# BitsFusion: 1.99 bits Weight Quantization of Diffusion Model

Yang Sui,  Yanyu Li,  Anil Kag,  Yerlan Idelbayev,  Junli Cao,  Ju Hu,

Dhritiman Sagar,  Bo Yuan,  Sergey Tulyakov,  Jian Ren

Snap Inc.     Rutgers University

# Text-to-Image Diffusion Model

## Stable Diffusion

# Challenge: Storage Size

**Quantization**

Our goal: Extremely Low-bit Text-to-Image Diffusion Model (i.e., 1.99 bits UNet)

# Overview of BitsFusion Pipeline

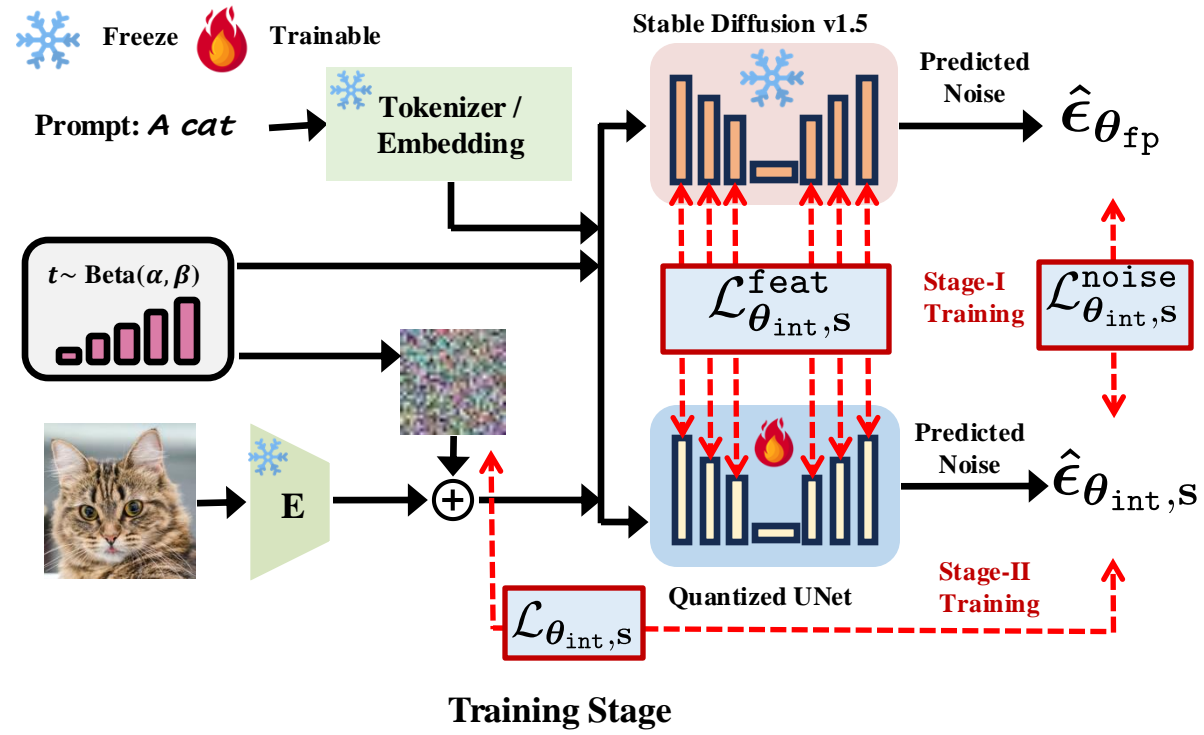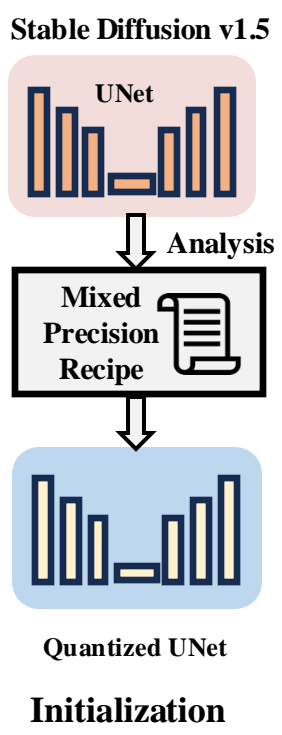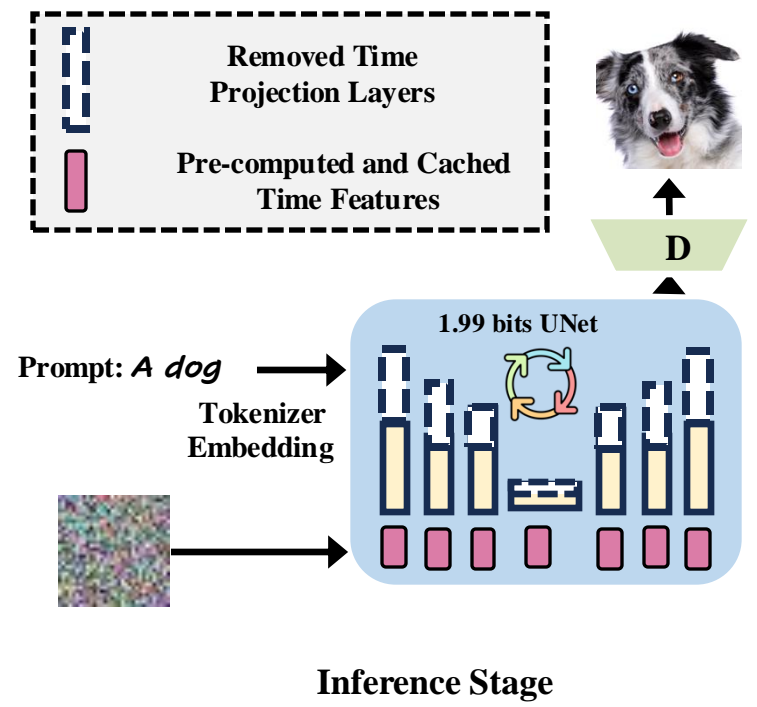**Stable Diffusion v1.5**

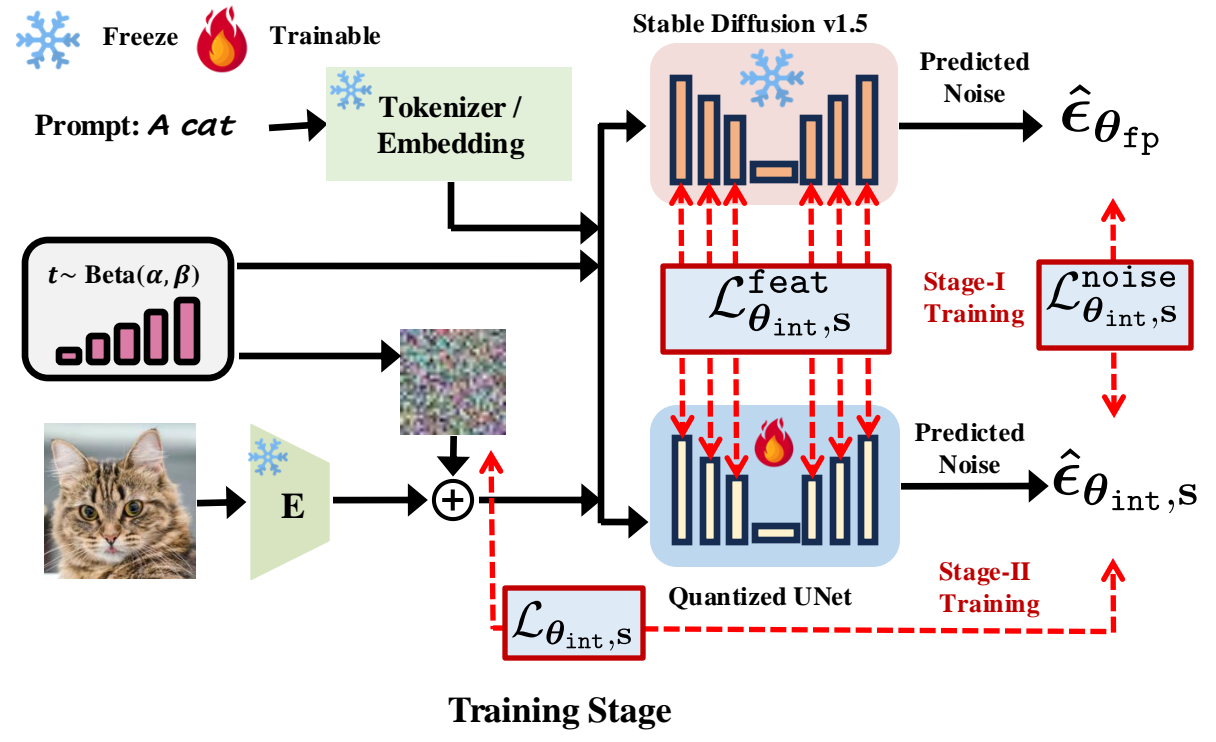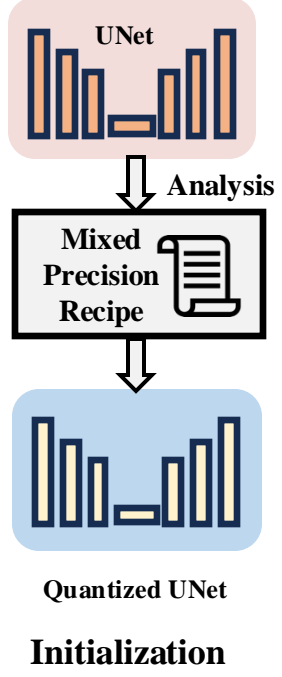UNet

Analysis

Mixed
Precision
Recipe

Quantized UNet

**Initialization**

# Overview of BitsFusion Pipeline



**Stable Diffusion v1.5**

UNet

$\Downarrow$ **Analysis**

**Mixed Precision Recipe**

$\Downarrow$

**Quantized UNet**

**Initialization**

❄️ **Freeze** 🔥 **Trainable**

**Prompt: *A cat*** → ❄️ **Tokenizer / Embedding**

$t \sim \text{Beta}(\alpha, \beta)$

**E**

$\oplus$

**Stable Diffusion v1.5** ❄️

**Predicted Noise** $\hat{\epsilon}_{\theta_{\text{fp}}}$

$\mathcal{L}^{\text{feat}}_{\theta_{\text{int}}, s}$

**Stage-I Training** $\mathcal{L}^{\text{noise}}_{\theta_{\text{int}}, s}$

🔥 **Quantized UNet**

**Predicted Noise** $\hat{\epsilon}_{\theta_{\text{int}}, s}$

$\mathcal{L}_{\theta_{\text{int}}, s}$ **Stage-II Training**

**Training Stage**

# Overview of BitsFusion Pipeline



**Initialization**

Stable Diffusion v1.5

UNet

Analysis

Mixed Precision Recipe

Quantized UNet

**Training Stage**

❄️ Freeze 🔥 Trainable

Prompt: *A cat*

Tokenizer / Embedding

Stable Diffusion v1.5

Predicted Noise $\hat{\epsilon}_{\theta_{\mathrm{fp}}}$

$t \sim \mathrm{Beta}(\alpha, \beta)$

E

$\mathcal{L}^{\mathtt{feat}}_{\theta_{\mathrm{int}},\mathbf{s}}$

Stage-I Training $\mathcal{L}^{\mathtt{noise}}_{\theta_{\mathrm{int}},\mathbf{s}}$

