# HarmoniCa: Harmonizing Training and Inference for Better Feature Caching in Diffusion Transformer Acceleration

Yushi Huang[1,2]*, Zining Wang[2,3]*, Ruihao Gong[2,3]+, Jing Liu[4], Xinjie Zhang[1], Jinyang Guo[3], Xianglong Liu[3], Jun Zhang[1]+

[1]Hong Kong University of Science and Technology  [2]SenseTime  [3]Beihang University  [4]Monash University

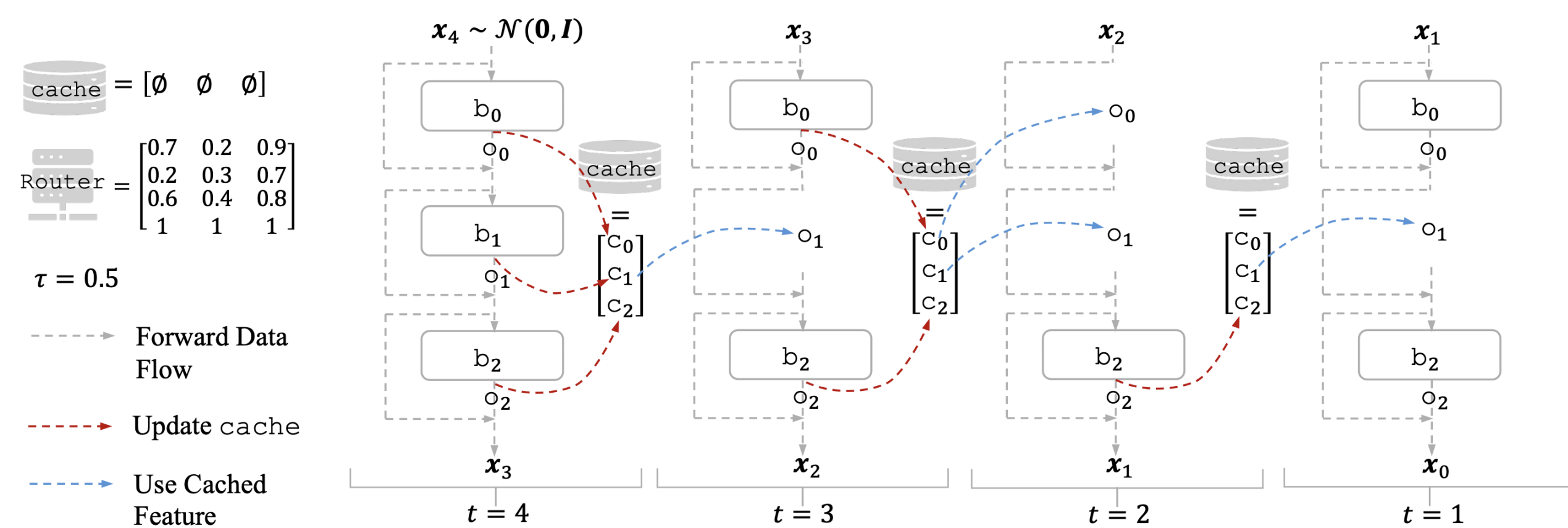*: Equal contribution
+: Corresponding authors

## Feature Caching

High similarity between features across denoising steps enables caching and reusing these features to avoid redundant computations.



**Q:** How to determine **when** and **where** to cache and reuse?
**A:** **Learn** an optimal `Router` (superior than heuristic ways).

