

DCC 2025

# COMPRESSION IN THE AGE OF GENAI

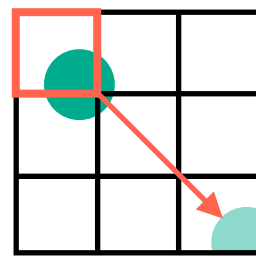
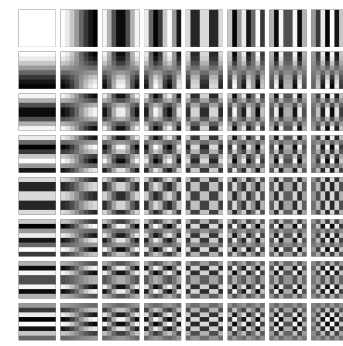
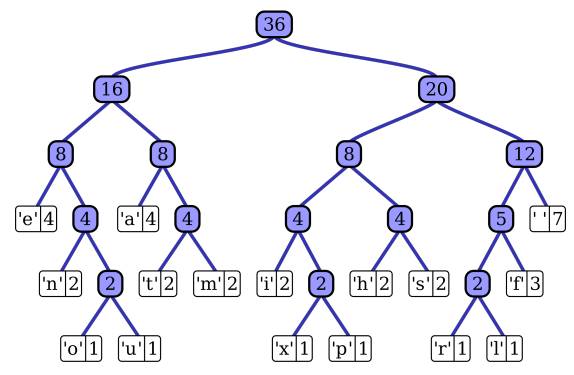
PERFECT REALISM AT EVERY BITRATE

Lucas Theis, Mabyduck

HFD: 0.0295 bpp

JPEG: 0.1102 bpp





Shannon  
1948

Lloyd  
1957

DCT  
1972

H.261  
1988

JPEG  
1992

AVC  
2004

HEVC  
2013

AV1  
2018

?

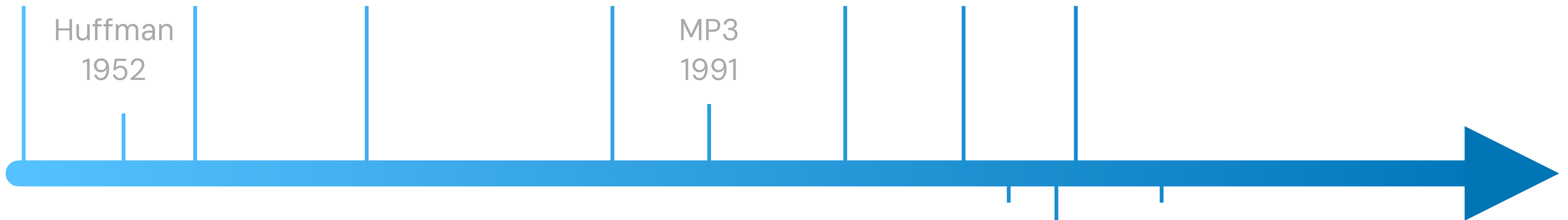
Huffman  
1952

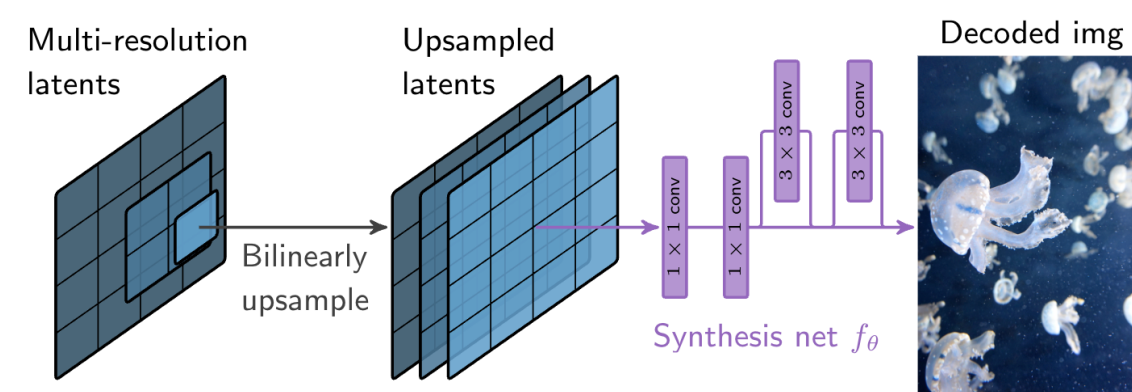
MP3  
1991

GANs  
2014

Diffusion  
2015

SD  
ChatGPT  
2022





Low-complexity  
neural compression

COOL-CHIC  
(Ladune et al., 2023)

C3  
(Kim et al., 2024;  
Ballé et al., 2024)

## Generative AI

Realism-distortion trade-offs

Reverse channel coding  
(a.k.a. channel simulation)

Diffusion compression

(Ho et al., 2020; Theis et al., 2022; Yang et al., 2025)

⚡ Rate-distortion theory

⚡ Quantization

⚡ Transform coding



Input



JPEG

0.068 bpp



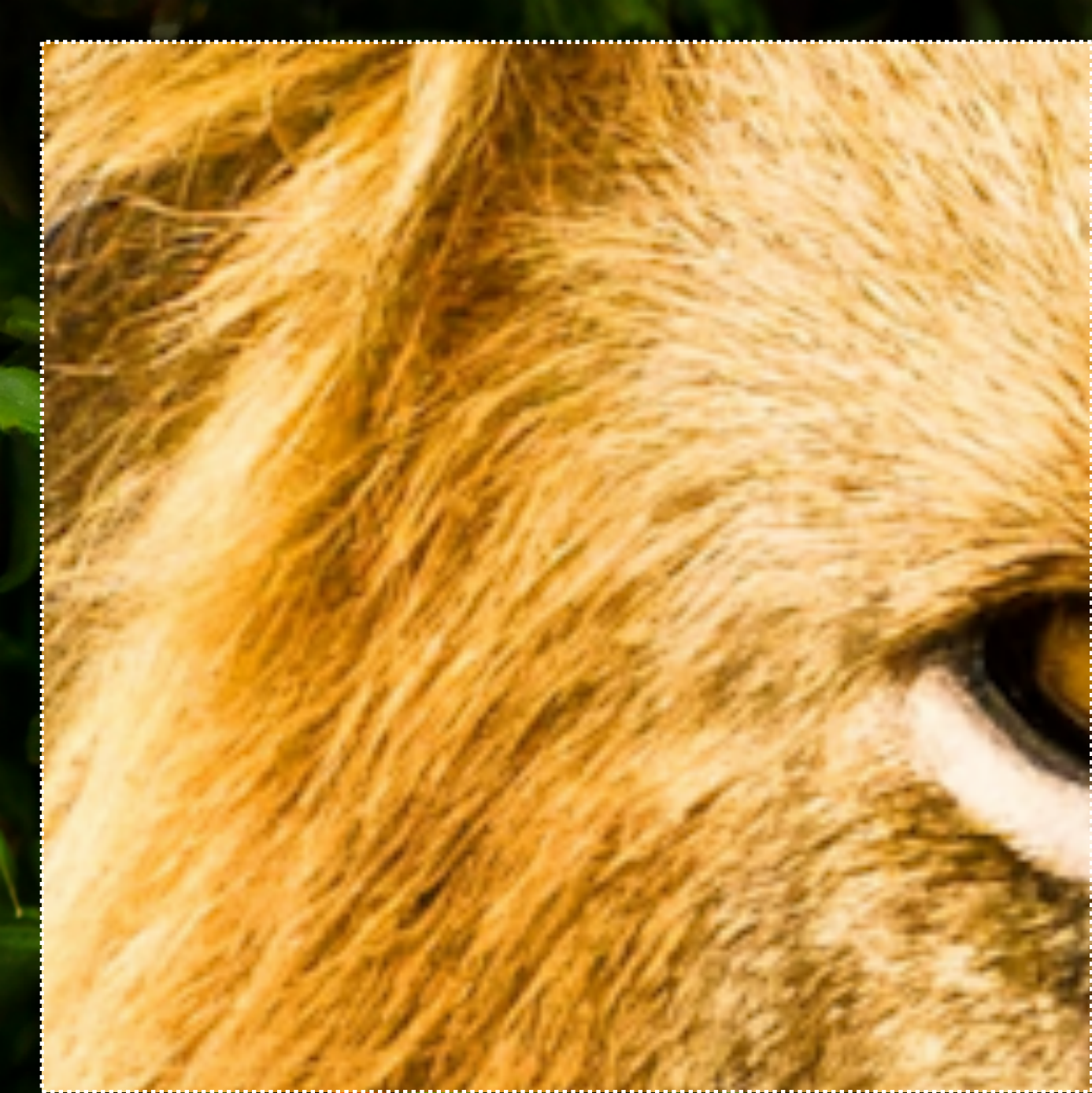
JPEG-XL

0.072 bpp



**ELIC** (MSE; based on He et al., 2022)

0.056 bpp



**HFD** (MSE + diffusion; Hoogeboom et al., 2023)

0.056 bpp





Input

# Overview

## Background

*What is the cost of perfect realism?*

## Diffusion I: High-fidelity diffusion (HFD)

*A transform coding approach*

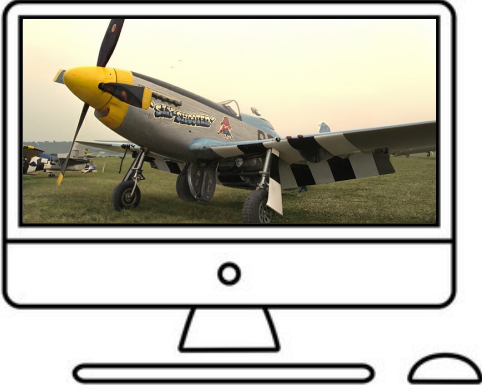
## Diffusion II: DiffC

*A novel approach using reverse channel coding*

# Realism

Accuracy of an ideal observer ( $\pm c$ )