Predicted Noise $\hat{\epsilon}_{\theta_{\mathrm{int}},\mathbf{s}}$

Quantized UNet

$\mathcal{L}_{\theta_{\mathrm{int}},\mathbf{s}}$ Stage-II Training

**Inference Stage**

Removed Time Projection Layers

Pre-computed and Cached Time Features

D

Prompt: *A dog*

Tokenizer Embedding
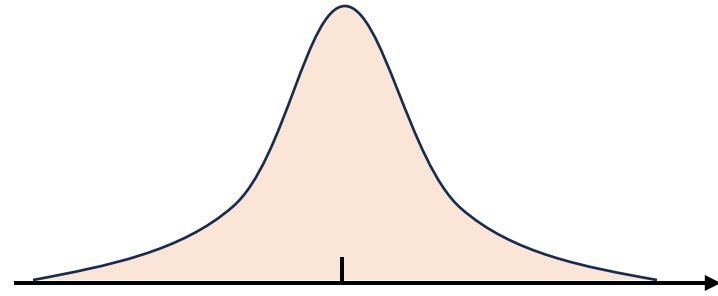
1.99 bits UNet

# Outline

- Mixed Precision Strategy
  - Per-Layer Quantization Error Analysis
  - Deciding the Mixed Precision

- Training Extreme Low-bit Diffusion Model
  - Initialization Schemes
  - Two-stage Training Pipeline

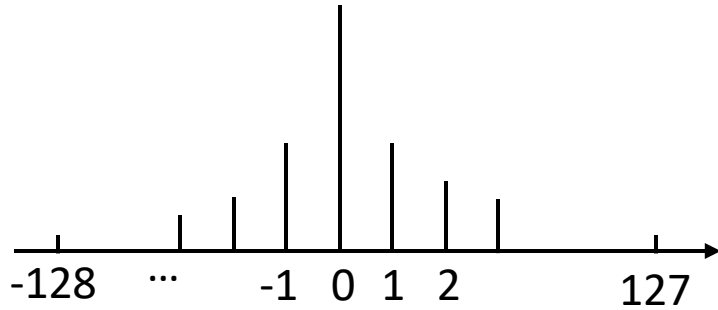- Results

# Mixed Precision Strategy

# Quantization

32 bits
Floating-point Weights

8 bits
Integer Weights

-128   ...   -1   0   1   2      127
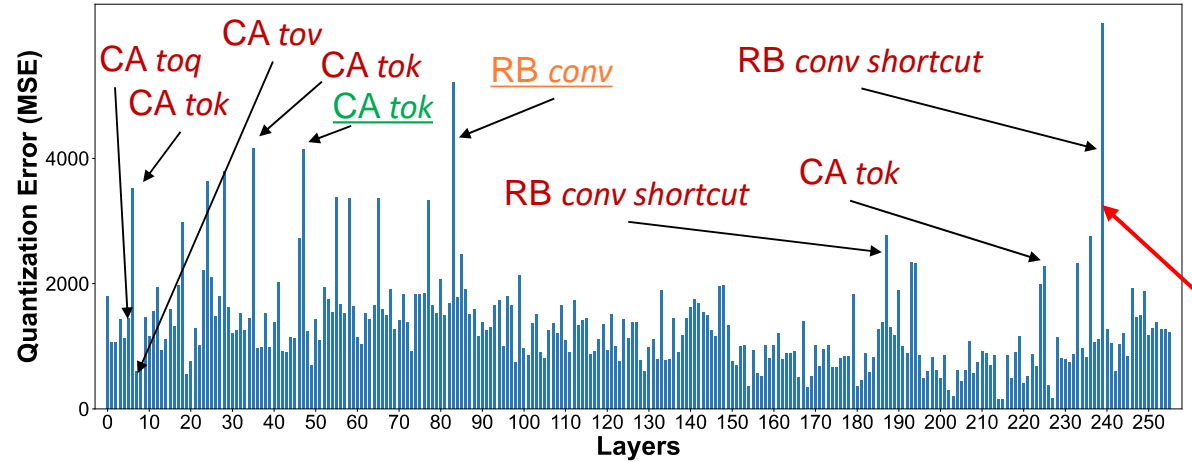
Saving Storage Size

# Per-Layer Quantization Error

## Mixed Precision Precision

Measure the impact when quantizing each layer:

- Quantize each single layer to 1,2,3 bits

# Per-Layer Quantization Error

**Mixed Precision Precision**

Measure the impact when quantizing each layer:

- Quantize each single layer to 1,2,3 bits
- Generate images from quantized models

# Per-Layer Quantization Error

**Mixed Precision Precision**

Measure the impact when quantizing each layer:

- Quantize each single layer to 1,2,3 bits
- Generate images from quantized models
- Calculate the metrics compared to full-precision model: MSE, CLIP Score, PSNR, LPIPS

# Per-Layer Quantization Error

## Mixed Precision Precision

# Per-Layer Quantization Error

**Which metrics should we use?**

Pearson correlation (absolute value) between different metrics

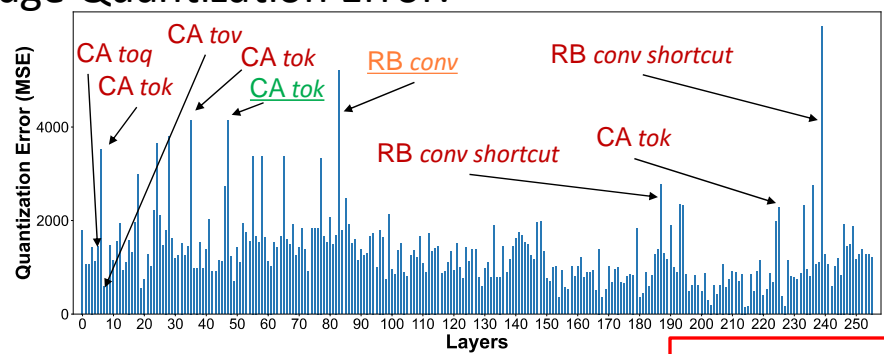|  | MSE vs. PSNR | MSE vs. LPIPS | MSE vs. CLIP Score |
|---|---|---|---|
| 1 bit | 0.870 | 0.984 | 0.733 |
| 2 bit | 0.882 | 0.989 | 0.473 |
| 3 bit | 0.869 | 0.991 | 0.535 |

Observation 1: MSE, PSNR, and LPIPS show strong correlation and they correlate well with the visual perception of image quality.

Conclusion: We adopt MSE as our main quantitative metric to represent the PSNR and LPIPS.

# Per-Layer Quantization Error

## Which metrics should we use?

Average Quantization Error:



MSE: *CA tok* < *RB conv*
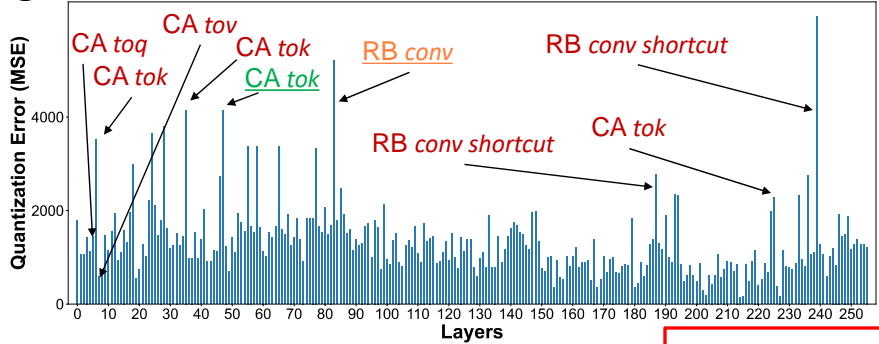CLIP Score drop: *CA tok* > *RB conv*

Observation 2: Although some layers show smaller MSE, they may experience larger semantic degradation, as reflected in larger CLIP score changes.
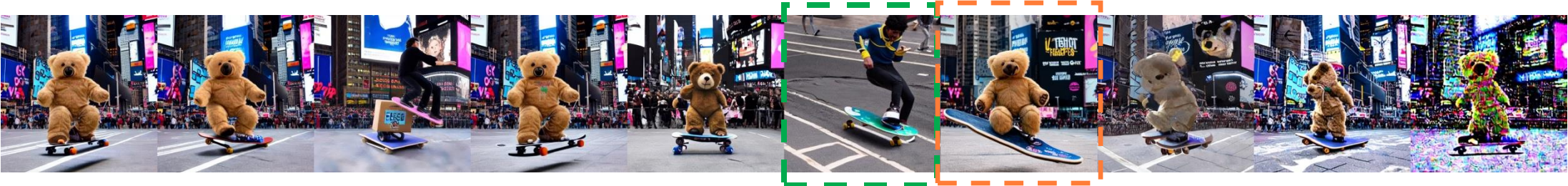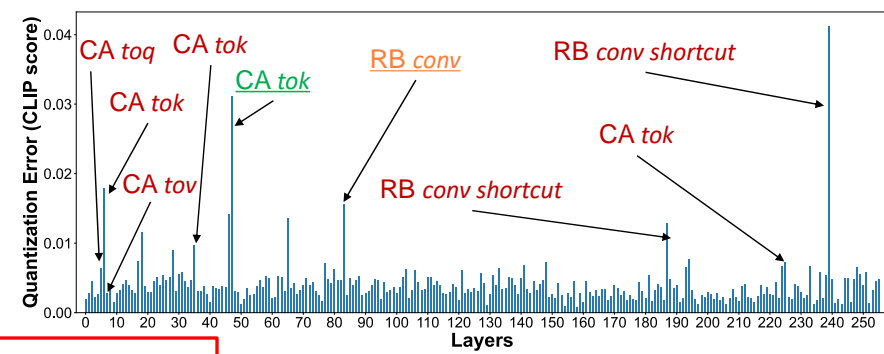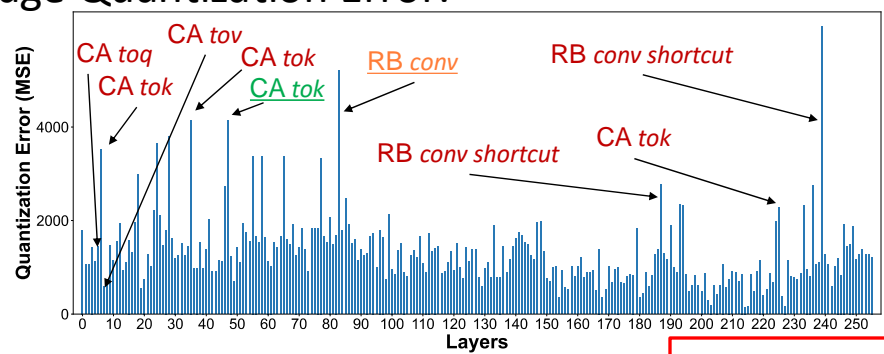
# Per-Layer Quantization Error

## Which metrics should we use?

Average Quantization Error:



MSE: CA *tok* < RB *conv*

CLIP Score drop: CA *tok* > RB *conv*

Observation 2: Although some layers show smaller MSE, they may experience larger semantic degradation, as reflected in larger CLIP score changes.