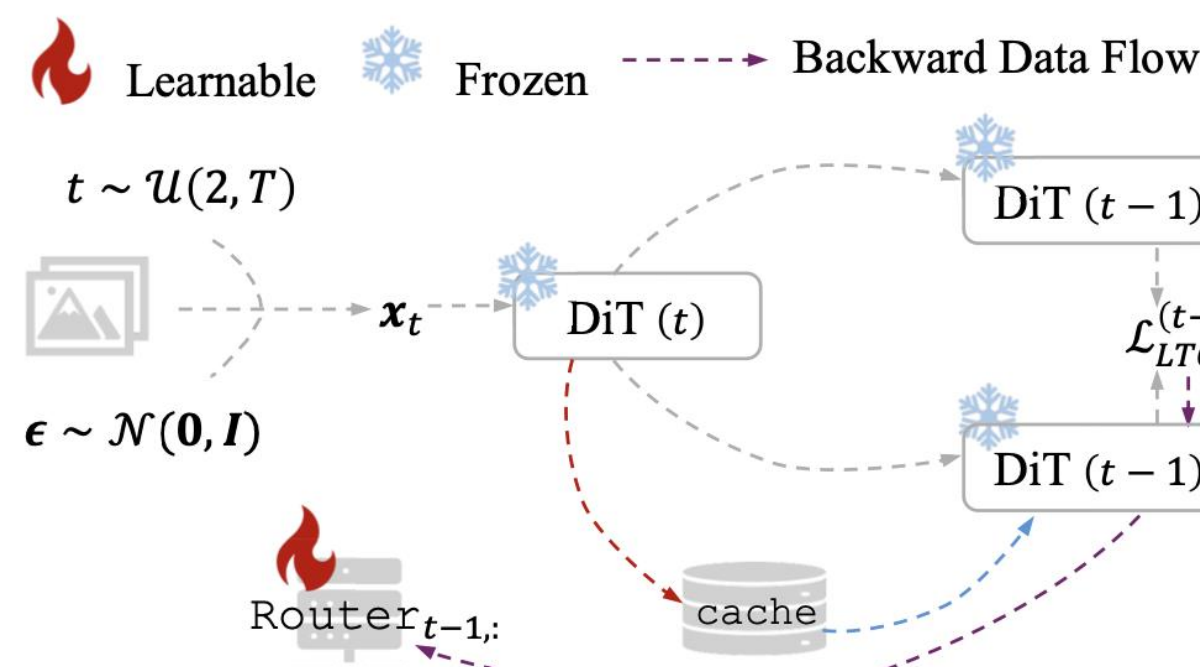## Discrepancy between Training & Inference

### 1. Prior timestep disregard

*Infer:* $x_1$ is pertubed by $o_0$ and $o_1$; Contents of the current `cache` are shaped by reuse and update.
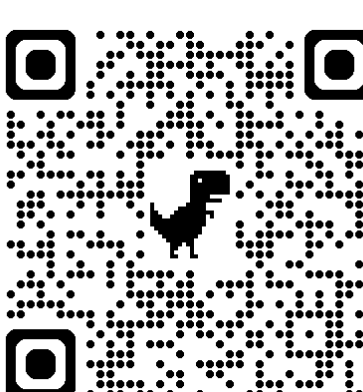*Train:* $x_{t-1}$ is accurate; contents of the current `cache` are only determined by DiT ($t$).