Divergence

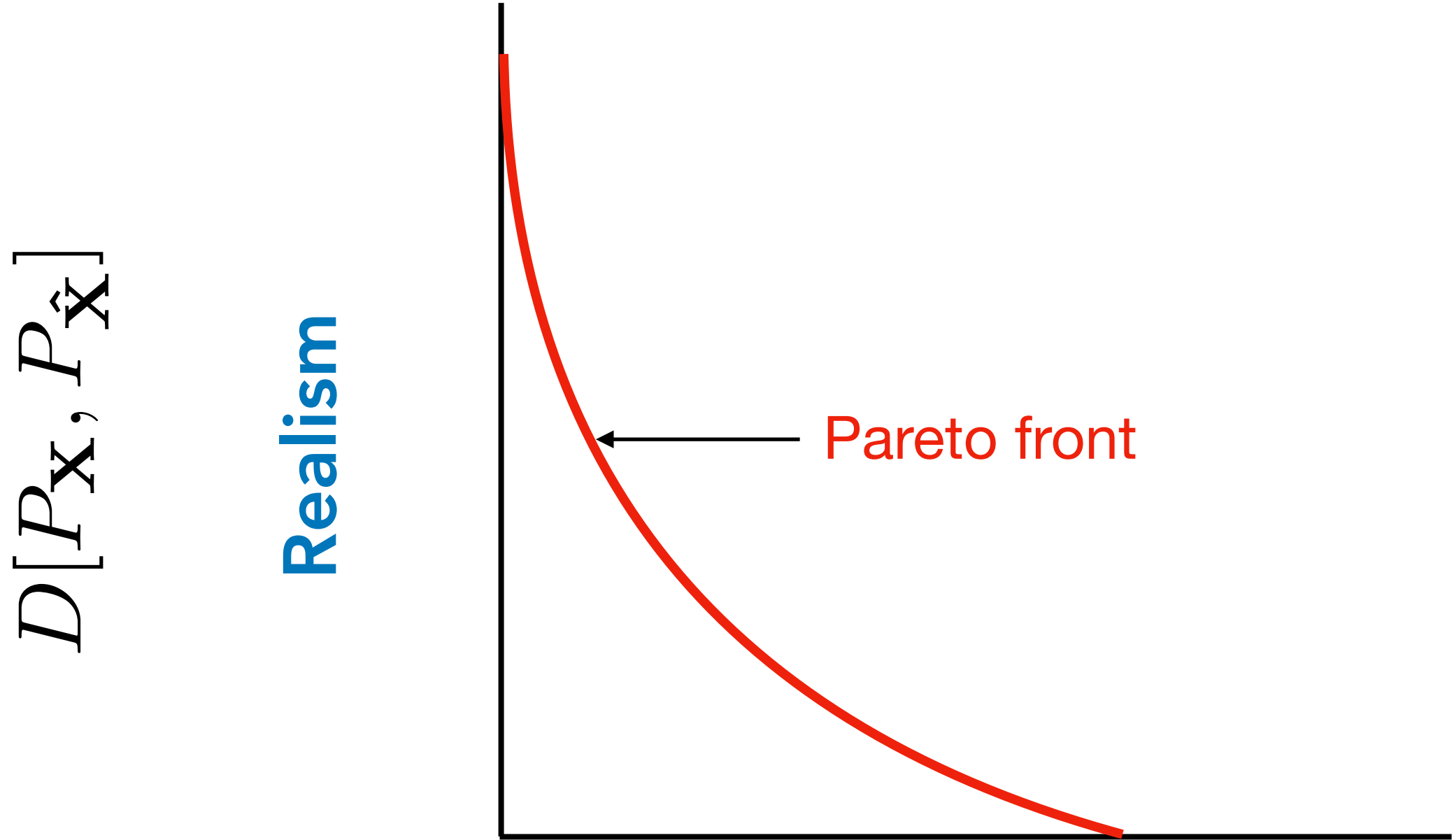


$$D[P_{\mathbf{x}}, P_{\hat{\mathbf{x}}}]$$

Image

Reconstructed image

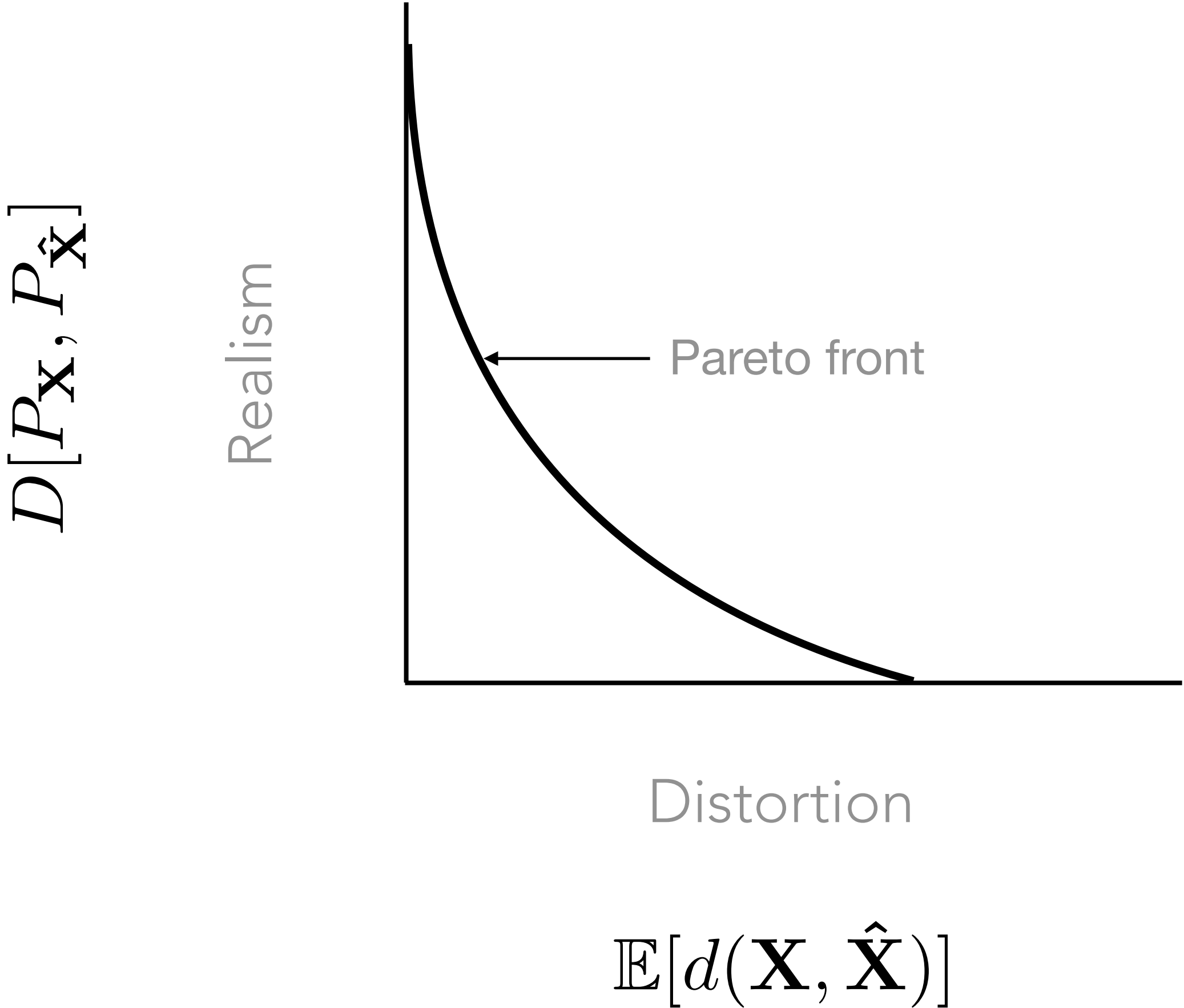
# Realism-distortion trade-off



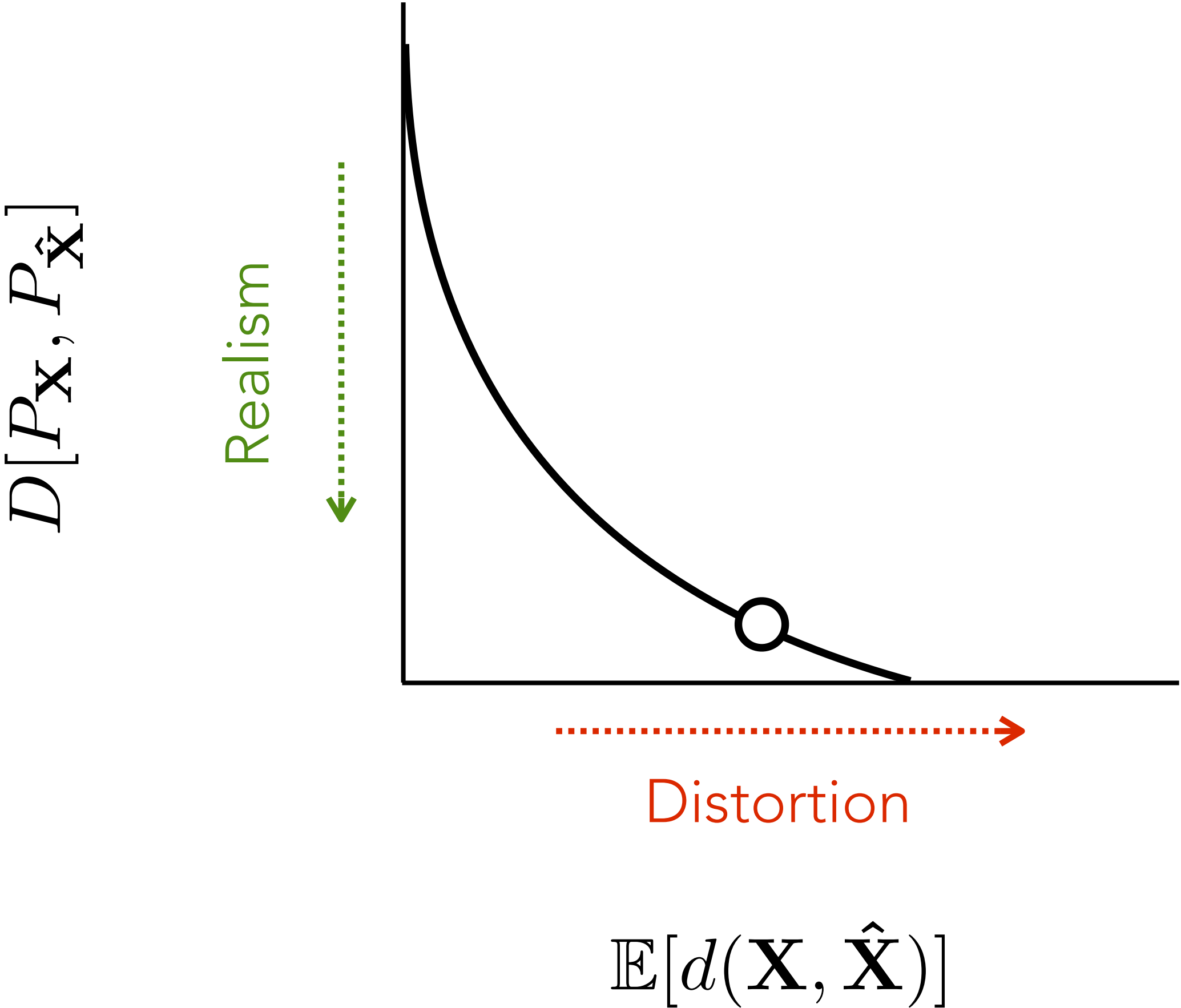
$$\mathbb{E}[d(\mathbf{X}, \hat{\mathbf{X}})]$$