*A teddy bear on a skateboard in Times Square, doing tricks on a cardboard box ramp*
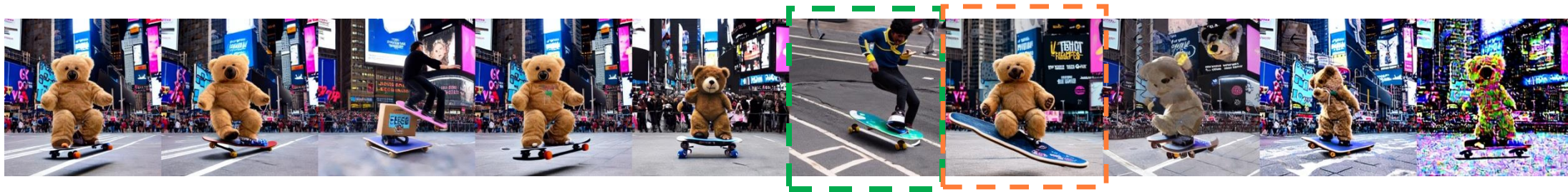
# Per-Layer Quantization Error

## Which metrics should we use?

Average Quantization Error:



MSE: CA *tok* < RB *conv*

CLIP Score drop: CA *tok* > RB *conv*

Observation 2: Although some layers show smaller MSE, they may experience larger semantic degradation, as reflected in larger CLIP score changes.



*A teddy bear on a skateboard in Times Square, doing tricks on a cardboard box ramp*

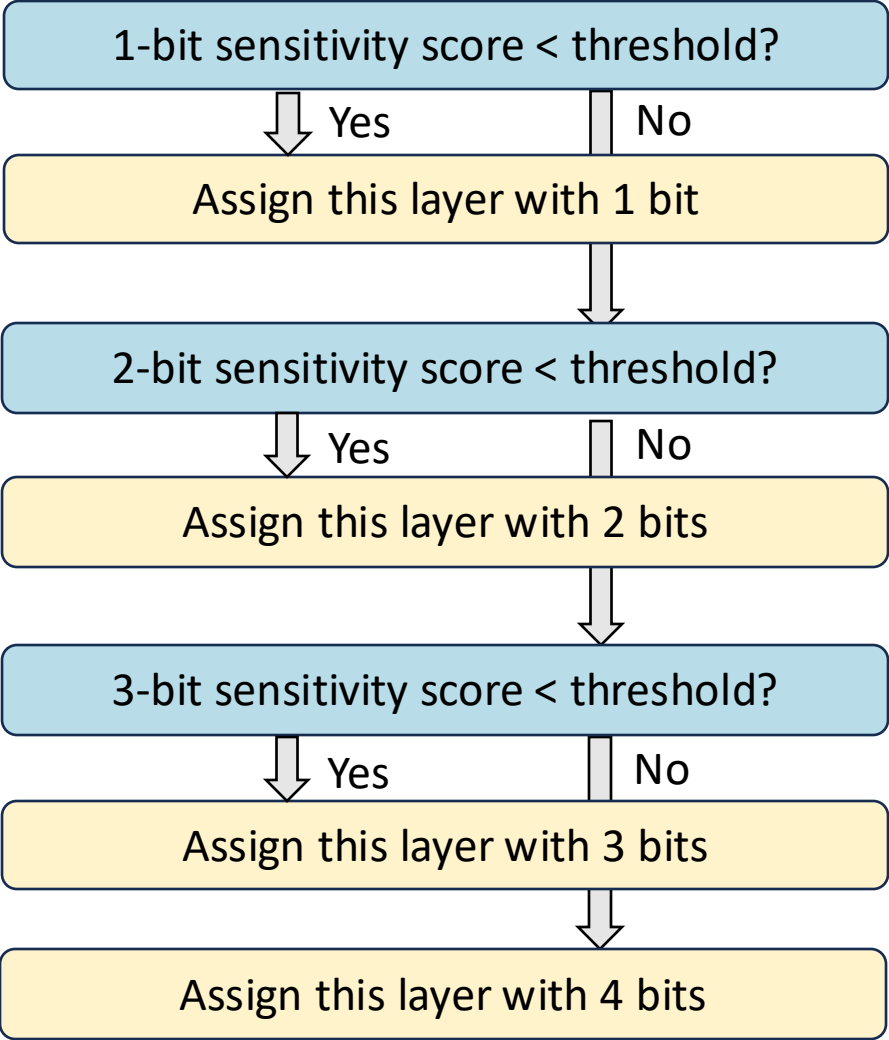Adopt CLIP score as our complementary quantitative metrics.

**Per-Layer Quantization Error**

Which metrics should we use?

Conclusion: MSE, CLIP Score

# Deciding the Optimal Precision

## 1. Assign bits based on MSE (For one layer):

Sensitivity score
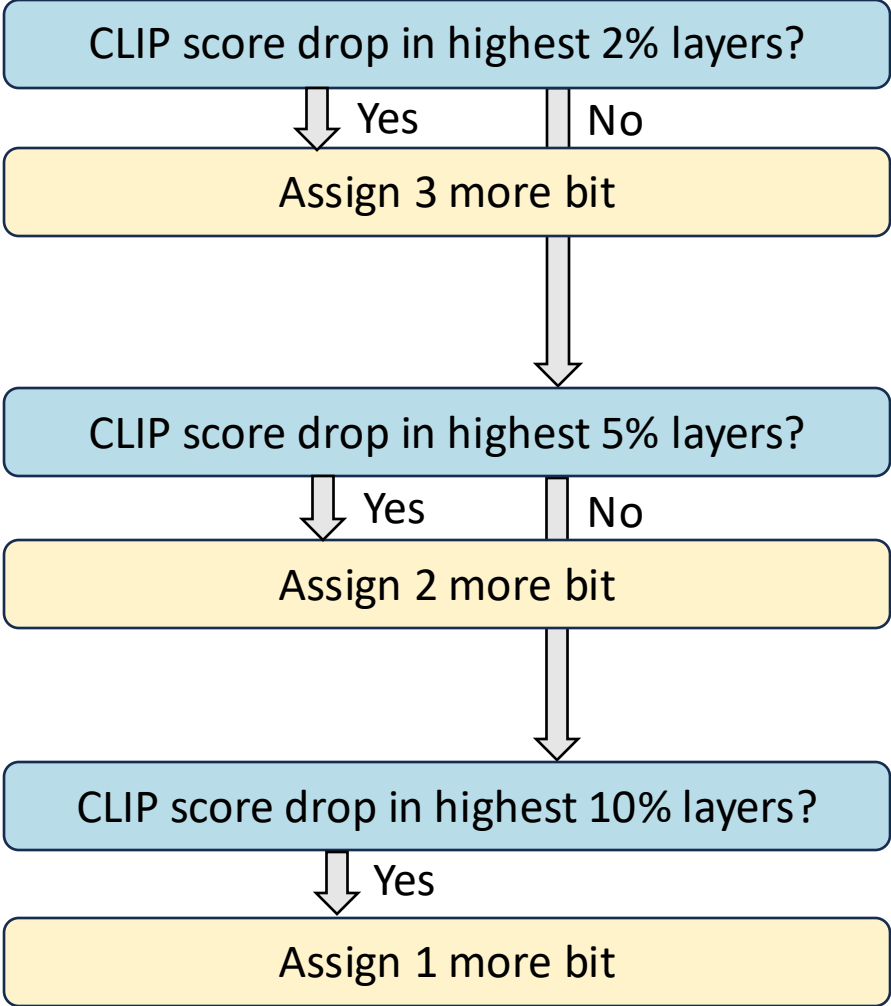
$$\mathcal{S}_{i,b} = M_{i,b} N_i^{-\eta}$$

$M_{i,b}$ : MSE

$N_i$ : Parameter size

$\eta$ : Parameter size factor

1-bit sensitivity score < threshold?

↓ Yes     No

Assign this layer with 1 bit

2-bit sensitivity score < threshold?

↓ Yes     No

Assign this layer with 2 bits

3-bit sensitivity score < threshold?

↓ Yes     No

Assign this layer with 3 bits

Assign this layer with 4 bits

# Deciding the Optimal Precision

2. Adjust bits based on CLIP scores (For one layer):

# Initialization

# Initialization

## Time Embedding Pre-computing and Caching

Stable Diffusion

# Initialization

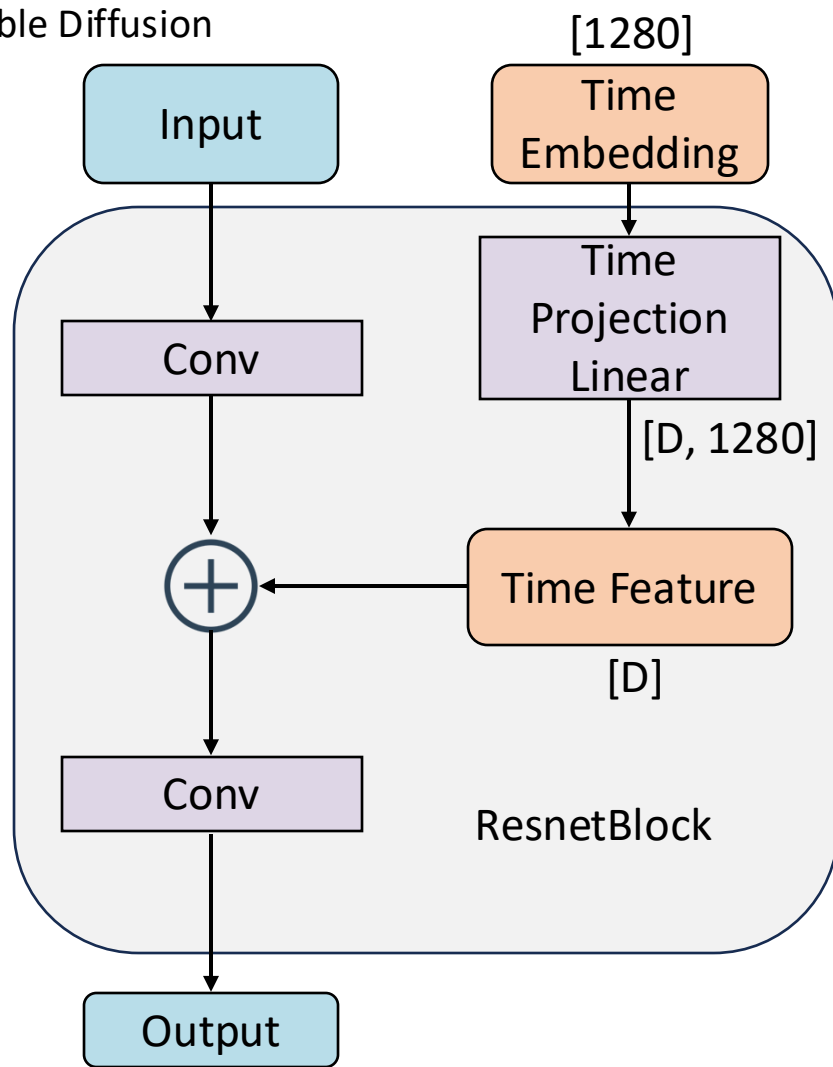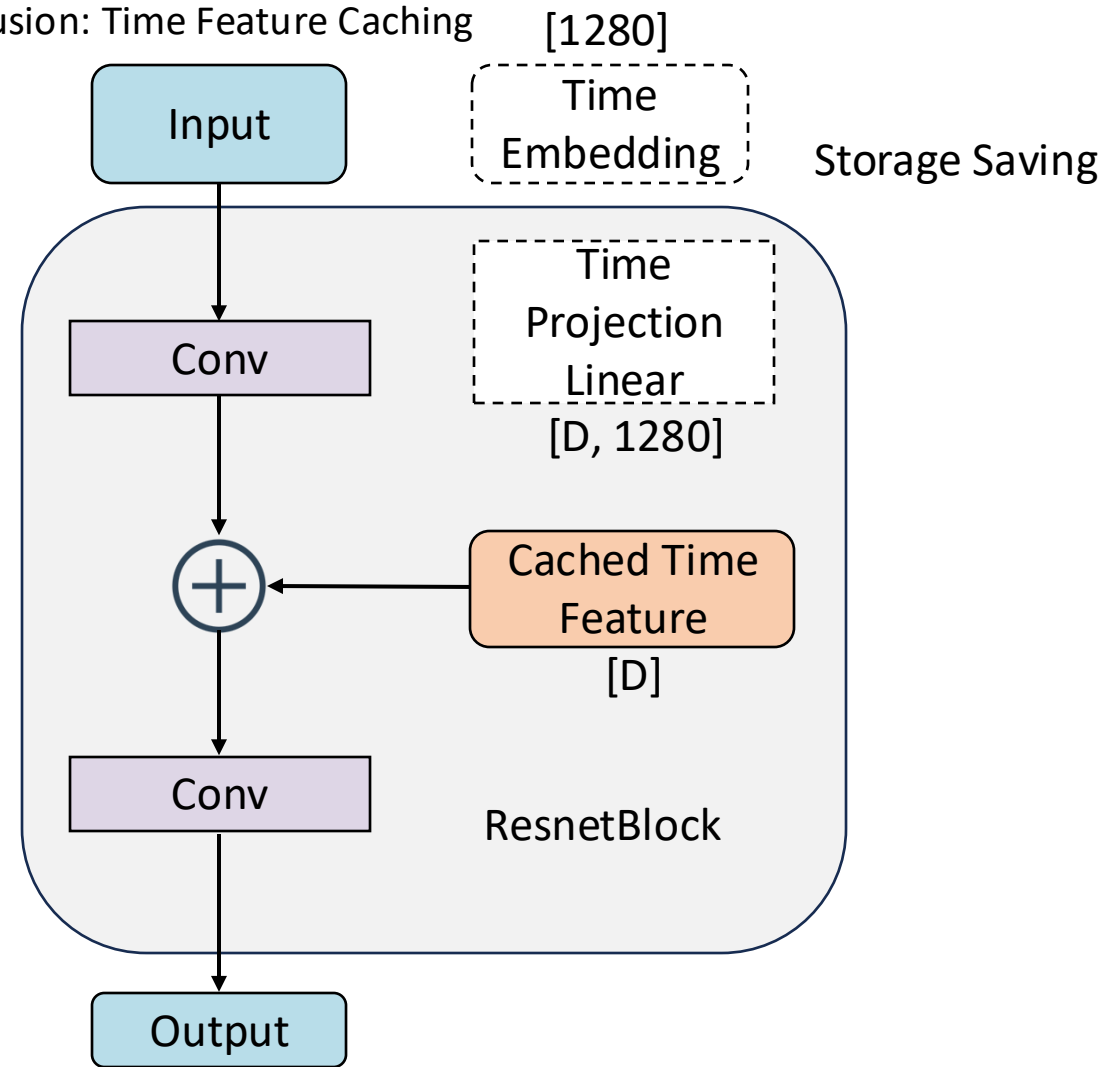## Time Embedding Pre-computing and Caching