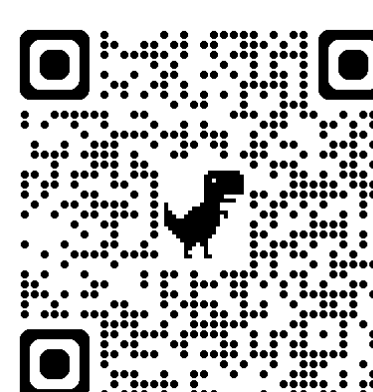### 2. Objective mismatch

*Infer:* Generate high quality $x_0$.
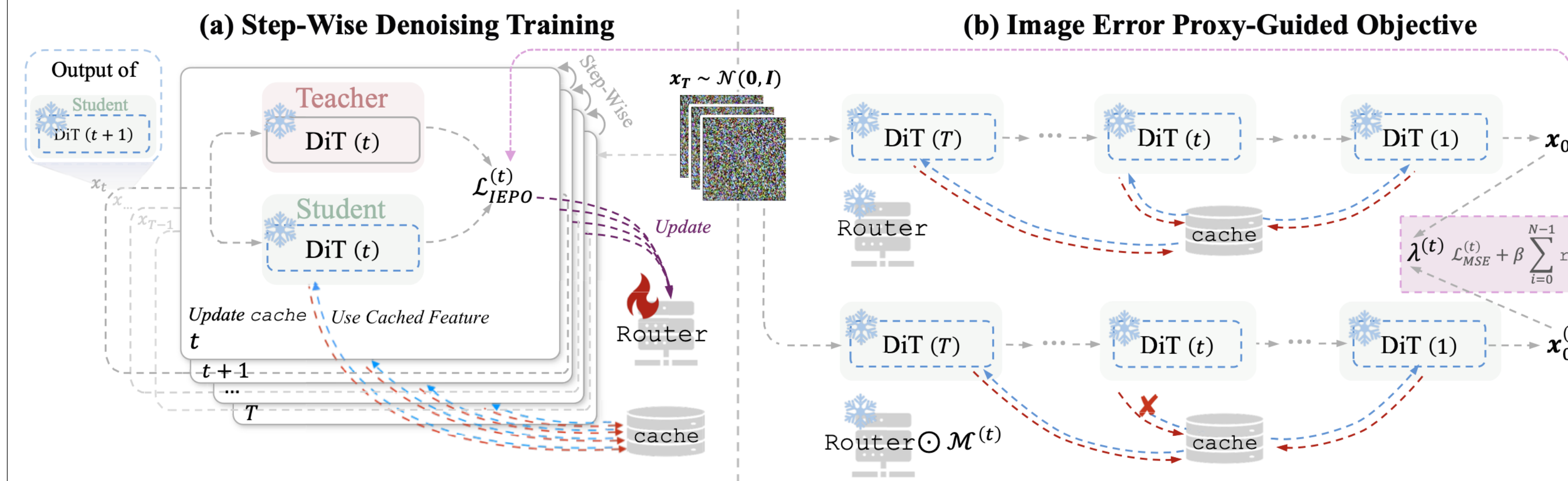*Train:* Solely align the predicted noise at each denoising step.



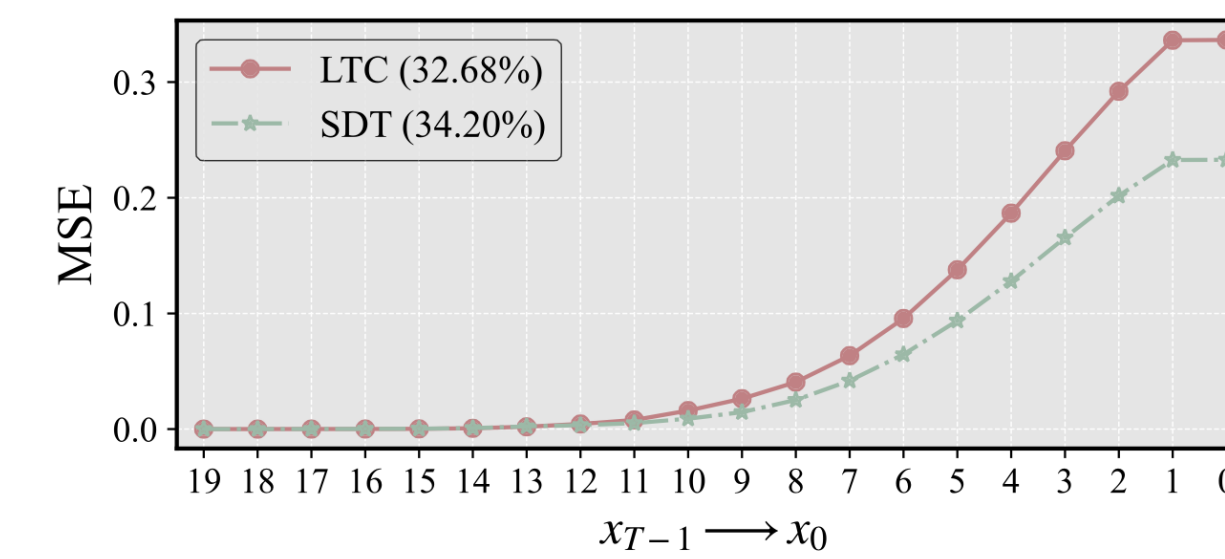$$\mathcal{L}_{LTC}^{(t)} = \mathcal{L}_{MSE}^{(t)} + \beta \sum_{i=0}^{N-1} \mathbf{r}_{t,i}$$