Reconstruction

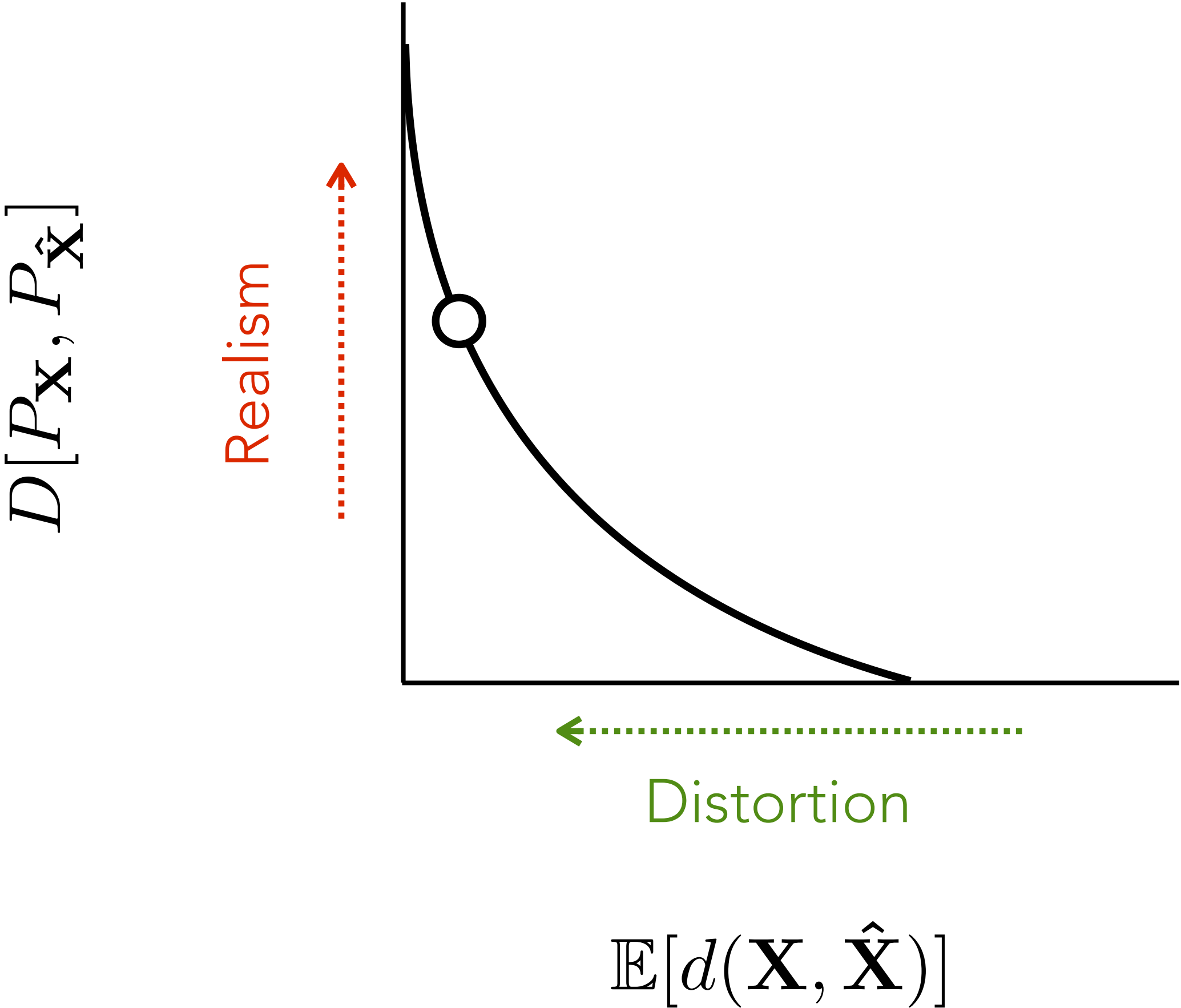
# Realism-distortion trade-off



# Realism-distortion trade-off



# Realism-distortion trade-off



# Perfect realism

$$D[P_{\mathbf{x}}, P_{\hat{\mathbf{x}}}] = 0$$

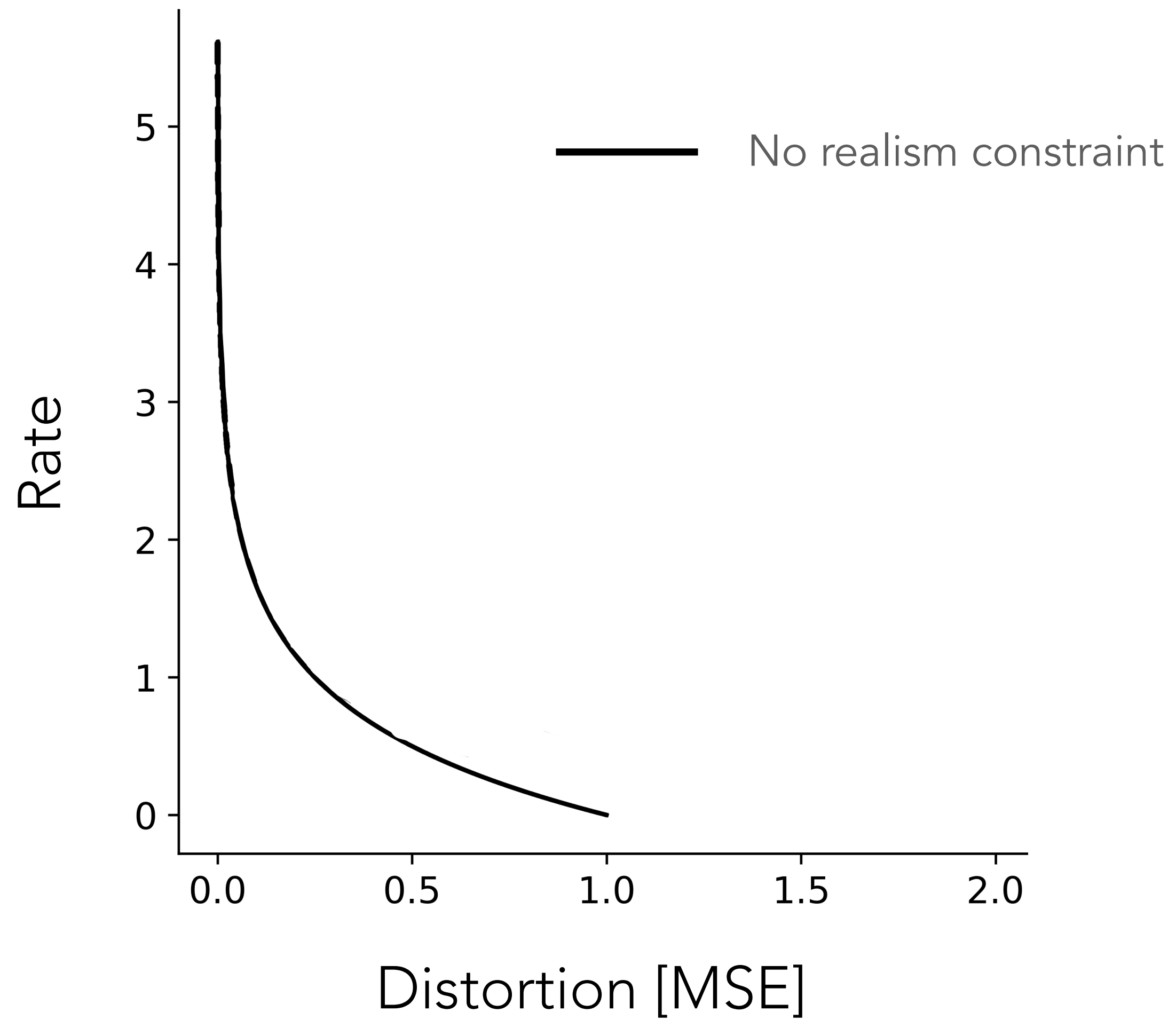
Image



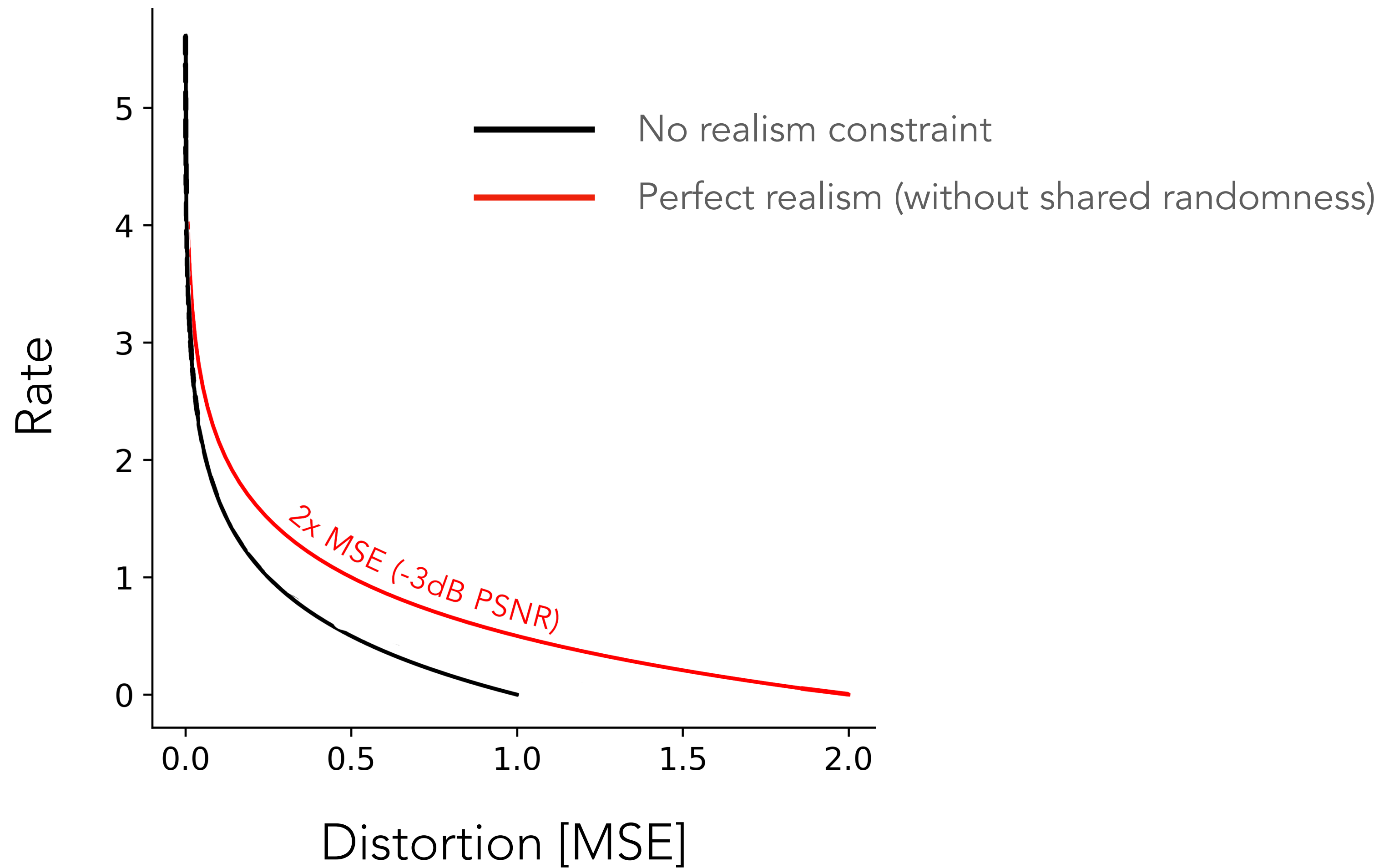
Reconstructed image



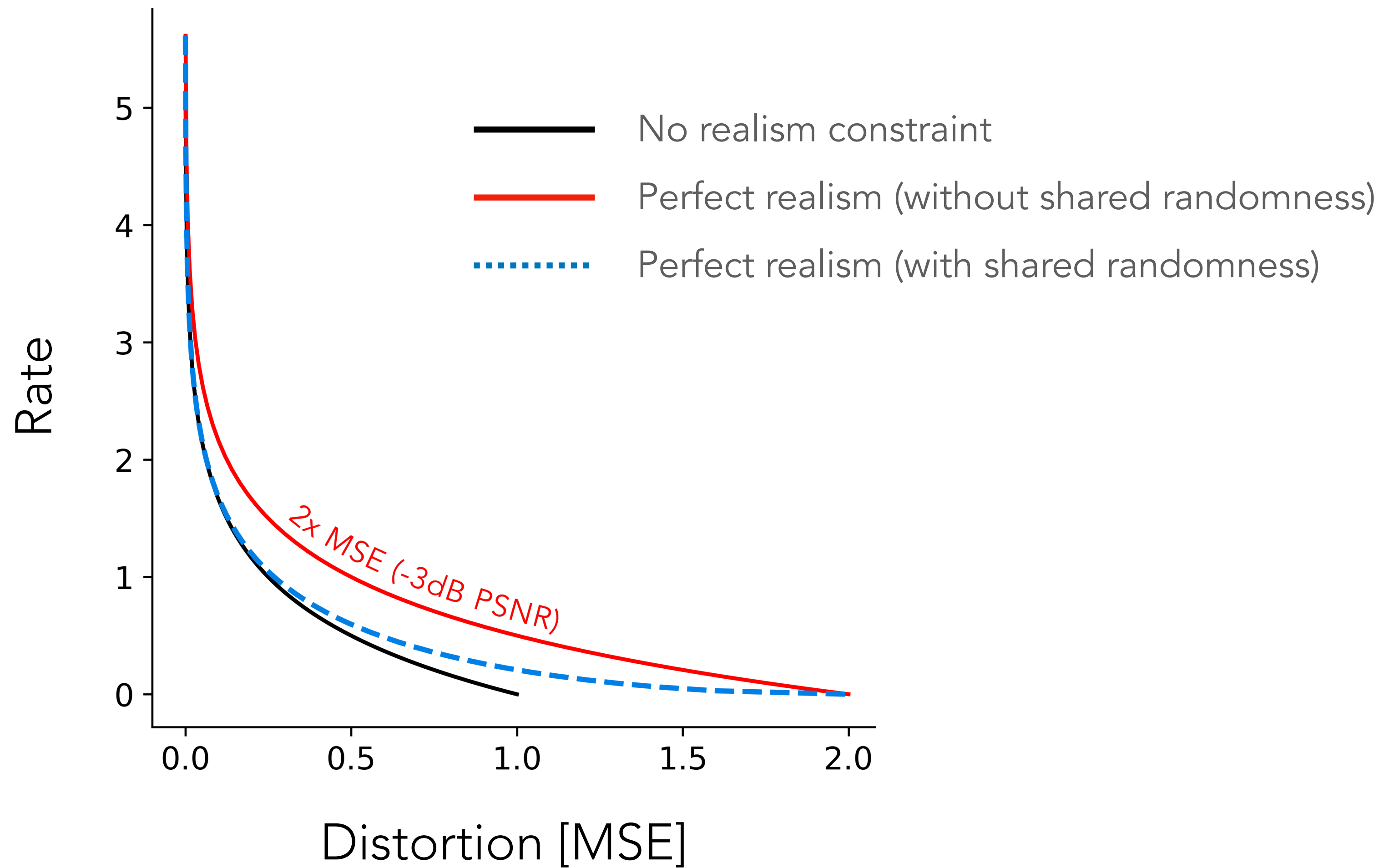
# The cost of perfect realism



# The cost of perfect realism



# The cost of perfect realism



# A mathematical theory of communication

systems; i.e., it must be possible to say of two systems represented by  $P_1(x,y)$  and  $P_2(x,y)$  that, according to our fidelity criterion, either (1) the first has higher fidelity, (2) the second has higher fidelity, or (3) they have equal fidelity. This means that a criterion of fidelity can be represented by a numerically valued function:

$$v(P(x,y))$$

whose argument ranges over possible probability functions  $P(x,y)$ .

We will now show that under very general and reasonable assumptions the function  $v(P(x,y))$  can be written in a seemingly much more specialized form, namely as an average of a function  $\rho(x,y)$  over the set of possible values of  $x$  and  $y$ :

DIFFUSION I:

# High-fidelity diffusion (HFD)

# Approach

1) Train neural compressor with MSE

2) Train generative model conditioned on output



$$\hat{\mathbf{X}}_{\text{MSE}}$$



$$\hat{\mathbf{X}} \sim P_{\mathbf{X}|\hat{\mathbf{X}}_{\text{MSE}}}$$



ELIC (0.1674 bpp)



HFD (0.1674 bpp)



# Justification

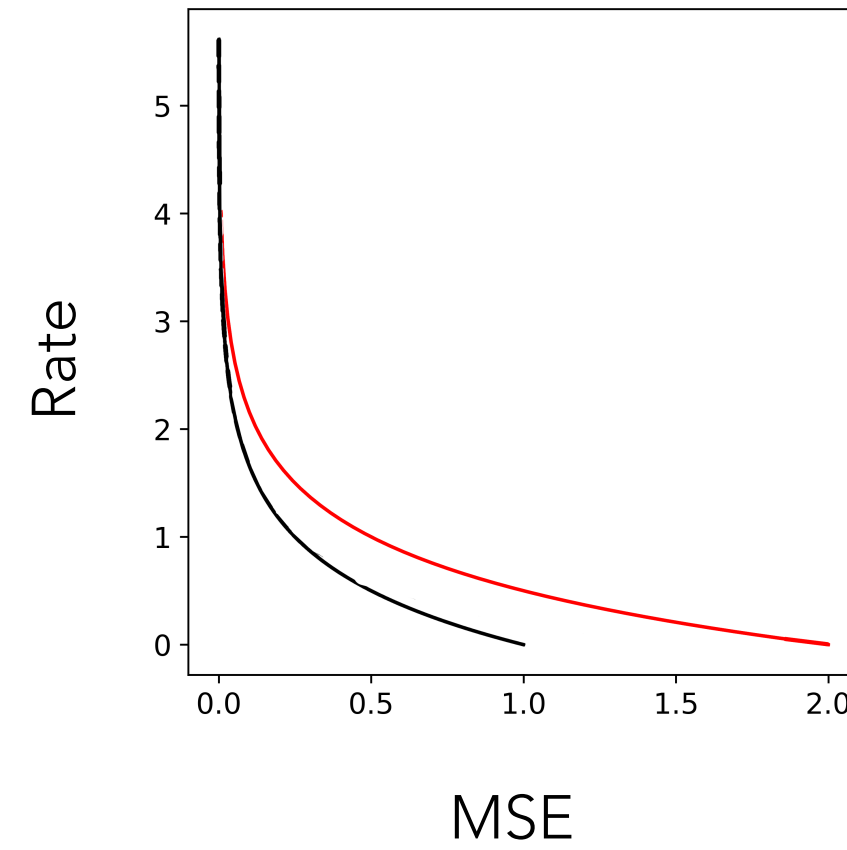
$$\mathbb{E}[\|\hat{\mathbf{X}} - \mathbf{X}\|^2]$$

Output of generative model  
(**realistic**)

$$\mathbb{E}[\|\hat{\mathbf{X}}^{\text{MSE}} - \mathbf{X}\|^2]$$

Output of neural compressor  
(**blurry**)

# Justification

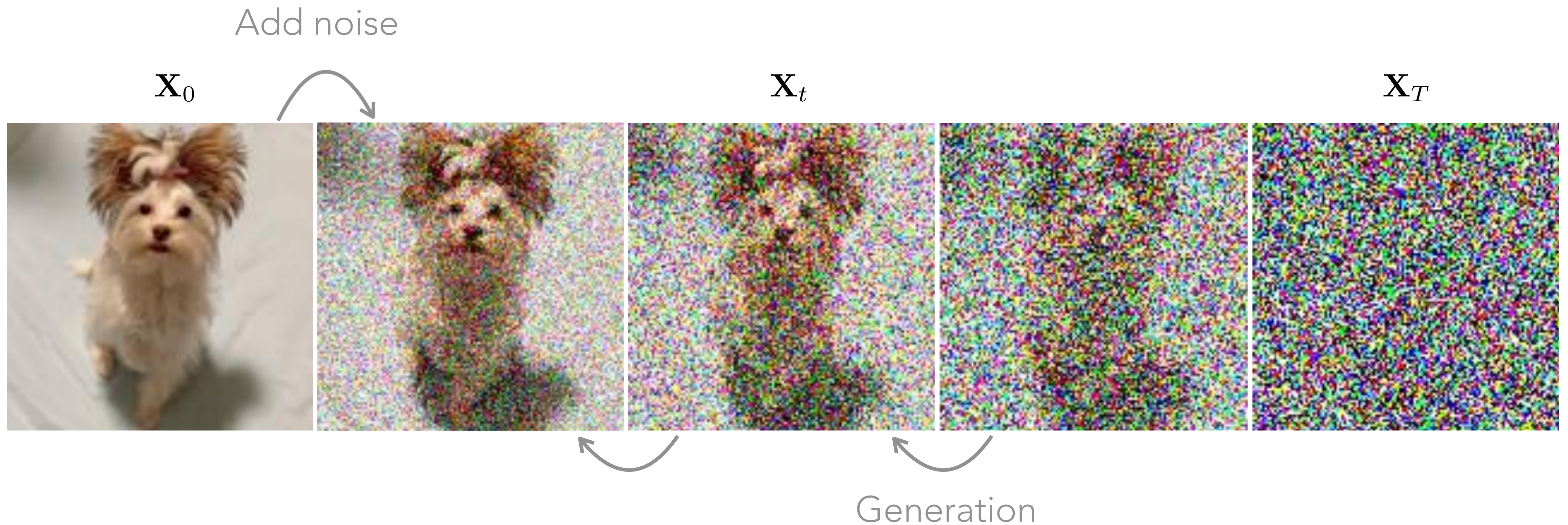


$$\mathbb{E}[\|\hat{\mathbf{X}} - \mathbf{X}\|^2] \leq 2 \mathbb{E}[\|\hat{\mathbf{X}}^{\text{MSE}} - \mathbf{X}\|^2]$$

Output of generative model  
(**realistic**)

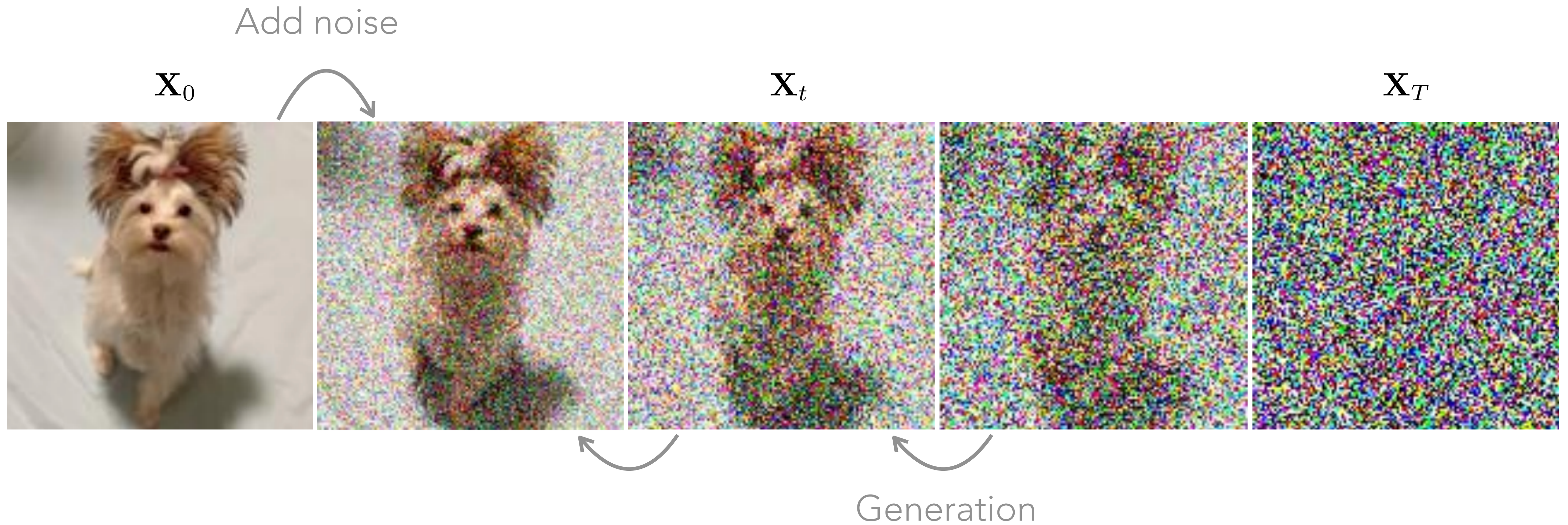
Output of neural compressor  
(**blurry**)

# Diffusion



💡  $P_{\mathbf{x}_{t-\delta}|\mathbf{x}_t}$  is approximately Gaussian (e.g., Feller, 1949; Anderson, 1982).