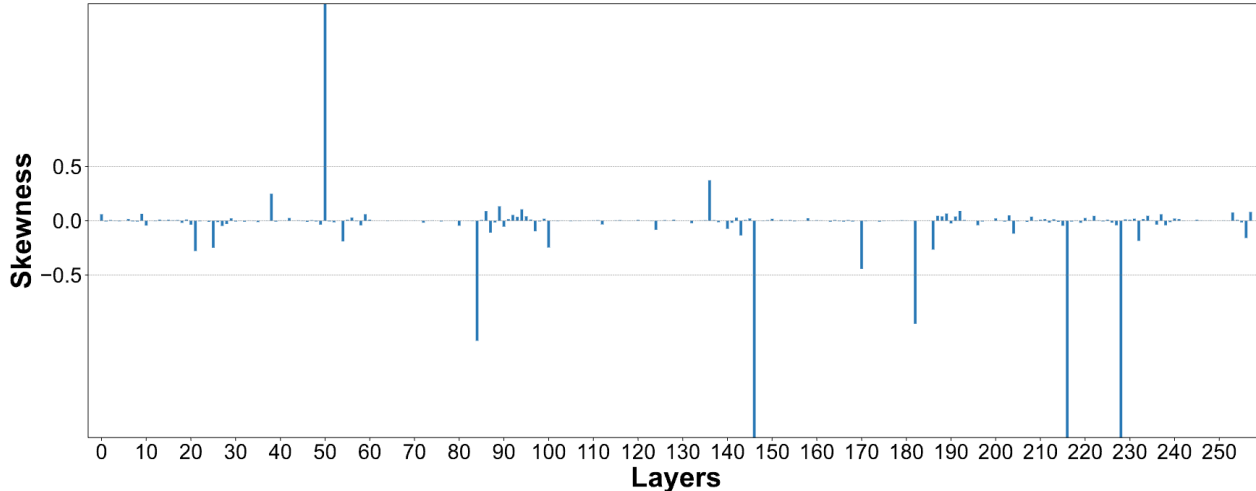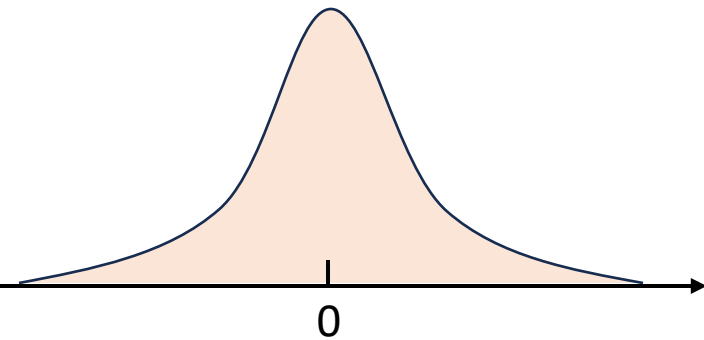# Initialization

## Time Embedding Pre-computing and Caching



During inference stage, only caching T time features (T is the sampling steps, T <= 50 in stable diffusion).

# Initialization

## Adding Balance Integer

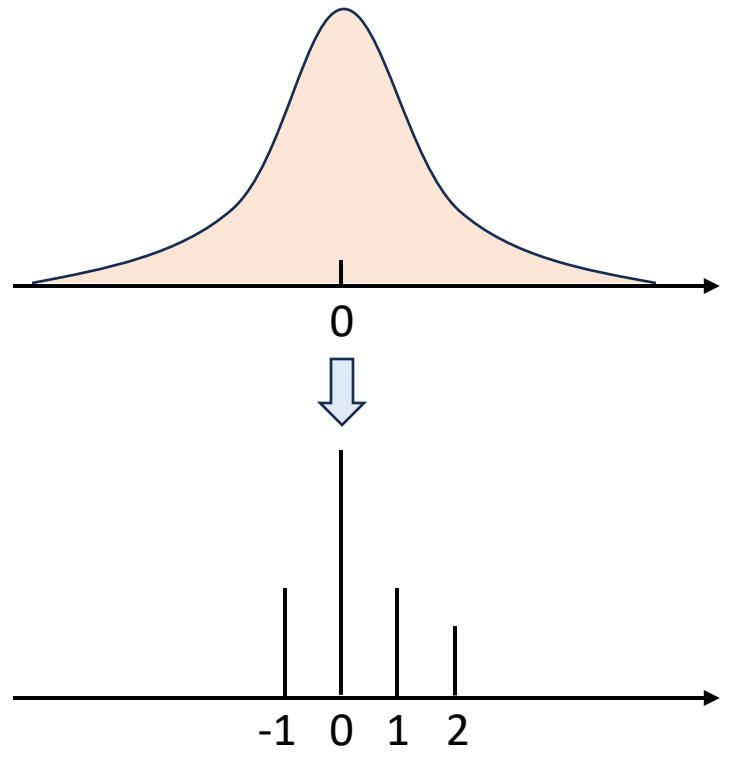Is weight distribution symmetric in Stable Diffusion?



97% of layers exhibiting skewness between [-0.5, 0.5]