**Paper**  **Code**

## Harmonizing Training & Inference

**(a) Step-Wise Denoising Training**

**(b) Image Error Proxy-Guided Objective**



### 1. Step-wise denoising training

This mimics the multi-timestep inference stage, which integrates the impact of prior timesteps at $t$. Moreover, this strategy requires **zero training images**.
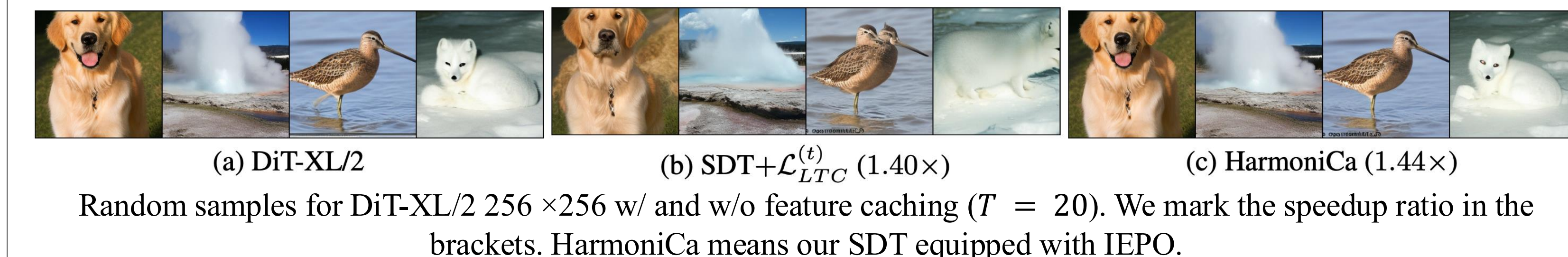
### 2. Image error proxy-guided objective

The straightforward solution (*i.e.,* direct optimize $x_0$) incurs $10\times$ memory and $5\times$ time consumption for training. Thus, we propose a proxy $\lambda^{(t)}$ to consider the trade-off between the error of $x_0$ and the `cache` usage at a certain denoising step.

$$\mathcal{L}_{IEPO}^{(t)} = \lambda^{(t)} \mathcal{L}_{MSE}^{(t)} + \beta \sum_{i=0}^{N-1} \mathbf{r}_{t,i}$$

$$\mathcal{M}^{(t)} = \begin{cases} 1, j \neq t \\ \frac{1}{\mathbf{r}_{j,k}}, j = t \end{cases} \Rightarrow \lambda^{(t)} = \left\| x_0 - x_0^{(t)} \right\|_F^2$$

The Mean Square Error (MSE) of $x_t$ for DiT-XL/2 256 ×256 induced by different feature caching methods.



Random samples for DiT-XL/2 256 ×256 w/ and w/o feature caching ($T = 20$). We mark the speedup ratio in the brackets. HarmoniCa means our SDT equipped with IEPO.

(a) DiT-XL/2     (b) SDT+$\mathcal{L}_{LTC}^{(t)}$ (1.40×)     (c) HarmoniCa (1.44×)

## Experiments

Text-to-image generation.