# Diffusion



💡 Optimize  $\mathbb{E}[\|\mathbf{X}_{t-\delta} - m_{\theta}(\mathbf{X}_t)\|^2]$  so that  $m_{\theta}(\mathbf{X}_t) \approx \mathbb{E}[\mathbf{X}_{t-\delta} \mid \mathbf{X}_t]$ .

# Overview

ELIC (He et al., 2022)

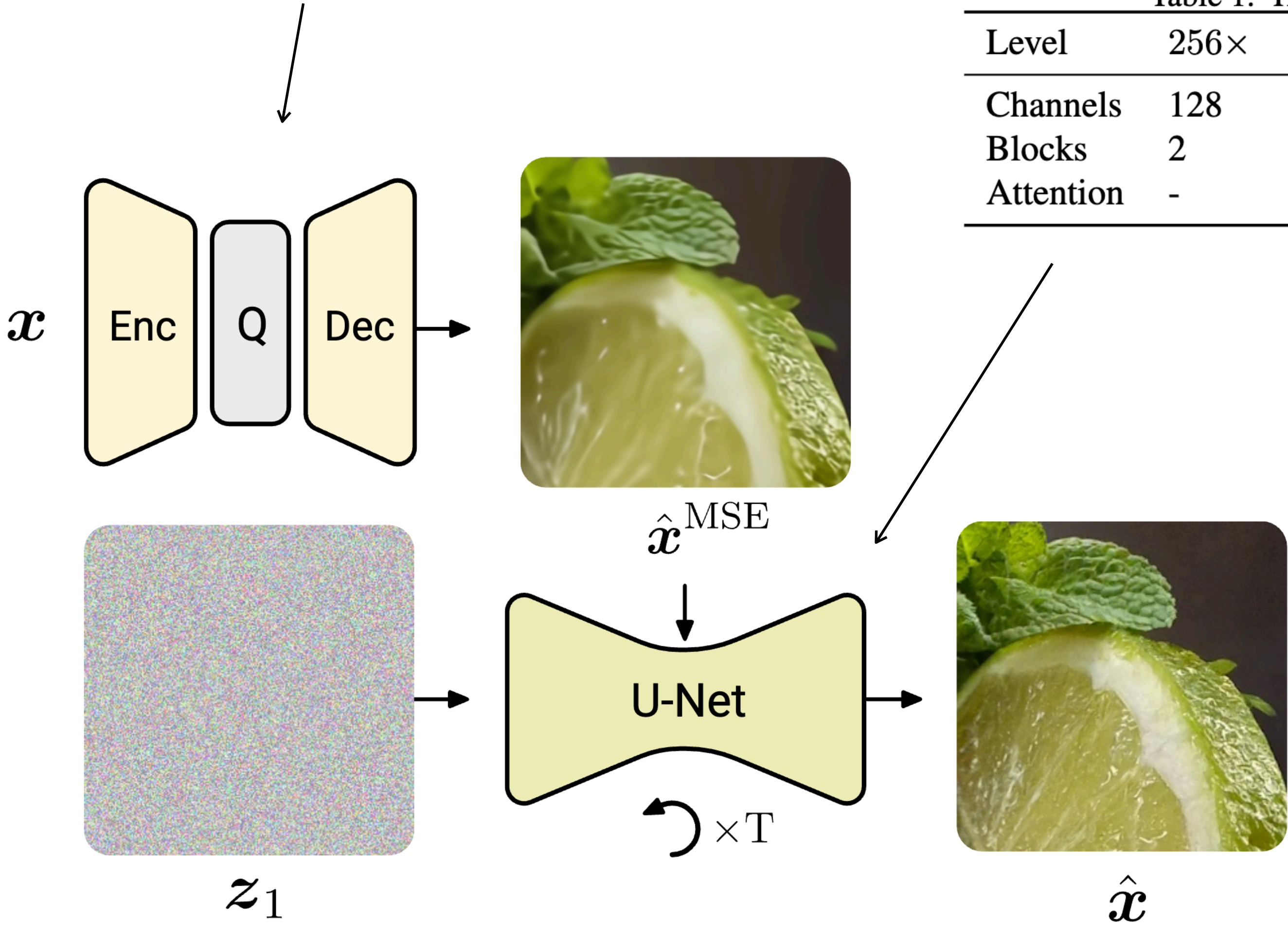
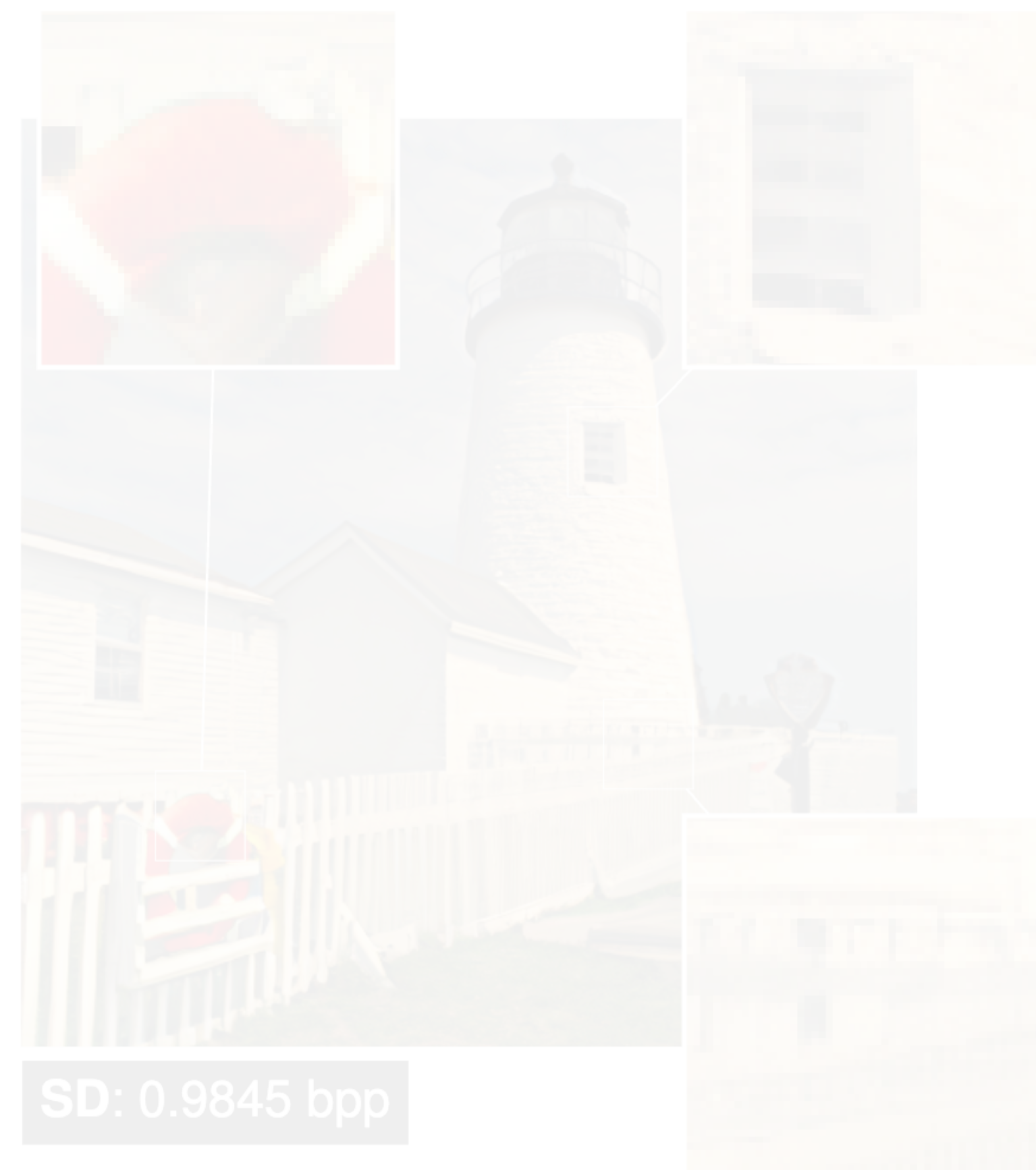
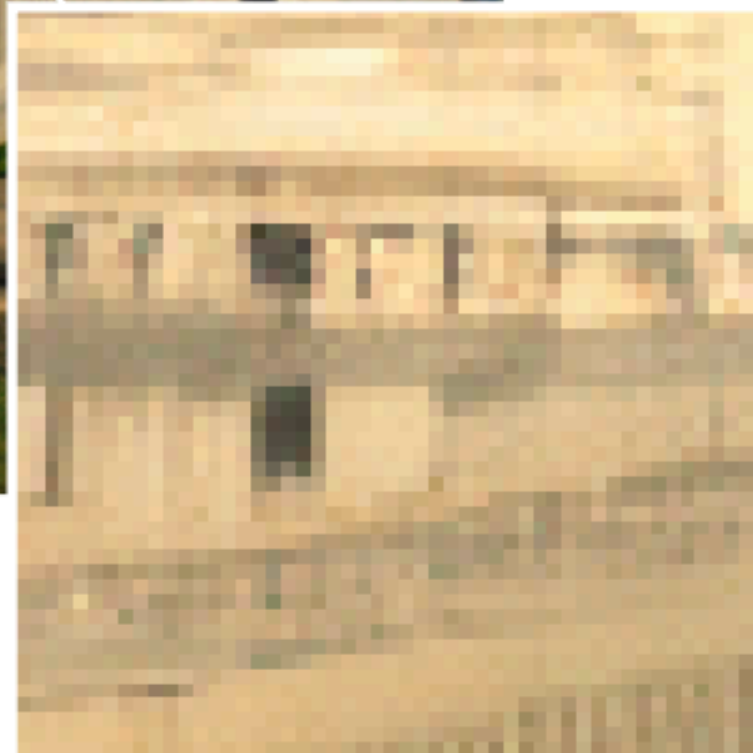
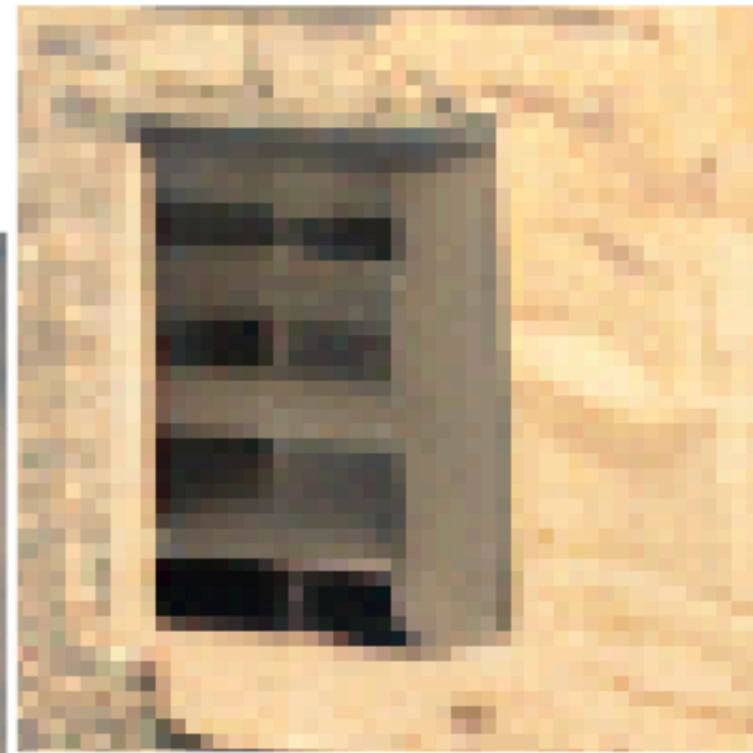
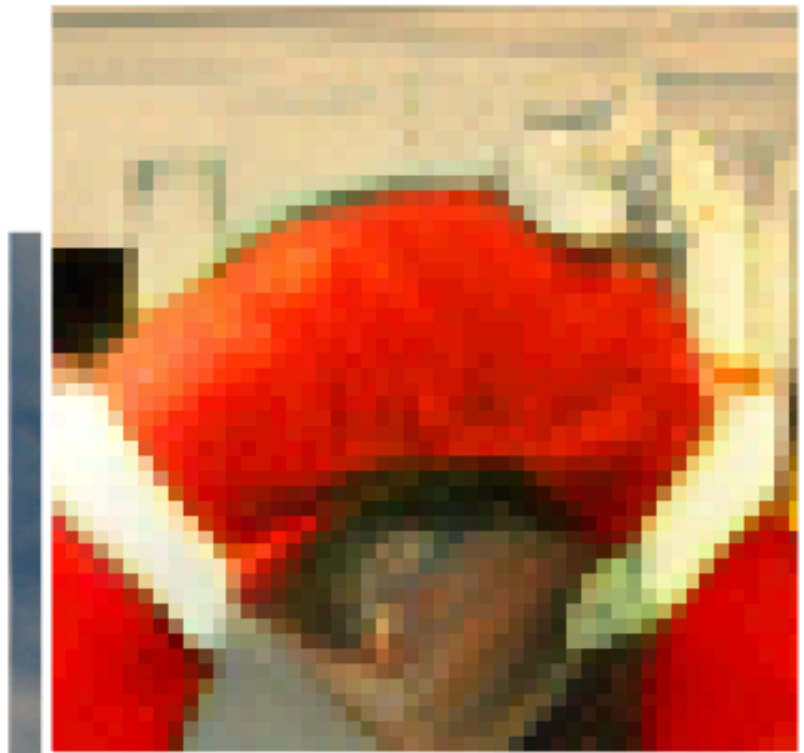


Table 1. HFD U-Net architecture

Level	256×	128×	64×	32×	16×
Channels	128	128	256	256	1024
Blocks	2	2	2	2	16
Attention	-	-	-	-	✓