Weight Distribution of layers in Stable Diffusion is symmetric

# Initialization
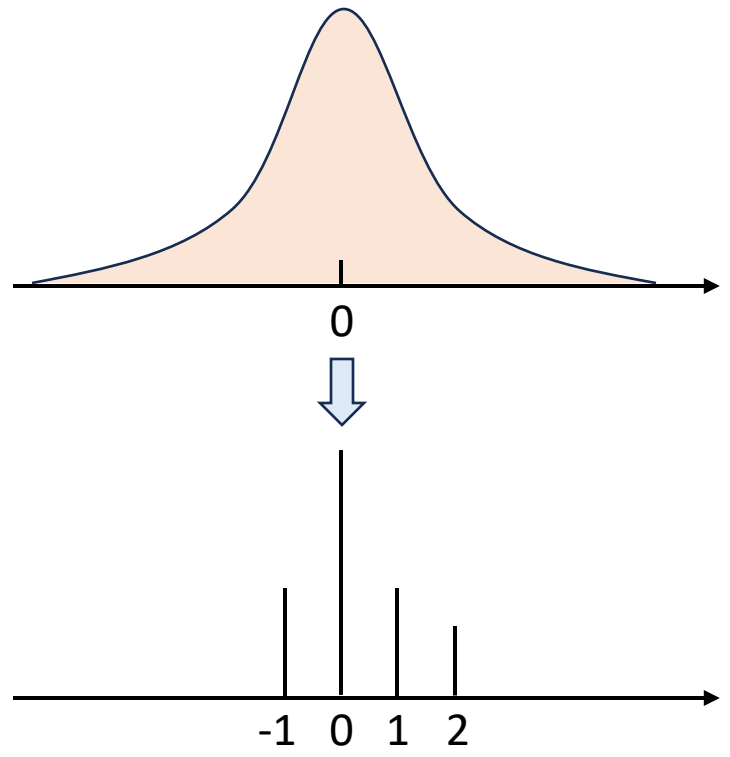
## Adding Balance Integer

2-bit mapping



0

-1  0  1  2

Unbalance in low-bit (e.g., 2 bits) quantization

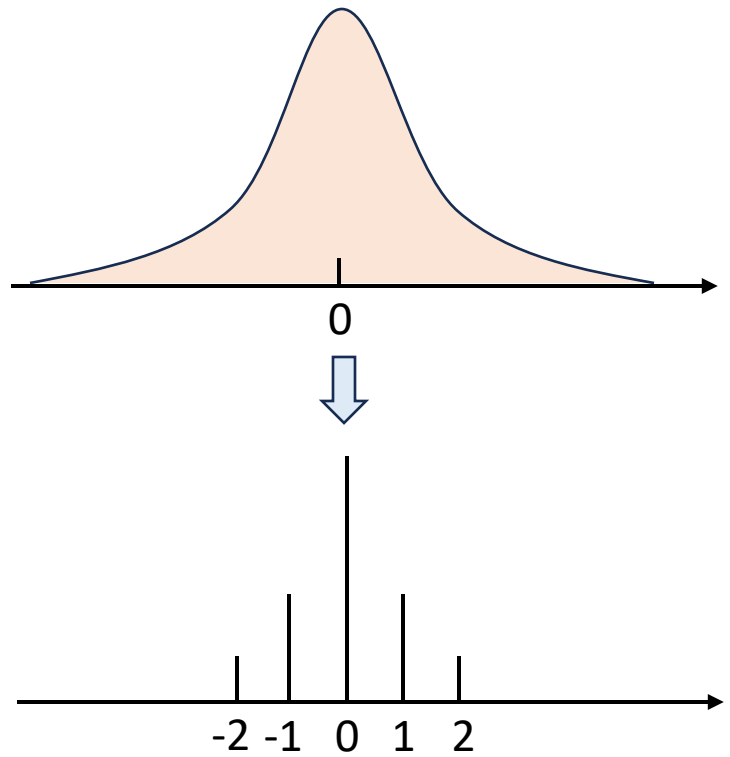# Initialization

## Adding Balance Integer



2-bit mapping

0

-1  0  1  2

Unbalance in low-bit (e.g., 2 bits) quantization

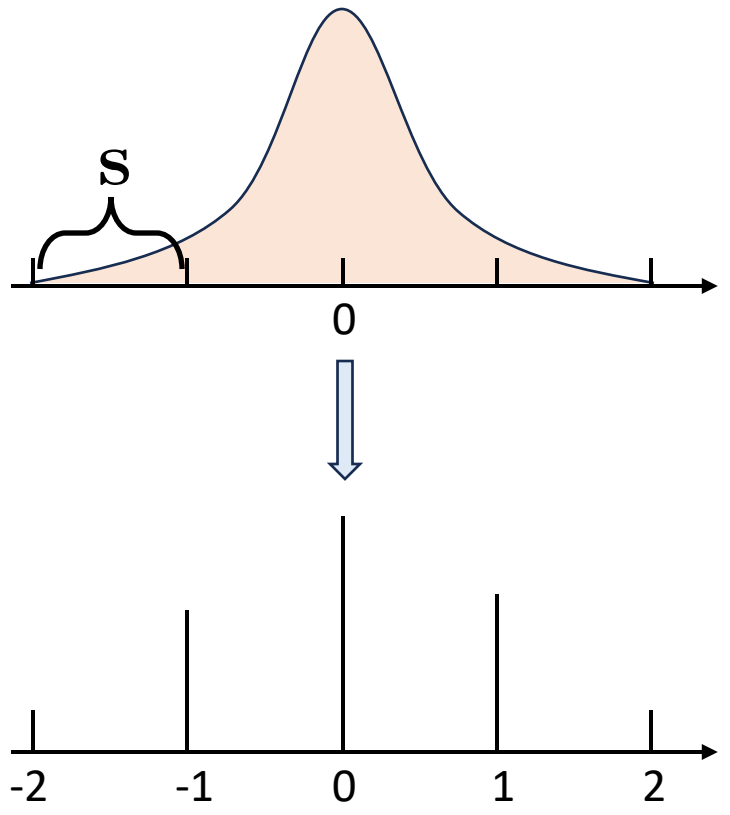Add one value

2-bit mapping

0

-2 -1  0  1  2

balanced values

# Initialization

## Scaling Factor Initialization via Alternating Optimization

2-bit mapping

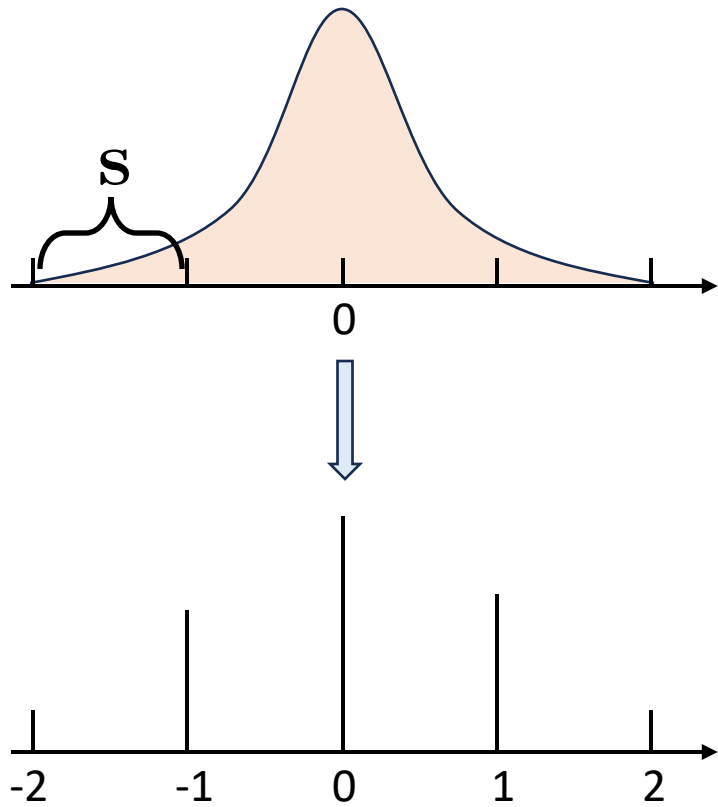Min-Max initialization mapping



Drawback:
Large quantization error in
low-bit (e.g., 2 bits) quantization

# Initialization

## Scaling Factor Initialization via Alternating Optimization
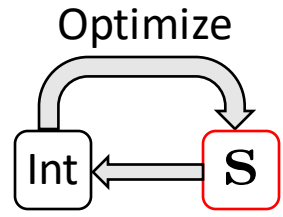
2-bit mapping

Min-Max initialization mapping



Minimize Initial Quantization Error
By Updating Scaling Factor

Optimize

Drawback:
Large quantization error in
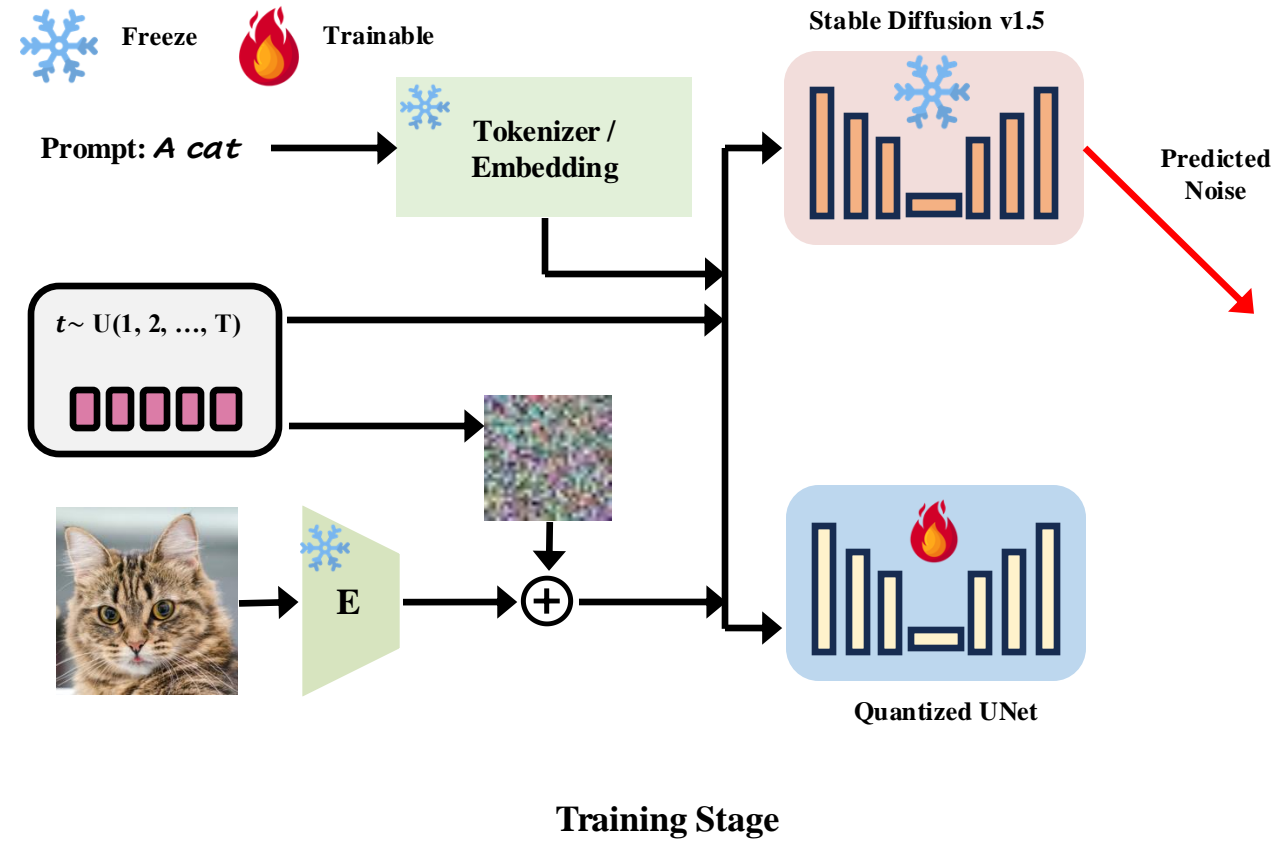low-bit (e.g., 2 bits) quantization
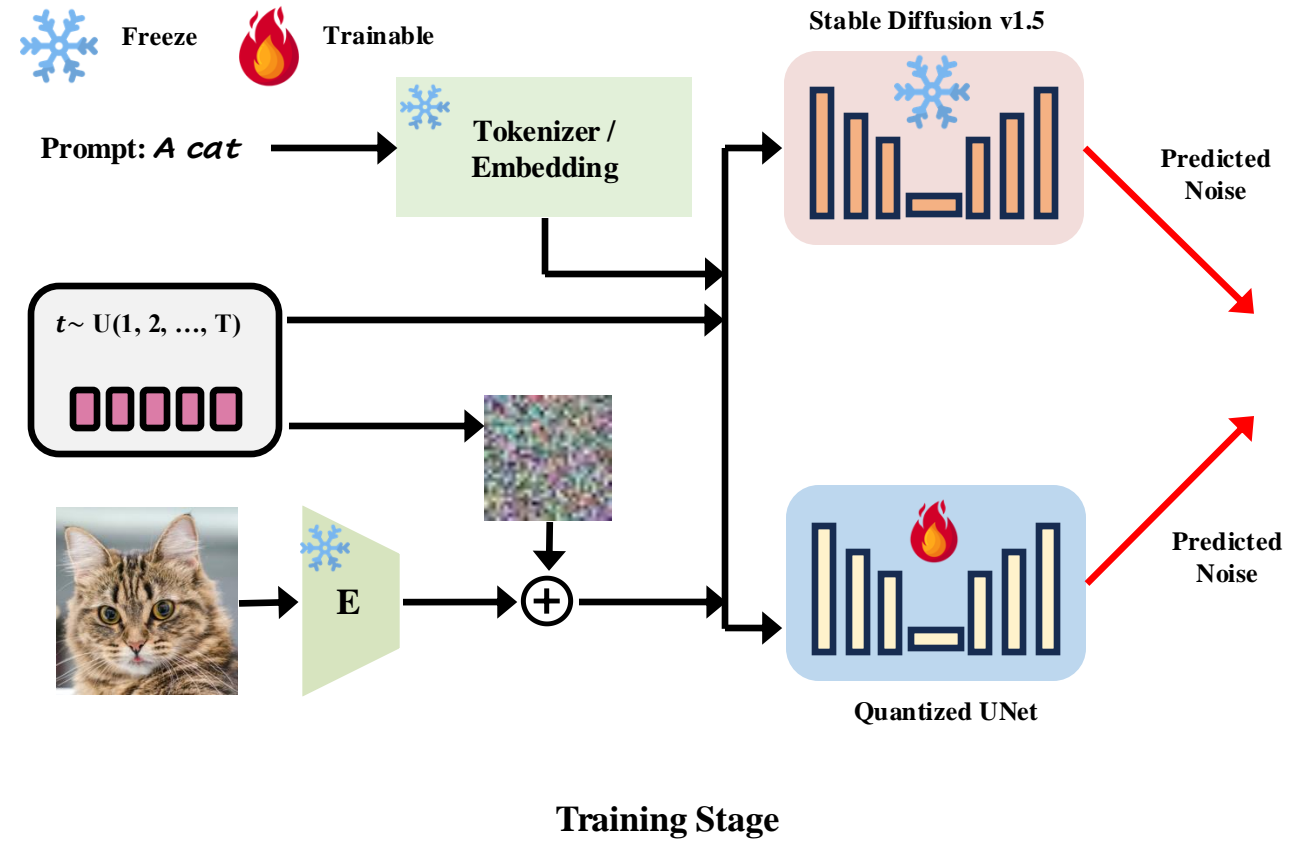
# Two-stage Training Pipeline