| Method | T | CLIP↑ | FID↓ | sFID↓ | CUR(%)↑ | Latency(s)↓ |
|---|---|---|---|---|---|---|
| PIXART-α 256 × 256 (cfg = 4.5) | | | | | | |
| DPM-Solver++ (Lu et al., 2022b) | 20 | 30.96 | 27.68 | 36.39 | - | 0.553 |
| DPM-Solver++ (Lu et al., 2022b) | 15 | 30.77 | 31.68 | 38.92 | - | 0.418(1.32×) |
| FORA (Selvaraju et al., 2024) | 20 | 31.10 | 27.42 | 37.98 | 50.00 | 0.364(1.52×) |
| HarmoniCa | 20 | **31.13** | **26.33** | **37.85** | **56.01** | 0.346(1.60×) |
| IDDPM (Nichol & Dhariwal, 2021) | 100 | 31.25 | 24.15 | 33.65 | - | 2.572 |
| IDDPM (Nichol & Dhariwal, 2021) | 75 | 31.25 | 24.17 | 33.73 | - | 1.868(1.37×) |
| FORA (Selvaraju et al., 2024) | 100 | 31.25 | 23.16 | 33.62 | 50.00 | 1.558(1.65×) |
| HarmoniCa | 100 | **31.17** | **23.73** | **32.23** | **53.24** | 1.523(1.69×) |
| SA-Solver (Xue et al., 2024) | 25 | 31.31 | 26.78 | 38.35 | - | 0.891 |
| SA-Solver (Xue et al., 2024) | 20 | 31.23 | 27.45 | 39.01 | - | 0.665(1.34×) |
| HarmoniCa | 25 | **31.27** | **27.07** | **38.62** | **54.19** | 0.561(1.59×) |
| PIXART-α 512 × 512 (cfg = 4.5) | | | | | | |
| DPM-Solver++ (Lu et al., 2022b) | 20 | 31.30 | 23.96 | 40.34 | - | 1.759 |
| DPM-Solver++ (Lu et al., 2022b) | 15 | 31.29 | 25.12 | 40.37 | - | 1.291(1.36×) |
| HarmoniCa | 20 | **31.29** | **24.81** | **40.18** | **54.64** | 1.072(1.64×) |
| SA-Solver (Xue et al., 2024) | 25 | 31.23 | 25.43 | 39.84 | - | 2.263 |
| SA-Solver (Xue et al., 2024) | 20 | 31.19 | 25.85 | 40.08 | - | 1.738(1.30×) |
| HarmoniCa | 25 | **31.20** | **25.74** | **39.99** | **54.24** | 1.406(1.61×) |
| PIXART-α 1024 × 1024 (cfg = 4.5) | | | | | | |
| DPM-Solver++ (Lu et al., 2022b) | 20 | 31.10 | 25.01 | 37.80 | - | 9.470 |
| DPM-Solver++ (Lu et al., 2022b) | 15 | 31.07 | 25.77 | 42.50 | - | 7.141(1.32×) |
| HarmoniCa | 20 | **31.09** | **23.02** | **36.24** | **55.06** | 5.786(1.63×) |
| SA-Solver (Xue et al., 2024) | 25 | 31.03 | 23.65 | 38.12 | - | 11.931 |
| SA-Solver (Xue et al., 2024) | 20 | 31.02 | 23.88 | 39.41 | - | 9.209(1.30×) |
| Harmonica | 25 | **31.07** | **23.77** | **38.93** | **53.98** | 7.551(1.58×) |

Class-conditional generation.