# JPEG + StableDiffusion





**SD: 0.9845 bpp**



**HFD (Ours): 0.0918 bpp**

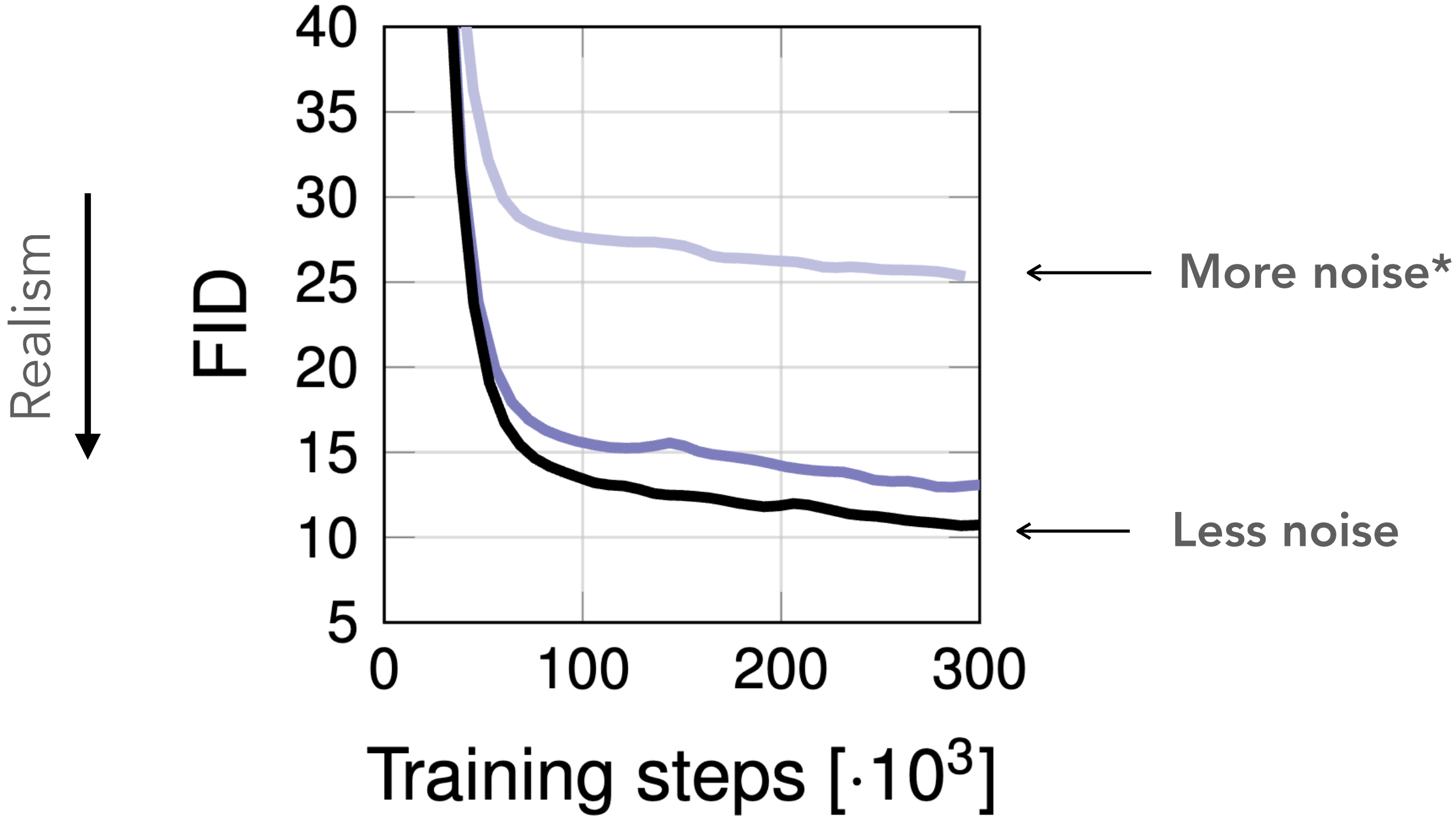


“A lighthouse in Maine behind a white fence with a red life buoy hanging on it.”



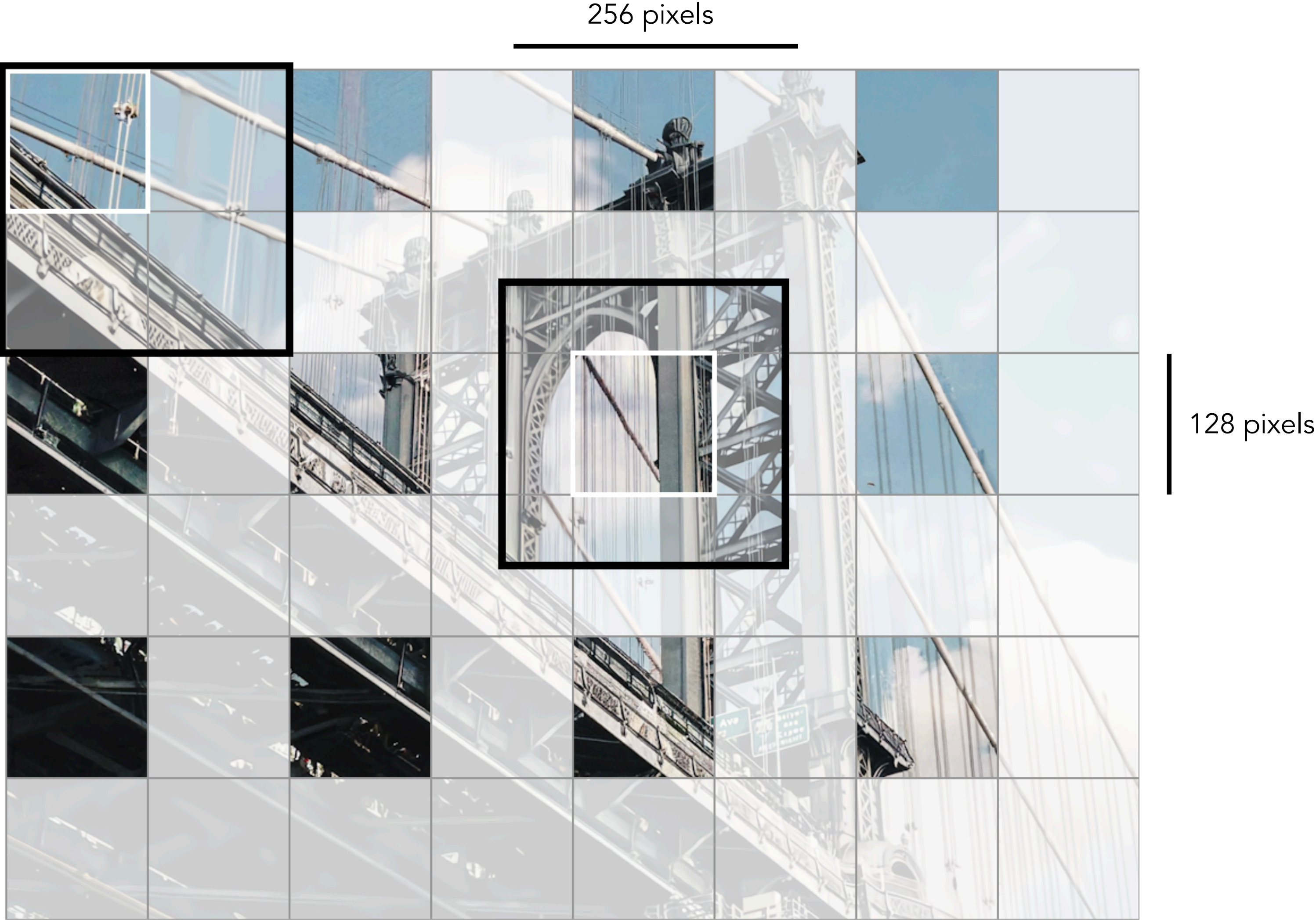
Sharper, but less consistent

# Noise schedule

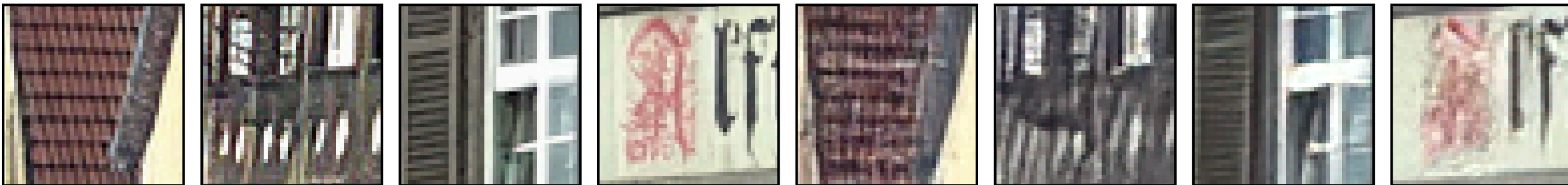


\* Recommended for text-to-image models (Chen et al., 2023; Hoogeboom et al., 2023a)

# Patch-wise generation





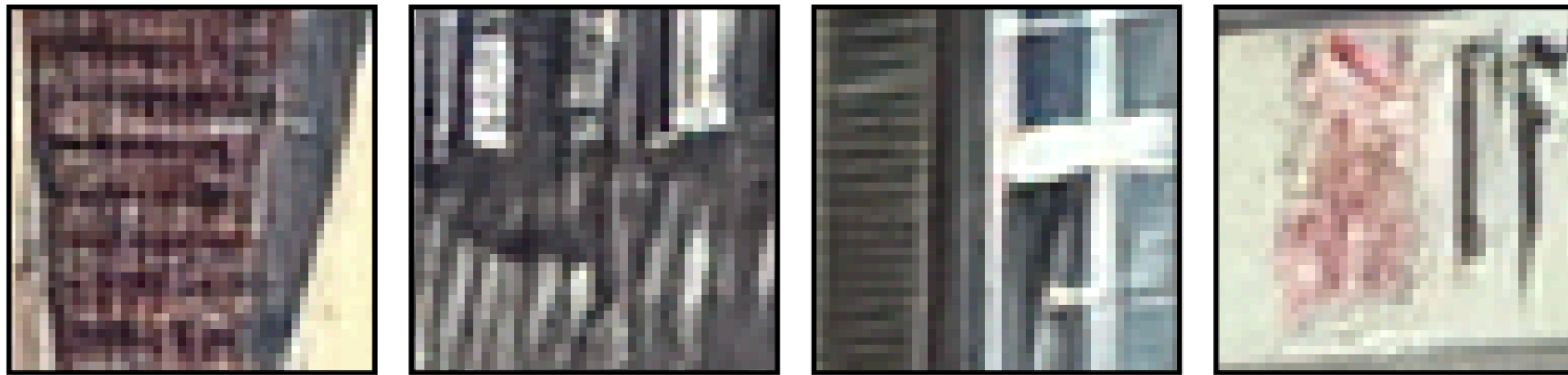


**HFD (Ours): 0.2639 bpp**

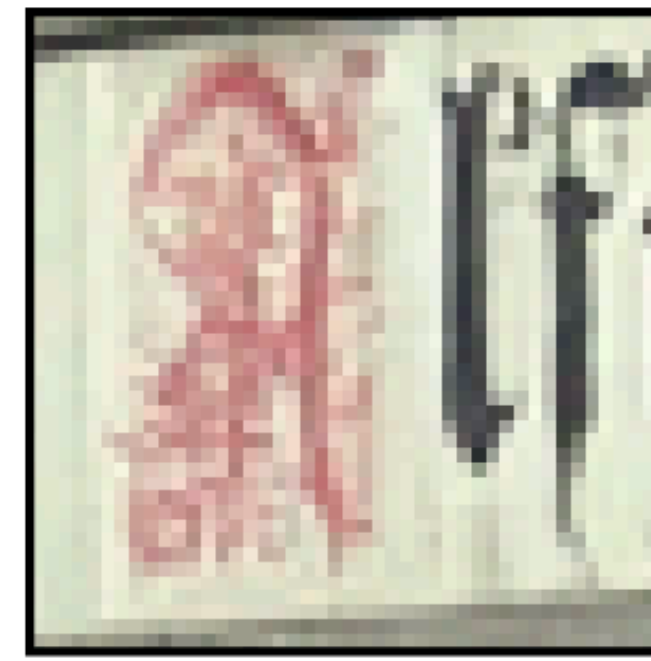
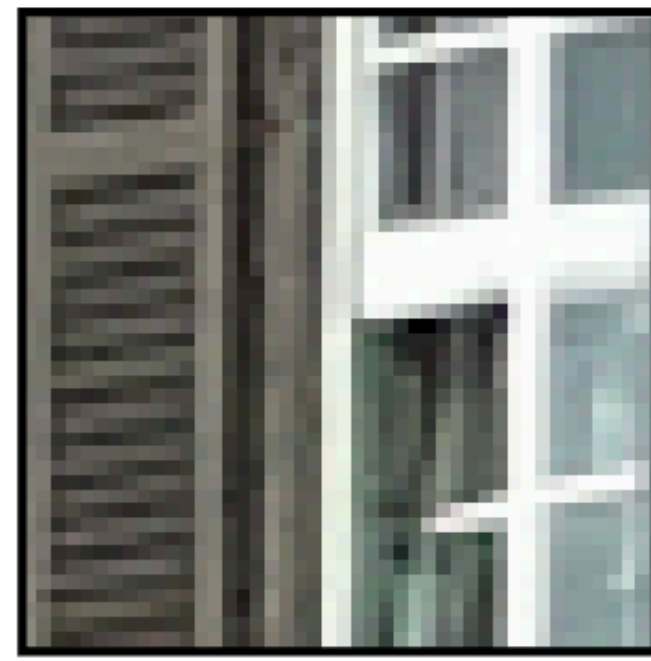
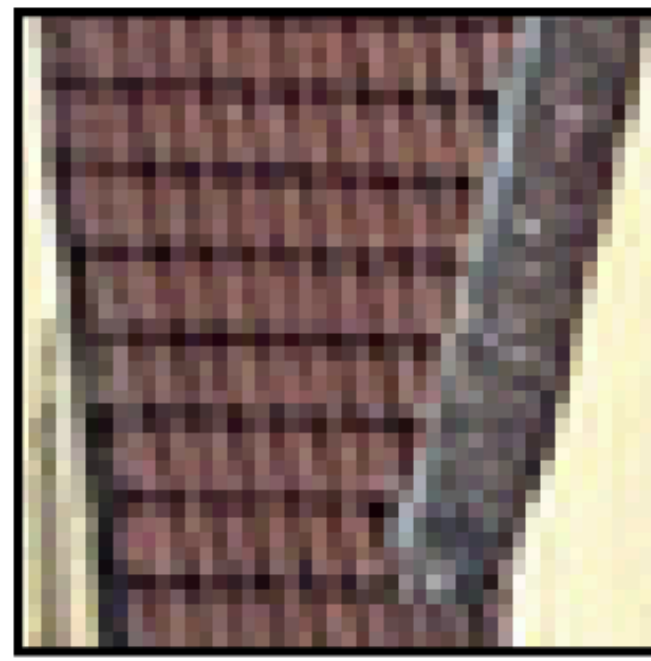