# Stage-I Training

### Loss Function

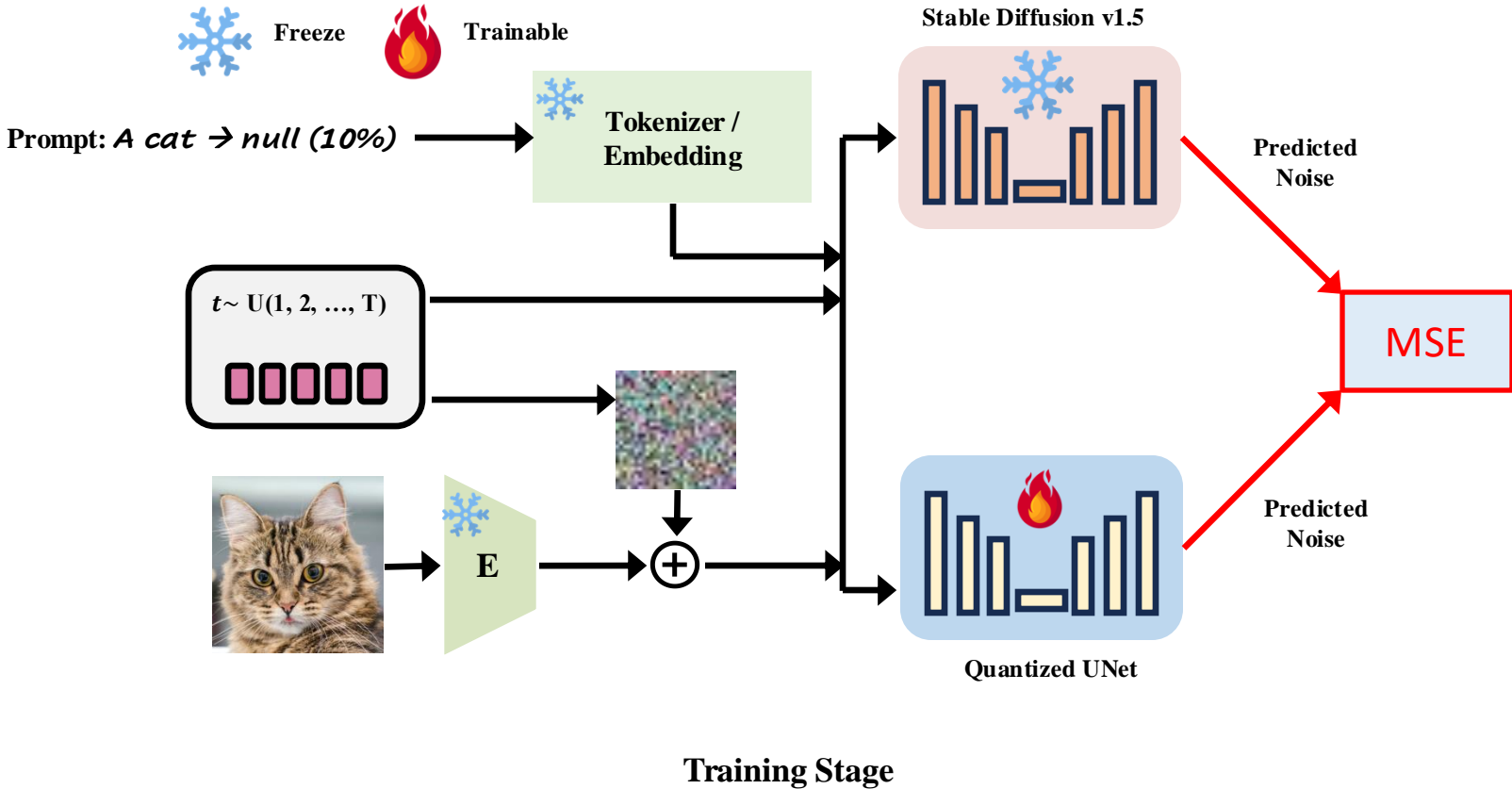**CFG-aware Quantization Distillation**

# Stage-I Training

**Loss Function**

**CFG-aware Quantization Distillation**



Training Stage

# Stage-I Training

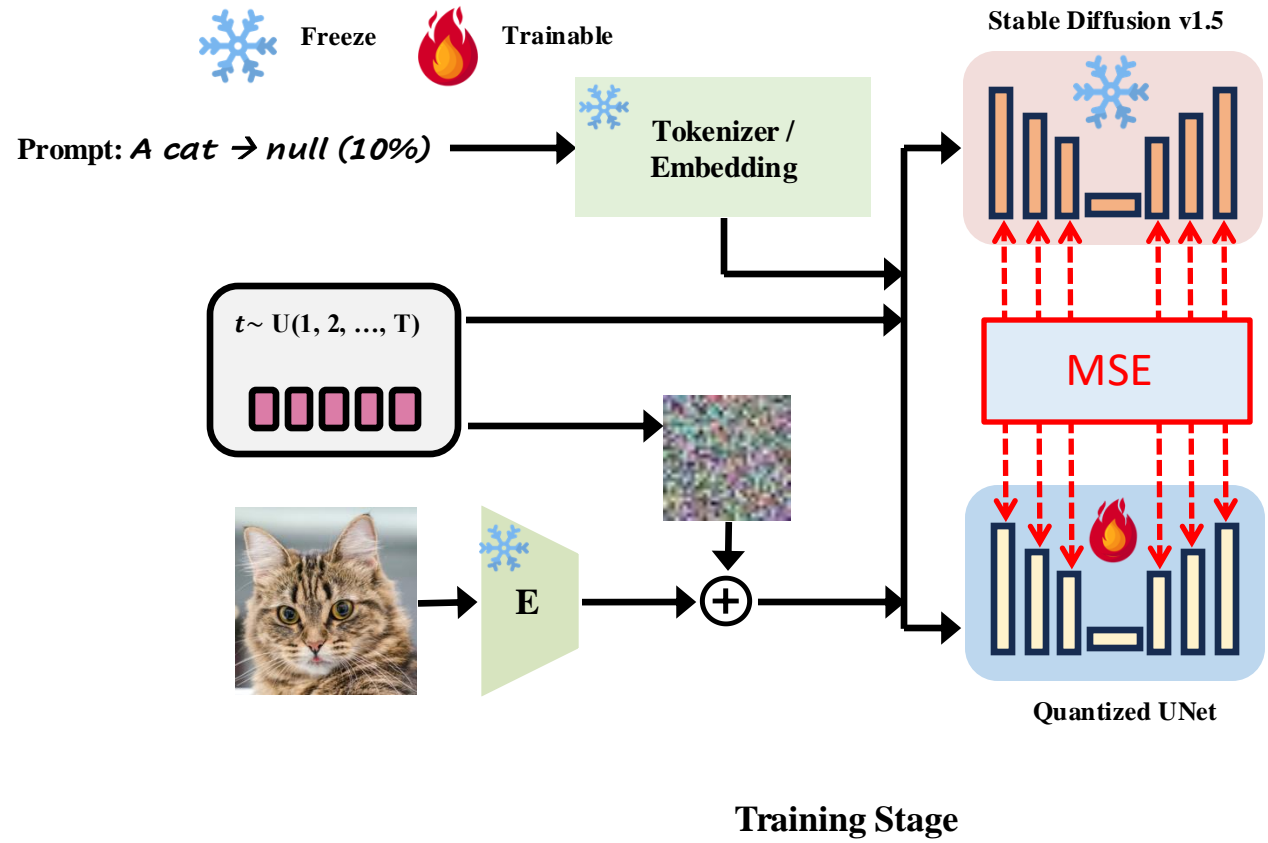### Loss Function

**CFG-aware Quantization Distillation**



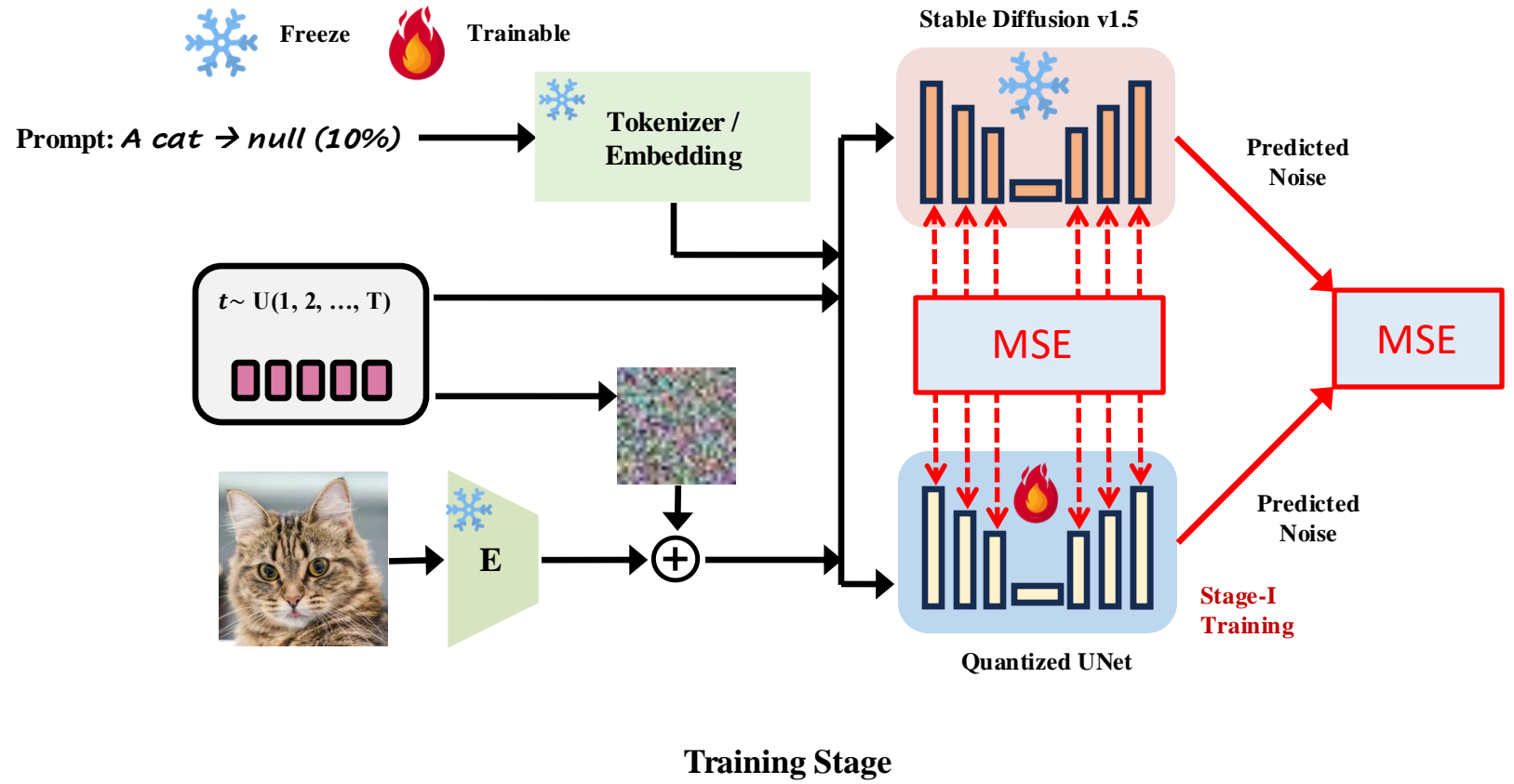**Training Stage**

# Stage-I Training

## Loss Function



Feature Distillation

# Stage-I Training

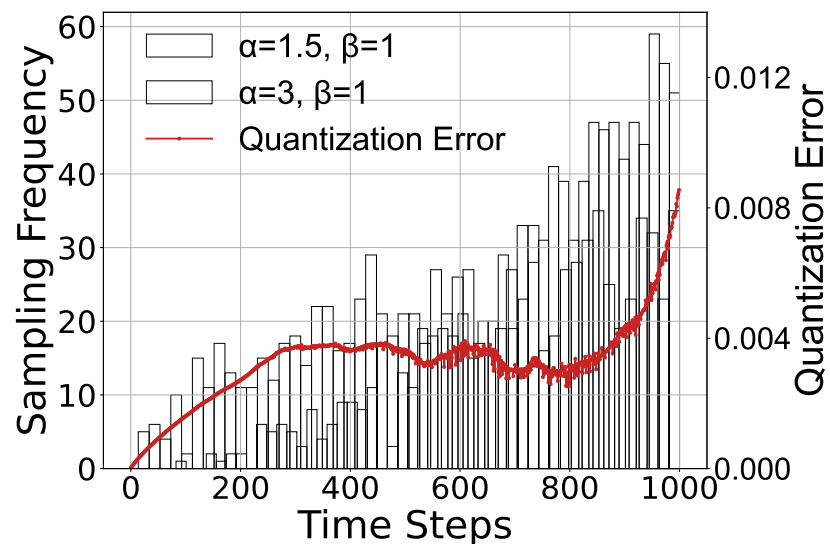### Loss Function

