# Qeios

Peer Review

# Review of: "Efficiency Meets Fidelity: A Novel Quantization Framework for Stable Diffusion"

Bhavani Malisetty[1]

1. Computer Science, University of Nebraska, Omaha, United States

## Summary

The paper proposes a novel quantization framework aimed at improving the deployability of Stable Diffusion models by achieving a fine balance between computational efficiency and generation fidelity. Stable Diffusion models, although widely used for high-quality text-to-image generation, are plagued by high computational costs and inference latencies due to their iterative denoising nature. This challenge becomes particularly acute in real-world applications where inference must be fast, scalable, and hardware-efficient—such as in mobile devices, embedded systems, or interactive creative tools.

Quantization offers a promising route to compress these large models by reducing the bit-widths used in computations. However, most prior quantization methods suffer from severe performance degradation at lower bit-widths (e.g., 4-bit) and often ignore the consistency between outputs of quantized and full-precision models, which is crucial for tasks that require stable and predictable image outputs.

To address these limitations, the authors introduce a Serial-to-Parallel (S2P) pipeline—a hybrid quantization training scheme that retains the data efficiency of serial training while leveraging the gradient stability of parallel training. In addition to this pipeline, the framework introduces:

Multi-time step activation quantization to account for varying activation distributions across the denoising schedule.

Time information precalculation to eliminate runtime computation for time embeddings.

A Hessian-based sensitivity-aware mixed-precision quantization strategy to allocate bit-widths more intelligently based on layer importance.

Empirical evaluations across three variants of Stable Diffusion—v1-4, v2-1, and XL 1.0—on the COCO and Stable-Diffusion-Prompts datasets show that the proposed approach significantly outperforms existing quantization methods like PCR, PTQ4DM, and Q-diffusion in both visual quality and training efficiency, even under aggressive low-bit configurations (W4A8).

## Strengths

### 1. Clear Motivation and Practical Relevance

The authors address a critical real-world limitation in deploying large generative models: the high memory, computation, and latency costs of inference. They go beyond simply reducing model size by emphasizing the need for output consistency—a key requirement in creative industries, UI/UX tools, and production environments where repeated prompts should produce predictable and stylistically coherent results. This consideration sets the work apart from traditional model compression techniques.

### 2. Novel Pipeline Design: Serial-to-Parallel (S2P)

The Serial-to-Parallel pipeline is an elegant synthesis of two existing training paradigms:

Serial: Good for mimicking full-precision model behavior but lacks gradient stability.

Parallel: Gradient stable due to multi-timestep training but introduces discrepancies from floating-point outputs.

By using latents generated by the floating-point model and training in a parallel fashion, the authors achieve the best of both worlds: *data-free training, gradient stability, and output fidelity*. This hybridization improves training dynamics and aligns better with the distribution of latents during real inference.

### 3. Multi-Time step Activation Quantization

This contribution acknowledges the reality that activations in diffusion models vary significantly across time steps due to their iterative nature. Instead of applying a uniform quantizer, the authors introduce individual quantization parameter sets per timestep, improving precision and reducing quantization noise. This granularity leads to visually consistent outputs even under low-bit quantization.

### 4. Time Information Precalculation

Time embeddings are a critical part of Stable Diffusion's denoising process, but quantizing these embeddings leads to quality loss. The authors sidestep this issue by precomputing and caching the time embeddings, thereby eliminating the need for quantizing these layers. This not only improves fidelity but also reduces inference complexity and runtime memory.

### 5. Mixed-Precision Quantization via Sensitivity Ranking

Rather than uniformly applying quantization across all layers, the paper proposes a Hessian-based sensitivity metric (approximated using Fisher Information) to rank layers by their importance to model output. The top 5% most sensitive layers are quantized at higher precision, while the least sensitive layers use more aggressive compression. This ensures that critical components retain quality, and compression is maximized elsewhere, resulting in smarter, adaptive quantization.

### 6. Strong Experimental Design and Evaluation

The framework is rigorously tested across:

Three Stable Diffusion versions (v1-4, v2-1, XL 1.0)

Two datasets (COCO, Stable-Diffusion-Prompts)

Multiple metrics (FID-to-FP, sFID, SSIM, LPIPS, PSNR, CLIP score)

The authors also conduct:

Ablation studies isolating the impact of each component

Sampling efficiency experiments under reduced sampling steps (25, 10, and 30)

Results show that their method achieves significantly lower FID-to-FP scores (up to 45% improvement) and better visual similarity, confirming its robustness and generalizability.

### Weaknesses

### 1. Limited Theoretical Grounding for Sensitivity Approximation

While the paper uses Fisher Information as a proxy for Hessian-based sensitivity, a deeper theoretical exploration—perhaps comparing it with other sensitivity estimation techniques—could add rigor and

broaden applicability. The approximation quality and its impact on different architectures or datasets remain underexplored.

## 2. Lack of Deployment Benchmarks on Real Hardware

Despite repeated mention of edge deployment motivation, the paper lacks experimental evidence on:

Actual latency improvements on edge devices (e.g., smartphones, Jetson, Raspberry Pi)

Memory consumption before and after quantization

Power or thermal efficiency

Real hardware validation would significantly boost the practical credibility of the proposed method.

## 3. Limited Reproducibility Information

The authors provide many experimental details (e.g., bit-width configurations, batch sizes, sampling steps) but do not publish code or precomputed latent datasets. Given the complexity of replicating the Serial-to-Parallel setup and sensitivity scoring, public access to these components would substantially enhance reproducibility and community adoption.

## Contributions to the Field

This paper is an important step forward in efficient generative modeling. It tackles an underexplored problem—maintaining fidelity and consistency in quantized diffusion models—with a thoughtfully designed solution that is empirically effective and practically motivated. Its contributions span:

Methodological innovation (Serial-to-Parallel training)

Algorithmic efficiency (precalculation and sensitivity-based mixed precision)

Practical deployment viability (data-free training, reduced memory, and compute)

It sets a new standard for low-bit quantization in diffusion models, paving the way for future research in portable, fast, and reliable generative AI systems.

## Personal Reflection

From a research perspective, this paper aligns closely with the broader goal of making large-scale generative models accessible and usable in constrained environments. The authors deliver not only technical novelty but also a pragmatic framework that acknowledges real-world deployment hurdles.

What makes this work particularly valuable is the attention to fidelity—an often neglected dimension in quantization research. The fact that the method can retain visual, structural, and stylistic similarity to the full-precision models underlines its practical relevance in domains like design, storytelling, and user-facing generative tools.

While the absence of on-device tests and a public implementation limits immediate application, the theoretical and empirical advancements make this a must-read paper for anyone working on efficient deep learning, LLM compression, or AI deployment on the edge.

## Declarations

**Potential competing interests:** No potential competing interests to declare.