**Yang & Mandt (2023): 0.2971 bpp**

End-to-end trained



Yang & Mandt (2023): 0.2971 bpp



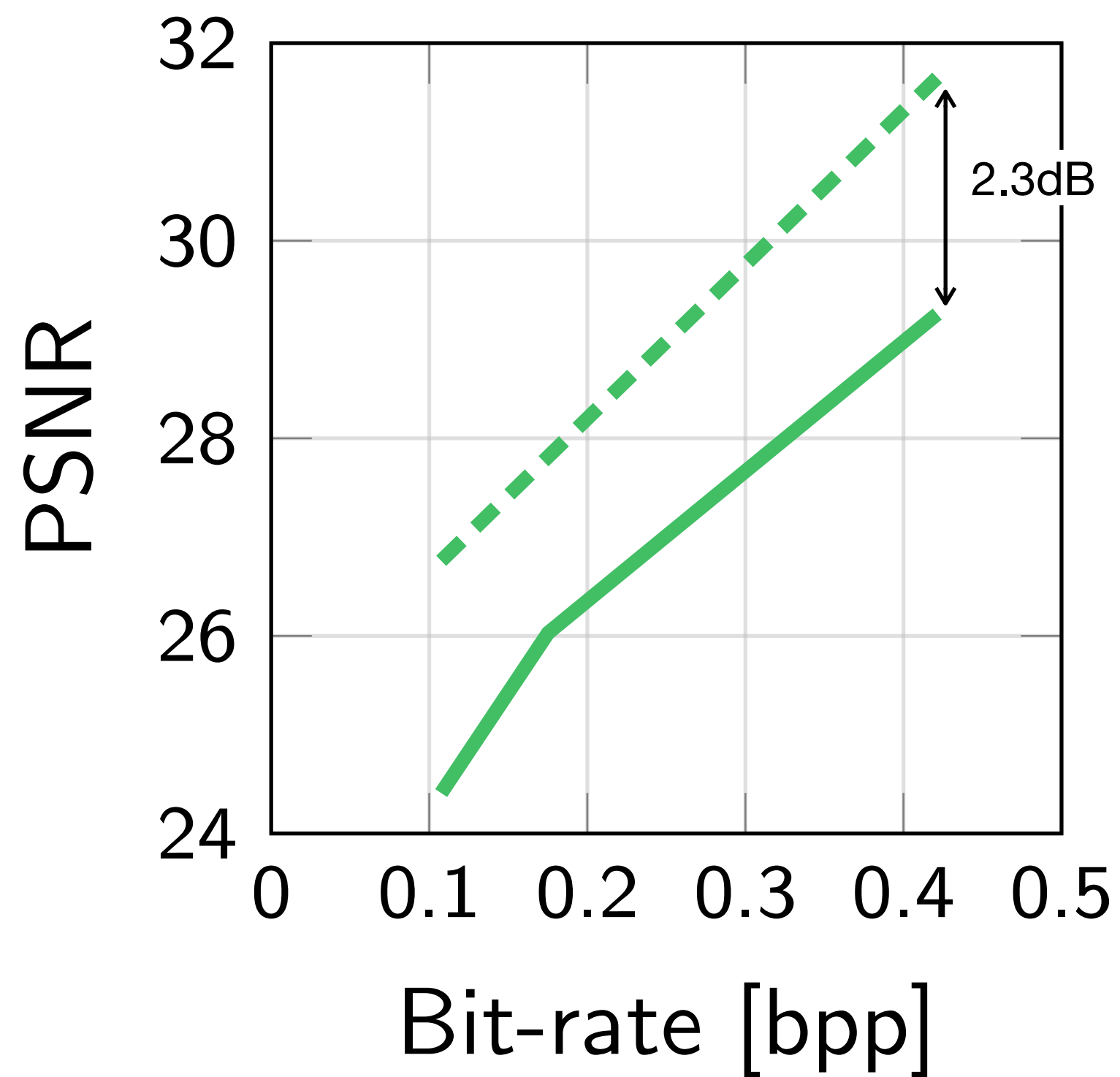
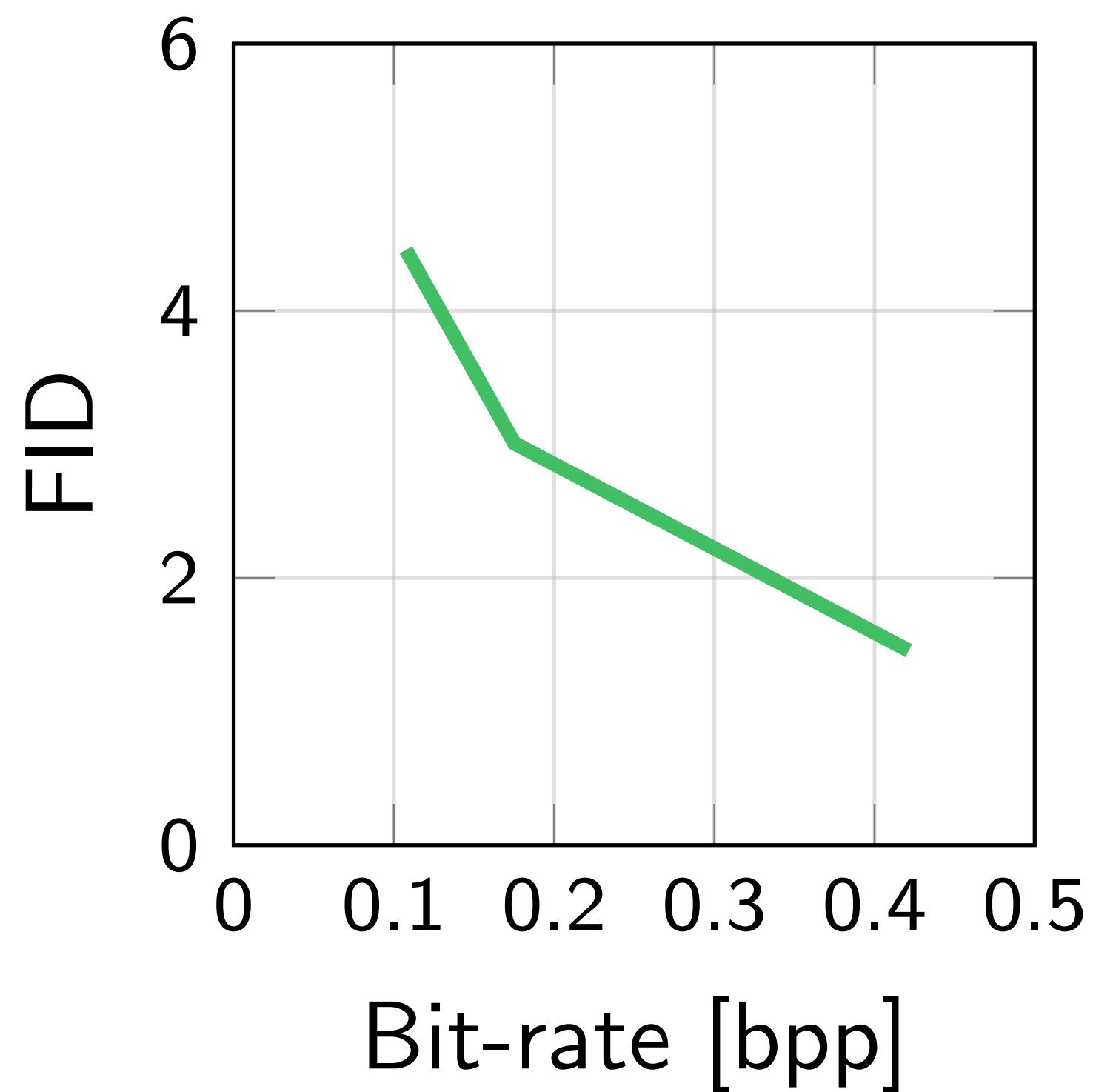
**HFD (Ours): 0.2639 bpp**

Realism

Distortion

MS-COCO 30k

MS-COCO 30k



- HFD/DDPM (Ours)
- - - ELIC

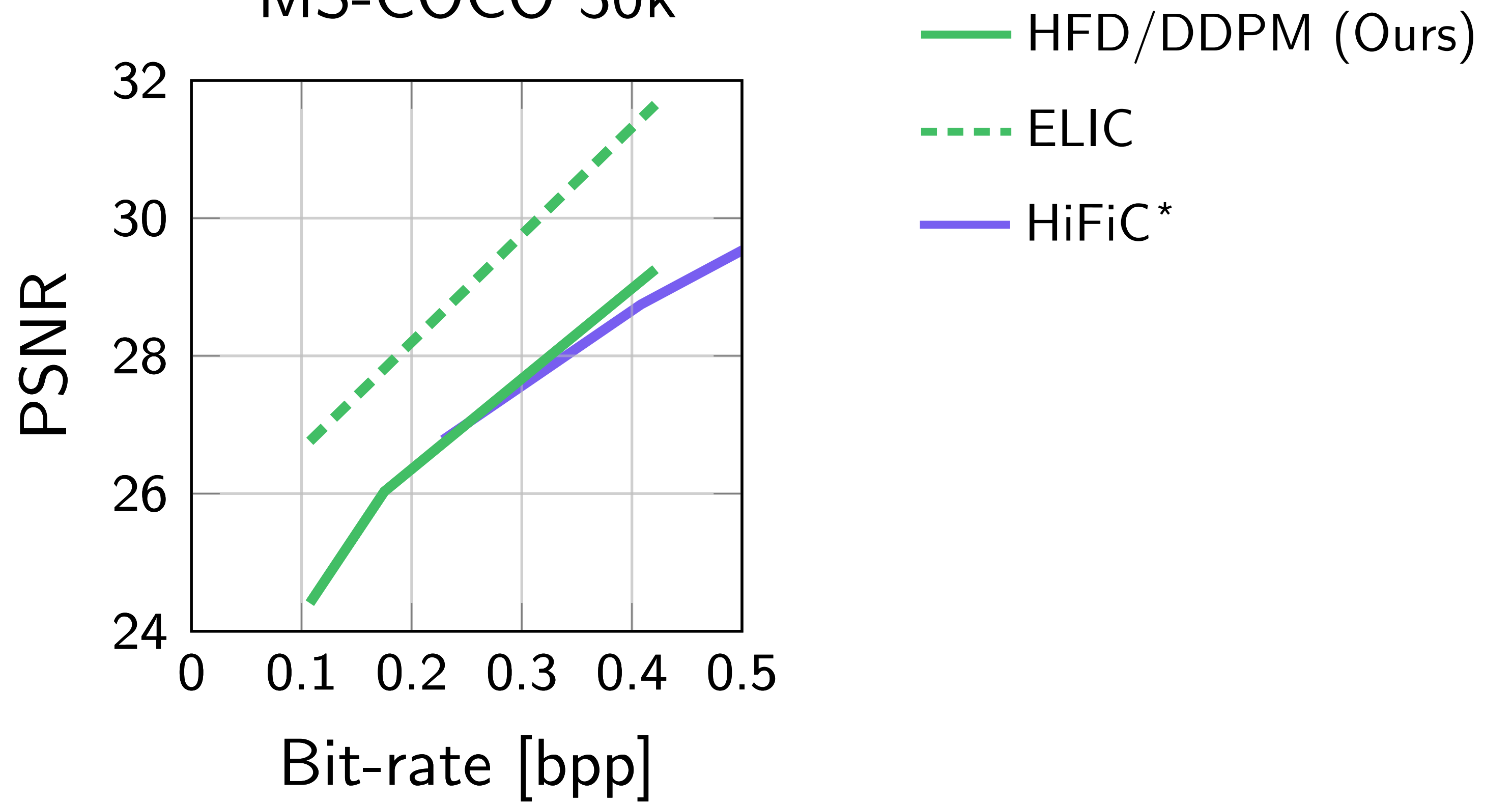
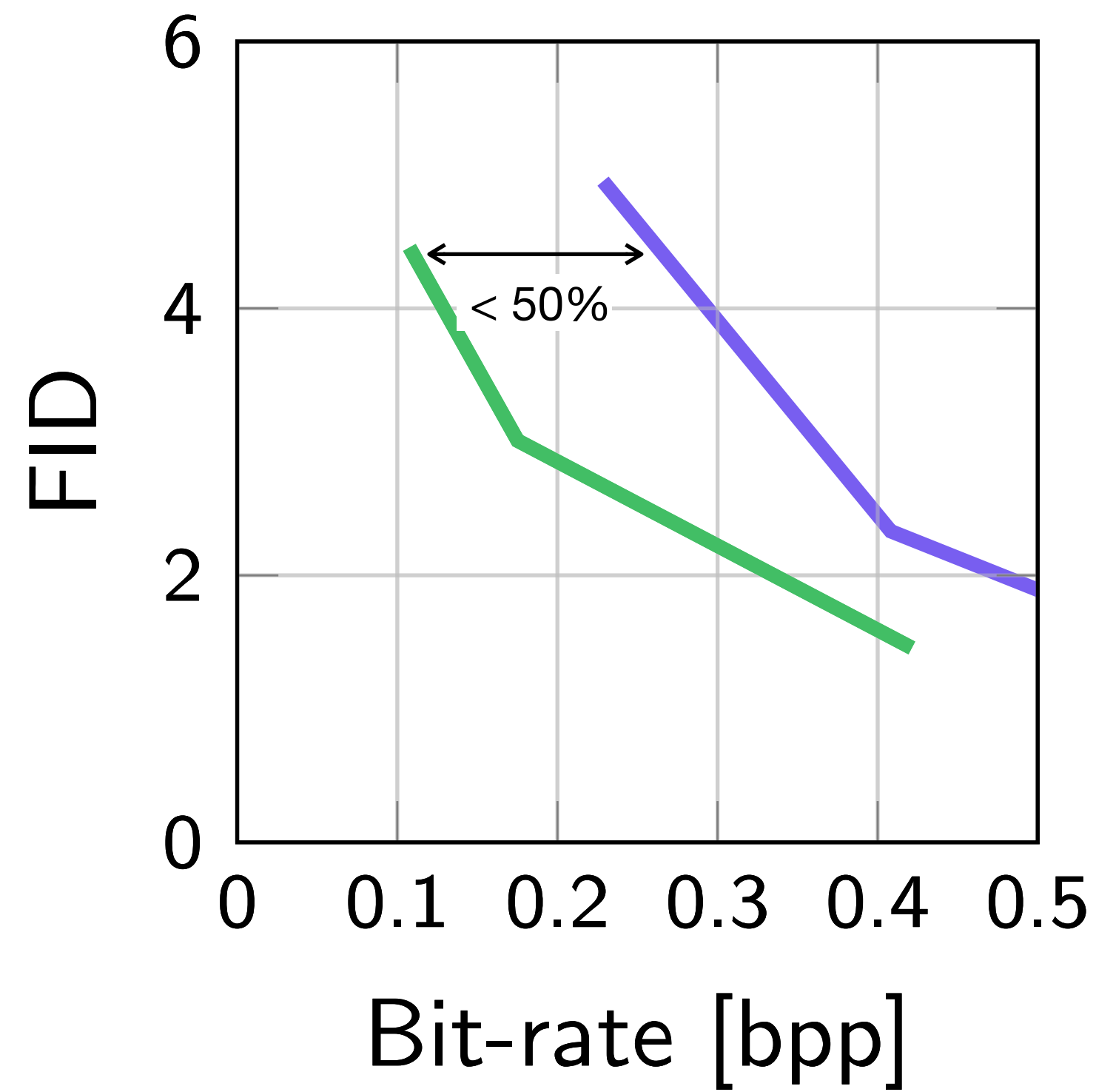


