = \mathbf{I}$ or $\mathcal{G}^{(k)} = \mathcal{J}^{-1} = \mathbf{I}$ where $\mathcal{J} =$
 1351 $\frac{\partial}{\partial \hat{\mathbf{X}}_{t_n}^{(k)}} \left(\hat{\mathbf{X}}_{t_n}^{(k)} - \mathcal{H} \left(t_n, t_{n-1}, \mathbf{X}_{t_{n-1}}, \hat{\mathbf{X}}_{t_N}(\mathbf{X}_{t_{n-1}}, \theta) \right) \right) = \mathbf{I}$
 1352

1353 In the worst case, convergence is achieved in $K = N$ iterations as the convergence front must
 1354 advance exactly one step per iteration. In practice, due to the immediate propagation of updated
 1355 components, convergence is typically achieved in $K \ll N$ iterations.
 1356

1357
 1358
 1359
 1360
 1361
 1362
 1363
 1364
 1365
 1366
 1367
 1368
 1369
 1370
 1371
 1372
 1373
 1374
 1375
 1376
 1377
 1378
 1379
 1380
 1381
 1382
 1383
 1384
 1385
 1386
 1387
 1388
 1389
 1390
 1391
 1392
 1393
 1394
 1395
 1396
 1397
 1398
 1399
 1400
 1401
 1402
 1403