## Overall Distillation

# Stage-I Training

## Quantization Error-aware Time Step Sampling

Motivation: Different Quantization Error at Different Time Steps

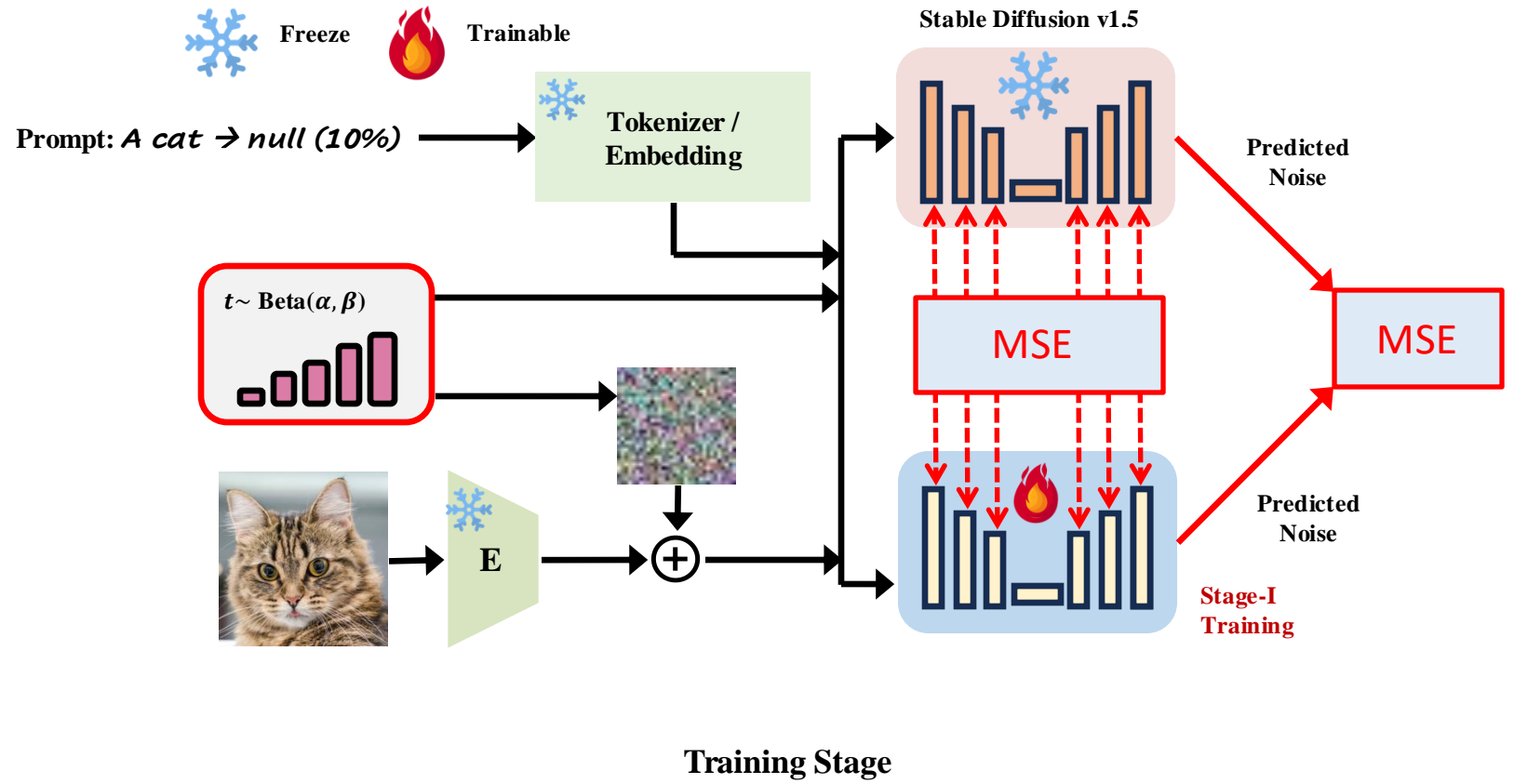Quantization error of predicted latent features between quantized model and FP model



Observation: the quantization error keeps increasing as the time steps approach t = 999.

Solution: Sample more time steps exhibiting the larger quantization errors near t = 999 by Beta distribution.
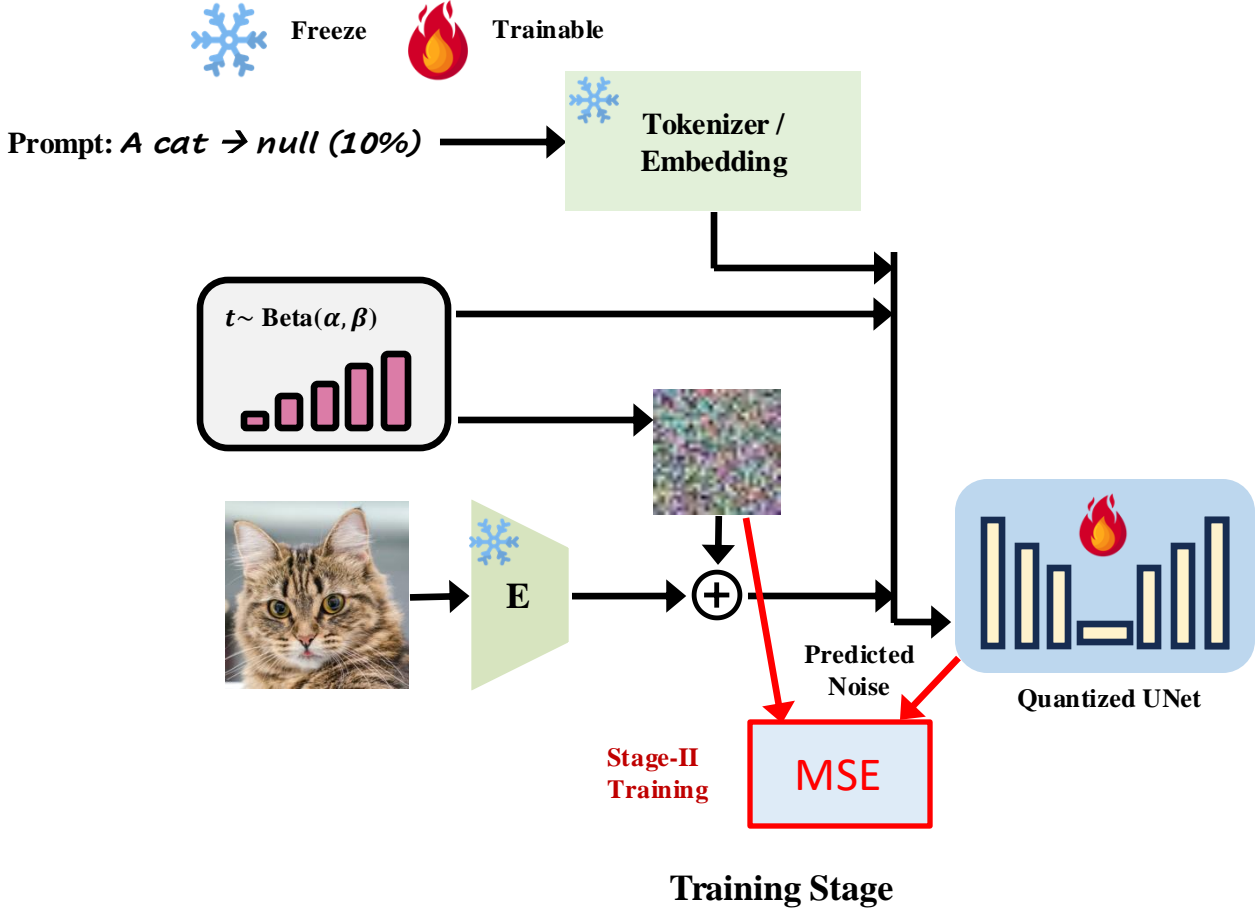
# Stage-I Training