Realism

Distortion

MS-COCO 30k

MS-COCO 30k

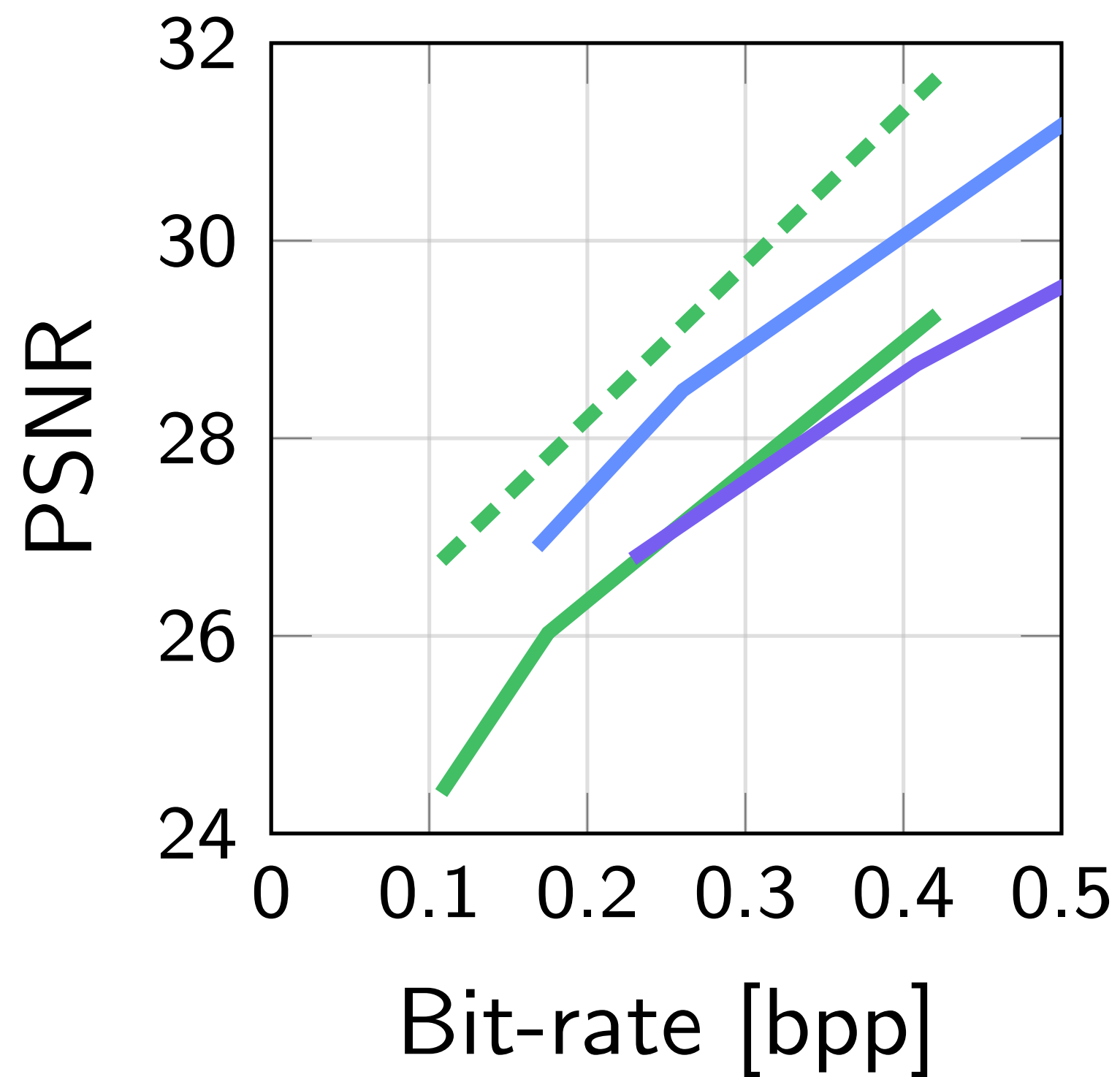
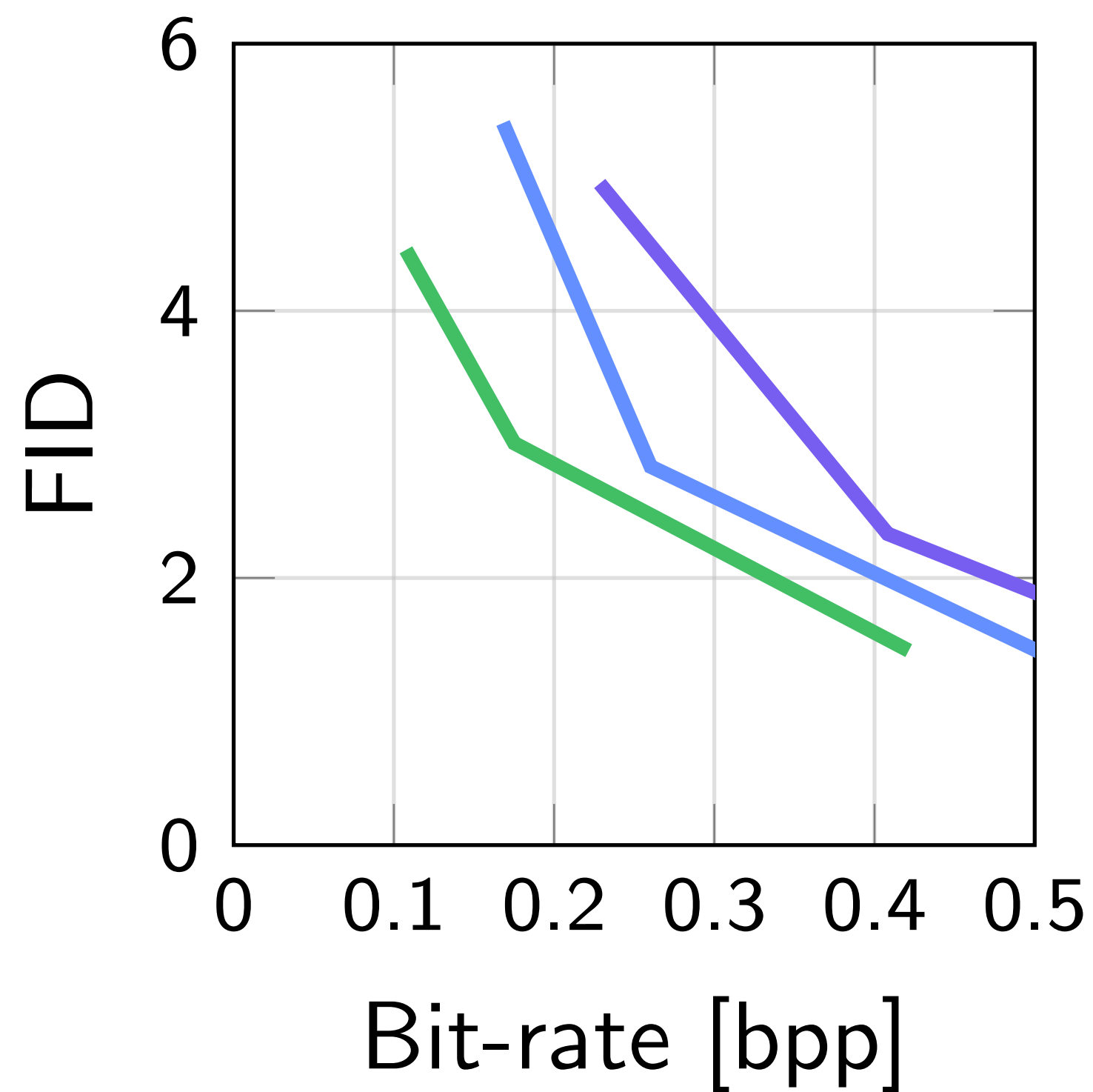


Realism

Distortion

MS-COCO 30k

MS-COCO 30k

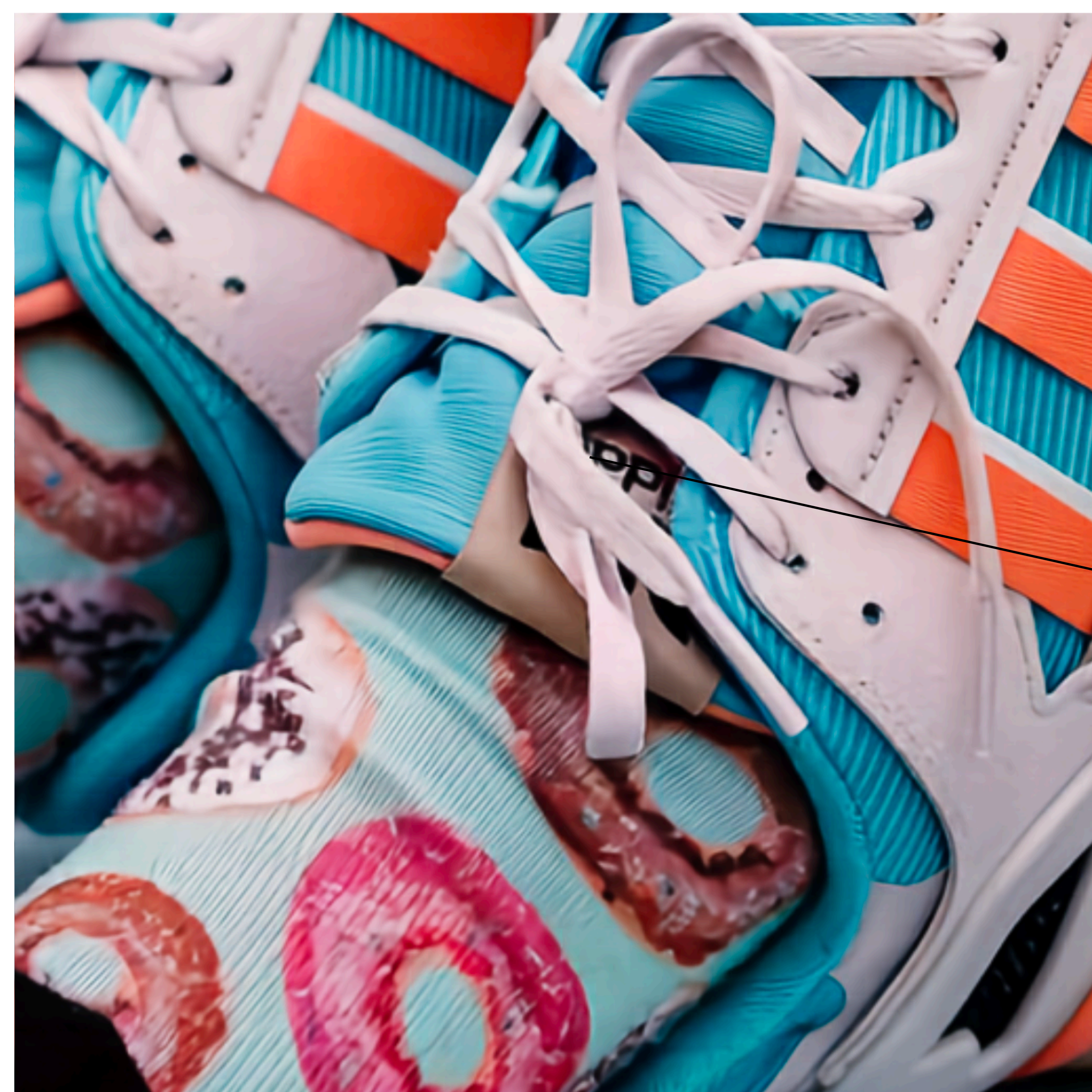


- HFD/DDPM (Ours)
- - - ELIC
- HiFiC
- MR\*



**HFD (Ours): 0.0538 (44.2%)**

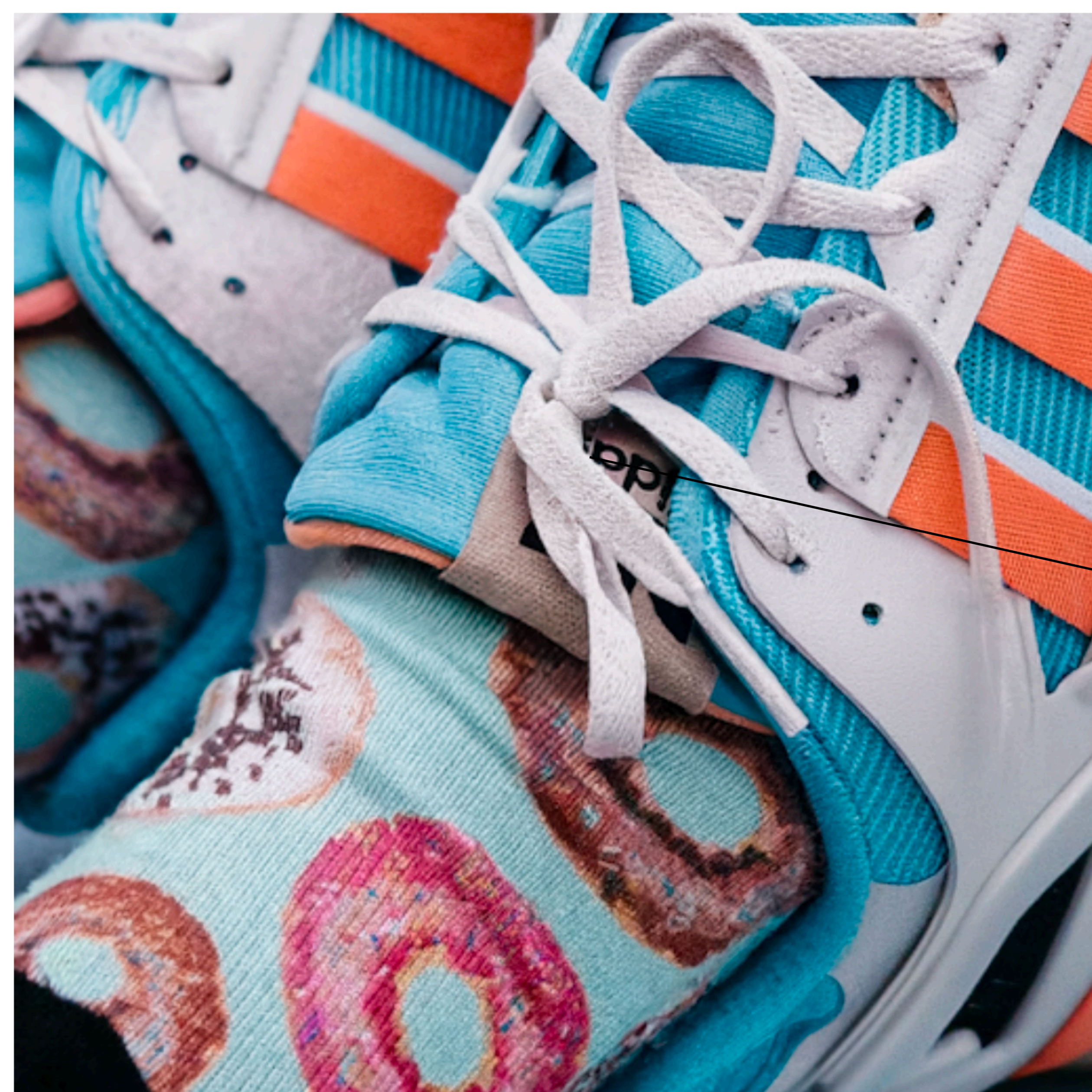
**PQ-MIM [8]: 0.124 (100%)**



**PQ-MIM [8]: 0.124 (100%)**

(similar FID)

(El-Nouby et al., 2023)



**HFD (Ours): 0.0538 (44.2%)**

(similar FID)

(Hoogeboom et al., 2023)

DIFFUSION II:

**DiffC**

# DiffC

## HFD (transform coding)

- Analysis transform
- Synthesis transform
- Conditional diffusion model
- Entropy model
- Quantization