| Method | T | IS↑ | FID↓ | sFID↓ | Prec.↑ | Recall↑ | CUR(%)↑ | Latency(s)↓ |
|---|---|---|---|---|---|---|---|---|
| DiT-XL/2 256 × 256 (cfg = 1.5) | | | | | | | | |
| DDIM (Song et al., 2020a) | 50 | 240.37 | 2.27 | 4.25 | 80.25 | 59.77 | - | 1.767 |
| DDIM (Song et al., 2020a) | 39 | 237.84 | 2.37 | 4.32 | 80.22 | 59.31 | - | 1.379(1.28×) |
| Learning-to-Cache (Ma et al., 2024a) | 50 | 233.26 | 2.62 | 4.50 | 79.40 | 59.15 | 23.39 | 1.419(1.25×) |
| HarmoniCa | 50 | **238.74** | **2.36** | **4.24** | **80.57** | **59.68** | **23.68** | 1.361(1.30×) |
| DDIM (Song et al., 2020a) | 20 | 224.37 | 3.52 | 4.96 | 78.47 | 58.08 | - | 0.658 |
| DDIM (Song et al., 2020a) | 14 | 201.83 | 5.77 | 6.61 | 75.14 | 55.08 | - | 0.466(1.41×) |
| Learning-to-Cache (Ma et al., 2024a) | 20 | 201.37 | 5.34 | 6.54 | 75.04 | 56.09 | 35.60 | 0.468(1.41×) |
| HarmoniCa | 20 | **206.57** | **4.88** | **5.91** | **75.20** | **58.74** | **37.50** | 0.456(1.44×) |
| DDIM (Song et al., 2020a) | 10 | 159.93 | 12.16 | 11.31 | 67.10 | 52.27 | - | 0.332 |
| DDIM (Song et al., 2020a) | 9 | 140.37 | 16.54 | 14.44 | 62.63 | 50.08 | - | 0.299(1.11×) |
| Learning-to-Cache (Ma et al., 2024a) | 10 | 145.09 | 14.59 | 11.58 | 64.03 | 52.06 | 19.11 | 0.279(1.19×) |
| HarmoniCa | 10 | **151.83** | **13.35** | **11.13** | **65.22** | **52.18** | **22.86** | 0.270(1.23×) |
| DiT-XL/2 512 × 512 (cfg = 1.5) | | | | | | | | |
| DDIM (Song et al., 2020a) | 20 | 184.47 | 5.10 | 5.79 | 81.77 | 54.50 | - | 3.356 |
| DDIM (Song et al., 2020a) | 16 | 173.31 | 6.47 | 6.67 | 81.10 | 51.30 | - | 2.688(1.25×) |
| Learning-to-Cache (Ma et al., 2024a) | 20 | 178.11 | 6.24 | 7.01 | 81.21 | 53.30 | 23.57 | 2.633(1.28×) |
| HarmoniCa | 20 | **179.84** | **5.72** | **6.61** | **81.33** | **55.80** | **25.98** | 2.574(1.30×) |



Comparison with Learning-to-Cache with the increase of the speedup ratio.

Comparison with quantization and pruning.

| Method | T | IS↑ | FID↓ | sFID↓ | Latency(s)↓ | Latency(s)↓* |
|---|---|---|---|---|---|---|
| DDIM (Zhang et al., 2022) | 20 | 224.37 | 3.52 | 4.96 | 0.658 | 1.217 |
| EfficientDM (He et al., 2024) | 20 | 172.70 | 6.10 | 4.55 | 0.591(1.11×) | 0.842(1.45×) |
| PTQ4DiT (Wu et al., 2024) | 20 | 17.06 | 71.82 | 23.16 | 0.577(1.14×) | 0.839(1.45×) |
| Diff-Pruning (Fang et al., 2023) | 20 | 180.18 | 8.22 | 6.20 | 0.458(1.44×) | 0.813(1.50×) |
| HarmoniCa | 20 | **206.57** | **4.88** | **5.91** | 0.456(1.44×) | 0.815(1.49×) |

Combination with quantization.