## Loss Function

## Overall Distillation



**Training Stage**

# Stage-II Training

## Fine-tuning with Noise Prediction

# Results

# Results

## Generated Images
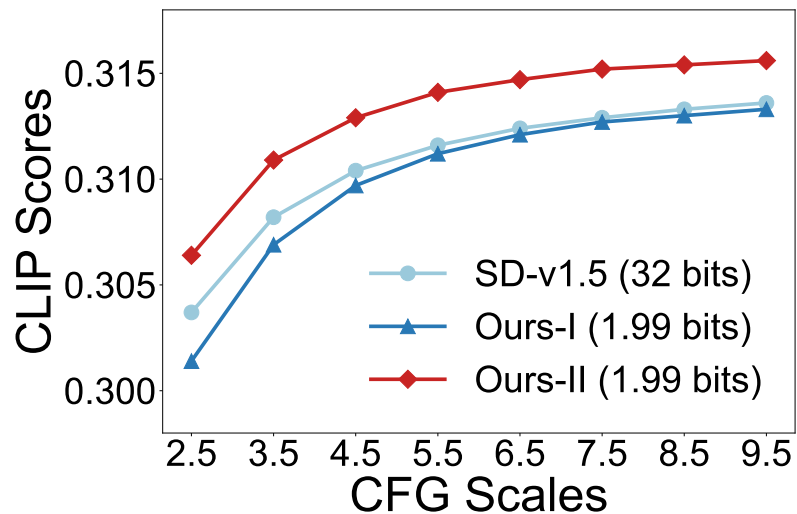
Sampler: PNDM

Steps: 50

Seed: 1024
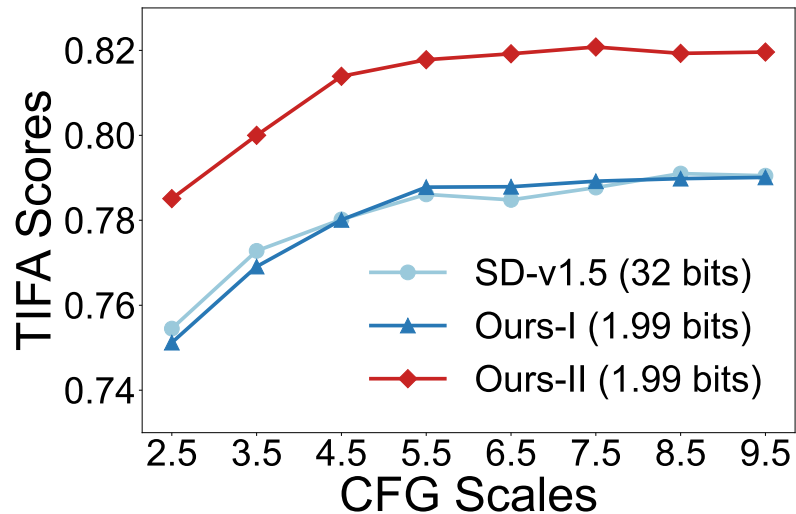
Stable Diffusion v1.5, 32 bits
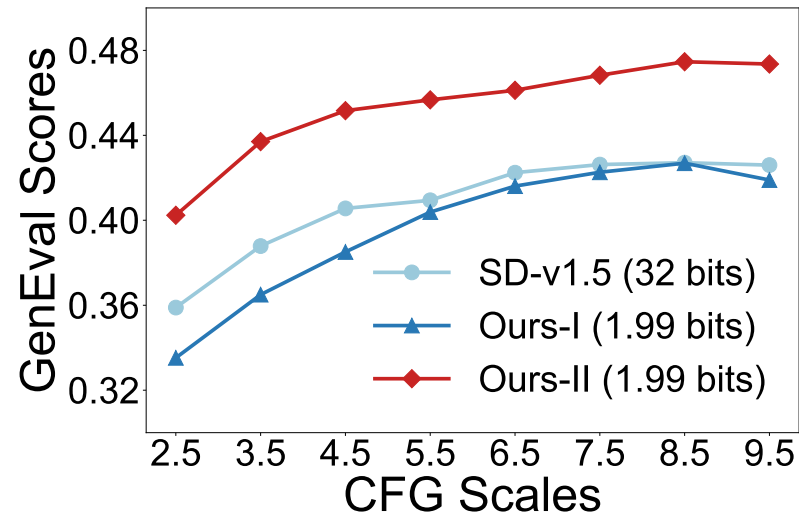


BitsFusion, 1.99 bits

# Results

## Quantitative performance



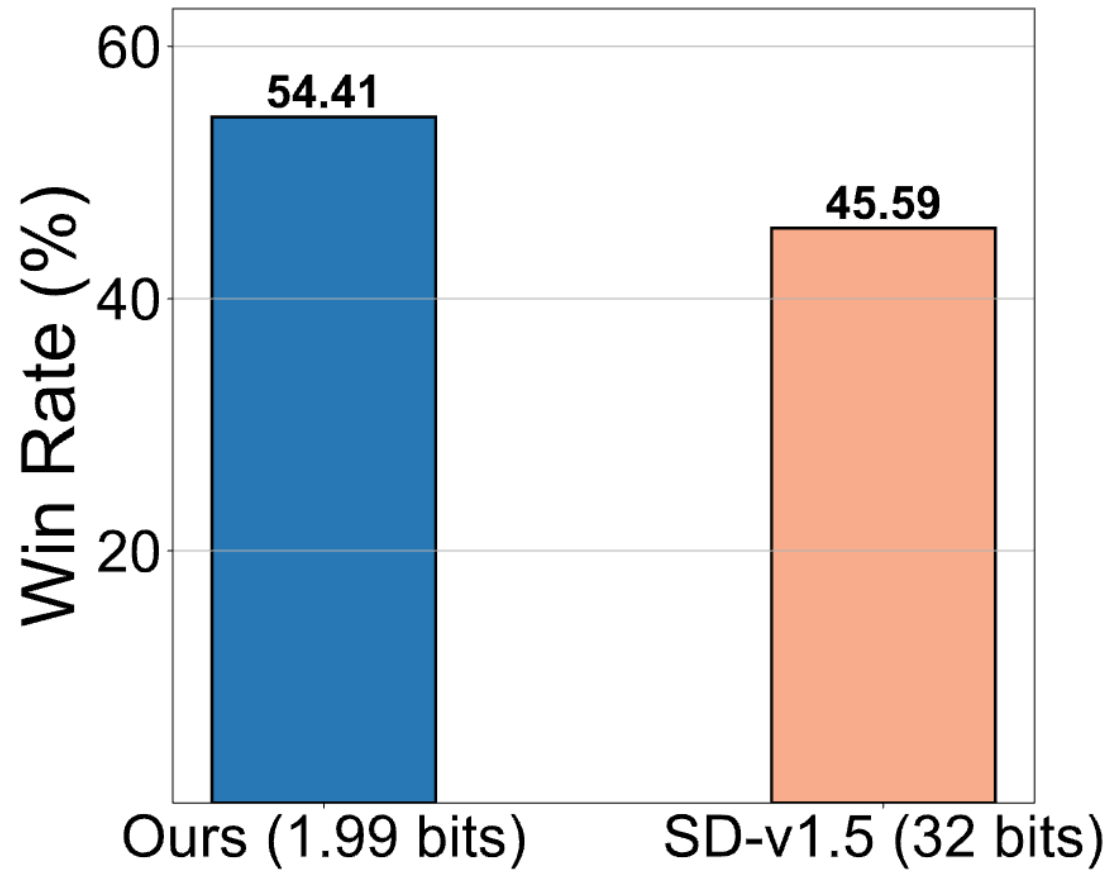CLIP Score on 30K MS-COCO.　　　　TIFA Scores　　　　GenEval Scores

Ours-I: Stage-I training

Ours-II: Stage-II training

BitsFusion consistently outperforms Stable Diffusion v1.5

# Results

## Human Evaluation

*Given a prompt, which image has better aesthetics and image-text alignment?*
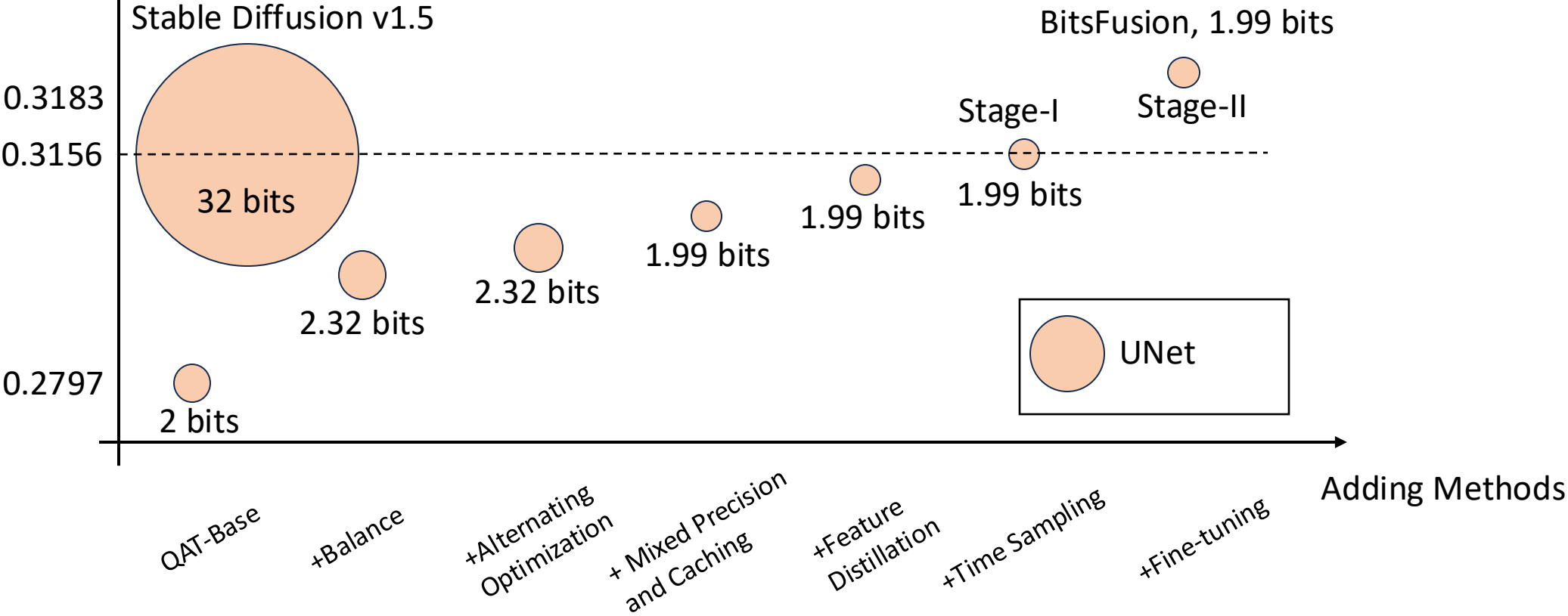


User preference of generated images from PartiPrompts (P2)

# Results

## Effect of each method



Average CLIP score across CFG scales 3.5, 5.5, 7.5, 9.5 on 1K PartiPrompts

Average CLIP score

Stable Diffusion v1.5

BitsFusion, 1.99 bits

Stage-I

Stage-II

0.3183

0.3156

32 bits

1.99 bits

1.99 bits

1.99 bits

2.32 bits

1.99 bits

2.32 bits

0.2797

2 bits

UNet

Adding Methods

QAT-Base

+Balance

+Alternating Optimization

+ Mixed Precision and Caching

+Feature Distillation

+Time Sampling

+Fine-tuning

# Results

## More comparisons

Stable Diffusion v1.5, 32 bits



BitsFusion, 1.99 bits

# Results

## More comparisons

Stable Diffusion v1.5, 32 bits



BitsFusion, 1.99 bits

# Results
## More comparisons

Stable Diffusion v1.5, 32 bits



BitsFusion, 1.99 bits

# Results

## More comparisons

Stable Diffusion v1.5, 32 bits



BitsFusion, 1.99 bits

# Results

## More comparisons

Stable Diffusion v1.5, 32 bits



BitsFusion, 1.99 bits

# Results

## More comparisons

Stable Diffusion v1.5, 32 bits



BitsFusion, 1.99 bits

# Results

## More comparisons

Stable Diffusion v1.5, 32 bits



BitsFusion, 1.99 bits

# Thank you