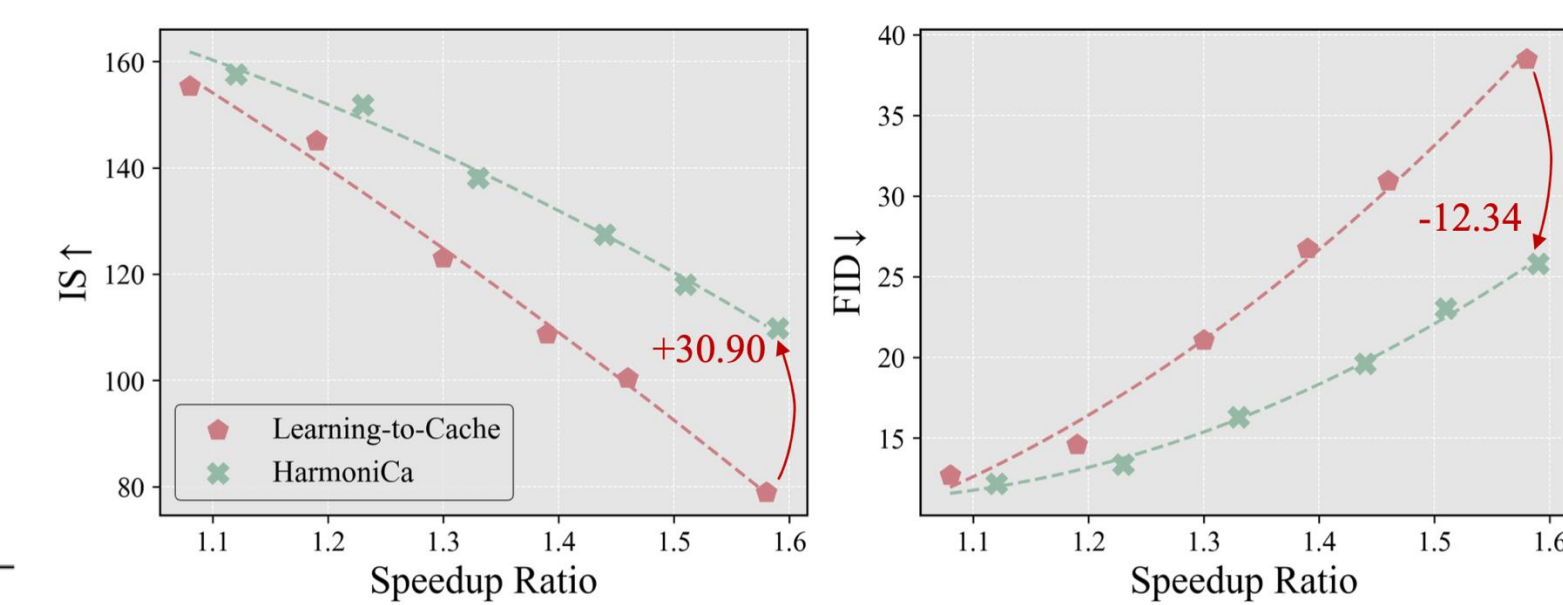| Method | IS↑/CLIP↑ | FID↓ | sFID↓ | CUR(%)↑ | Latency(s)↓ | #Size(GB)↓ |
|---|---|---|---|---|---|---|
| DiT-XL/2 256 × 256 (cfg = 1.5) | | | | | | |
| EfficientDM (He et al., 2024) | 172.70 | 6.10 | 4.55 | - | 0.591(1.11×) | 0.64(3.93×) |
| w/ HarmoniCa (β = 4e⁻⁸) | 168.16 | 6.48 | 4.32 | 26.25 | 0.473(1.40×) | 0.64(3.93×) |
| PIXART-α 256 × 256 (cfg = 4.5) | | | | | | |
| EfficientDM (He et al., 2024) | 30.09 | 34.84 | 30.34 | - | 0.469(1.18×) | 0.59(1.98×) |
| w/ HarmoniCa | 30.15 | 34.96 | 30.55 | 53.34 | 0.299(1.85×) | 0.59(1.98×) |
| PIXART-α 512 × 512 (cfg = 4.5) | | | | | | |
| EfficientDM (He et al., 2024) | 30.71 | 25.82 | 41.64 | - | 0.461(1.20×) | 0.59(1.98×) |
| w/ HarmoniCa | 30.75 | 26.15 | 41.99 | 53.11 | 0.281(1.97×) | 0.59(1.98×) |

(Left) Comparison with addition caching methods on U-ViT; (Right) Training costs for DiT-XL/2 256 ×256.

| Method | T | FID↓ | Latency(s)↓ |
|---|---|---|---|
| DPM-Solver (Lu et al., 2022a) | 20 | 2.57 | 7.60 |
| Faster Diffusion (Li et al., 2023a) | 20 | 2.82 | 5.95(1.28×) |
| DeepCache (Ma et al., 2024b) | 20 | 2.70 | 4.68(1.62×) |
| HarmoniCa | 20 | **2.61** | **4.60(1.65×)** |

| Method | #Images | Time(h) | Memory(GB/GPU) |
|---|---|---|---|
| Learning-to-Cache | 1.22M | 2.15 | 33.33 |
| SDT+$\mathcal{L}_{LTC}^{(t)}$ | 0 | 1.47 | 33.28 |
| HarmoniCa | 0 | 1.63 | 33.28 |



(a) 20-step DPM-Solver
(b) 15-step DPM-Solver (1.36×)
(c) FORA (1.53×)
(d) HarmoniCa (1.64×)

Visualization results.