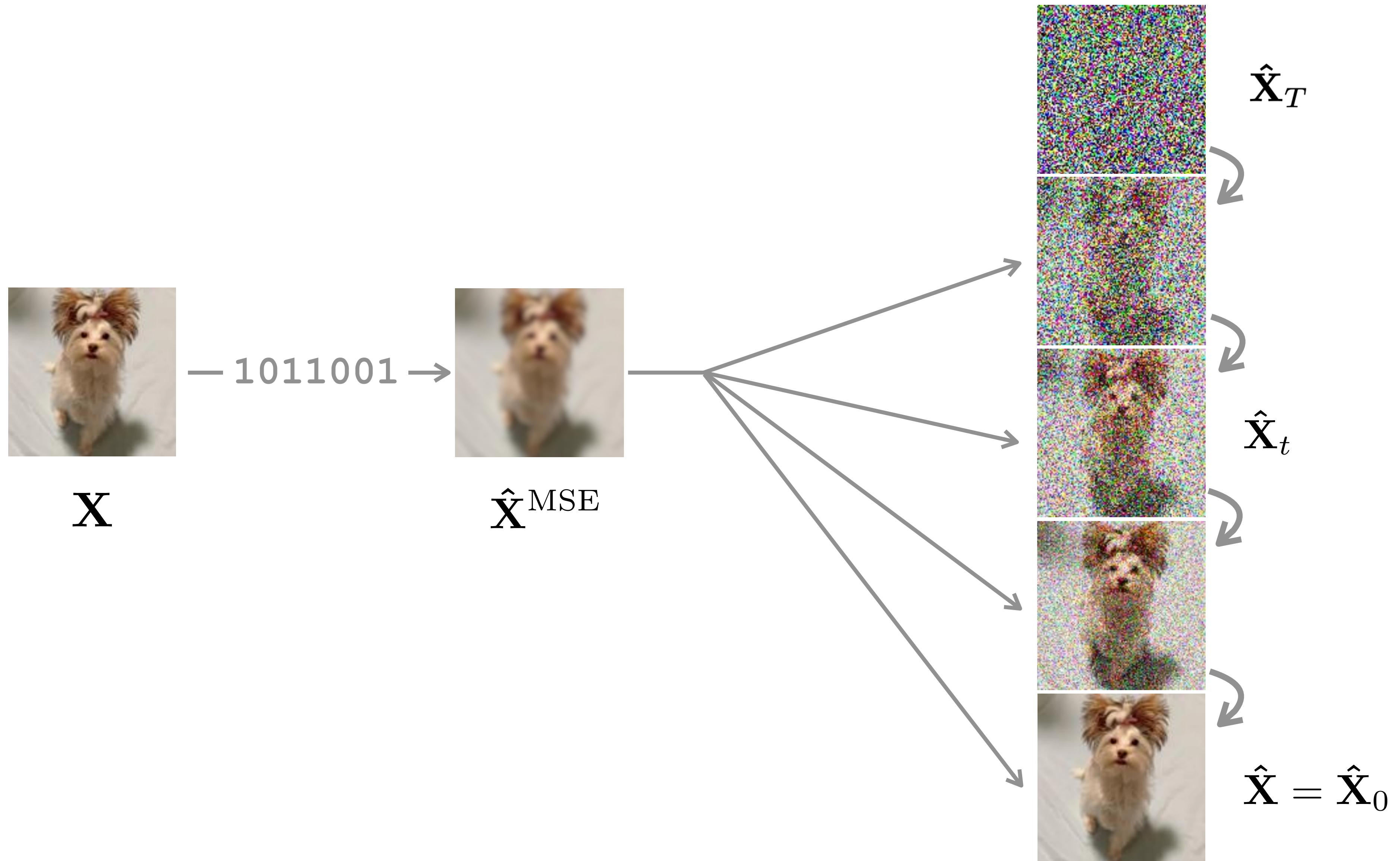


**x N** (potentially different models for different bit-rates)

## DiffC

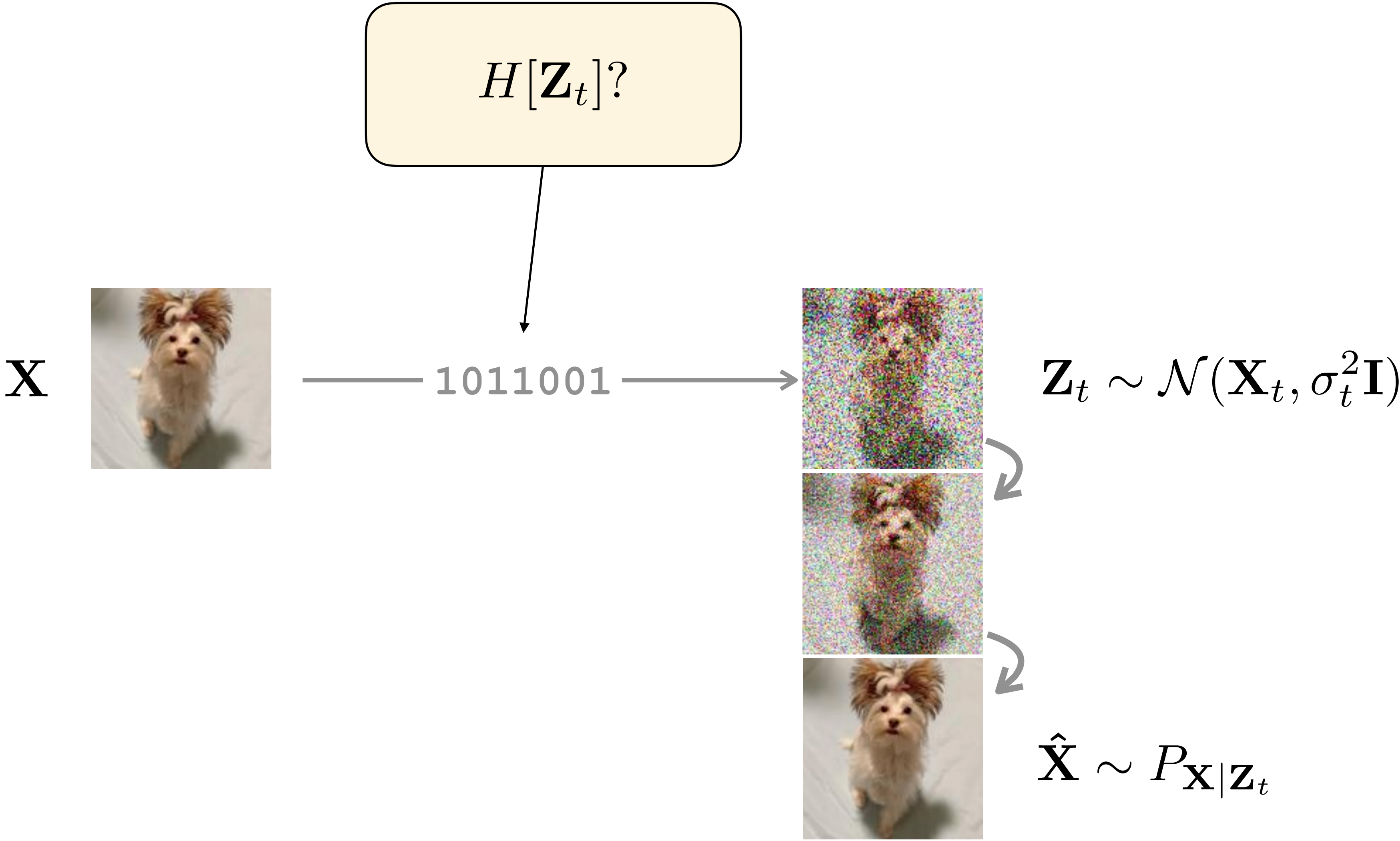
- A single unconditional diffusion model for all bit-rates
- *Reverse channel coding*

# HFD

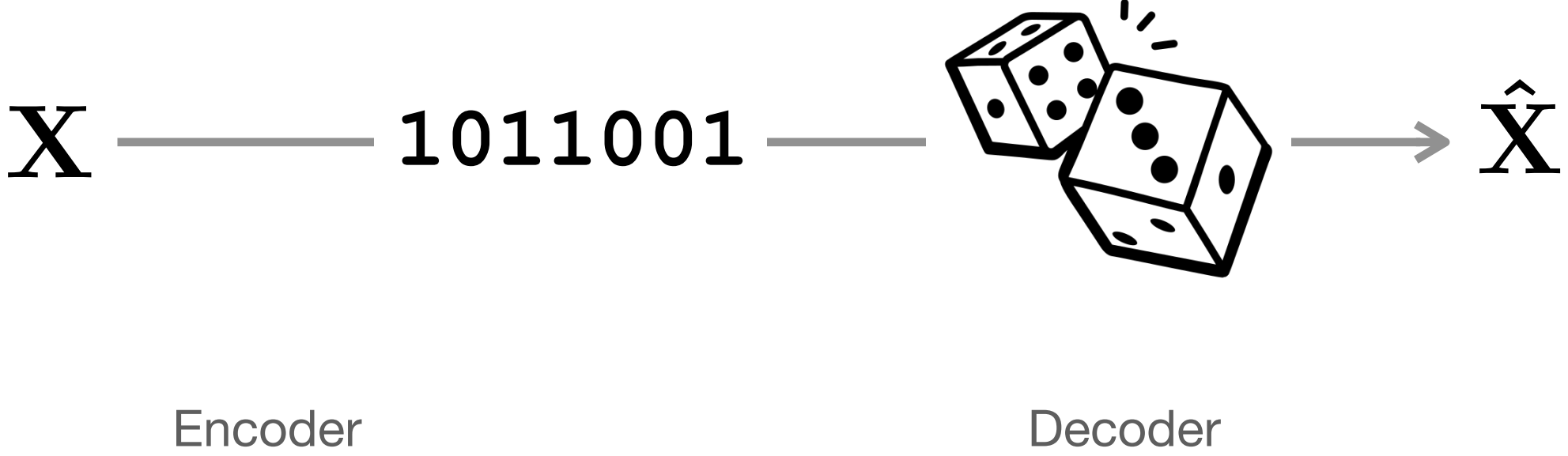




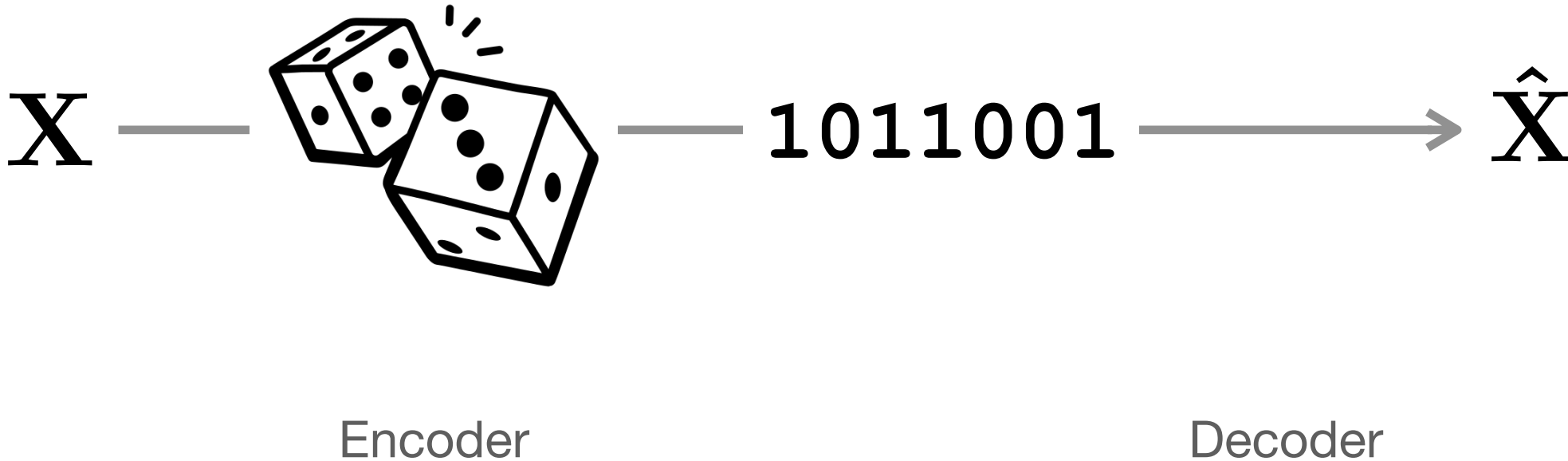
DiffC



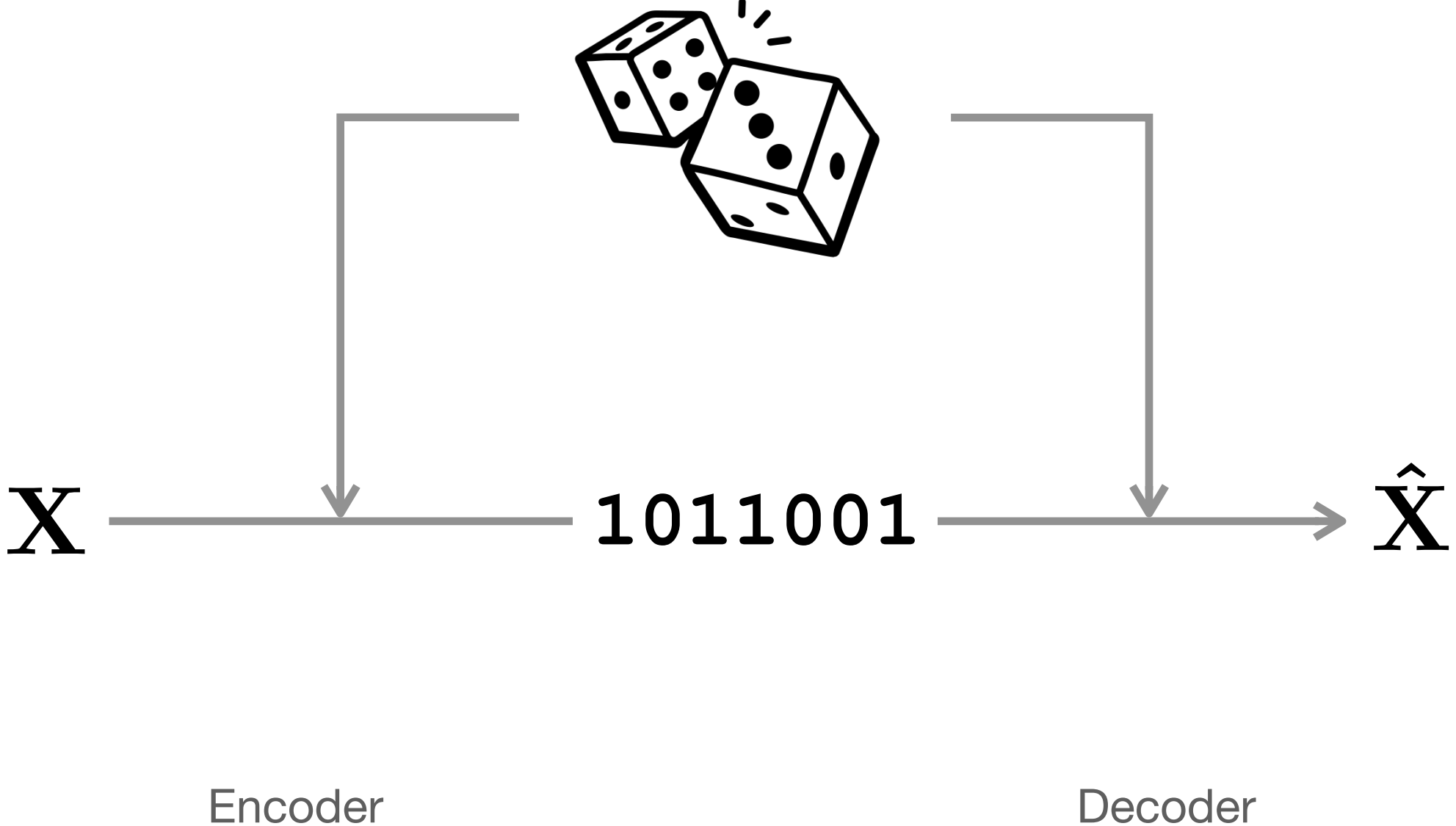
# Decoder randomness



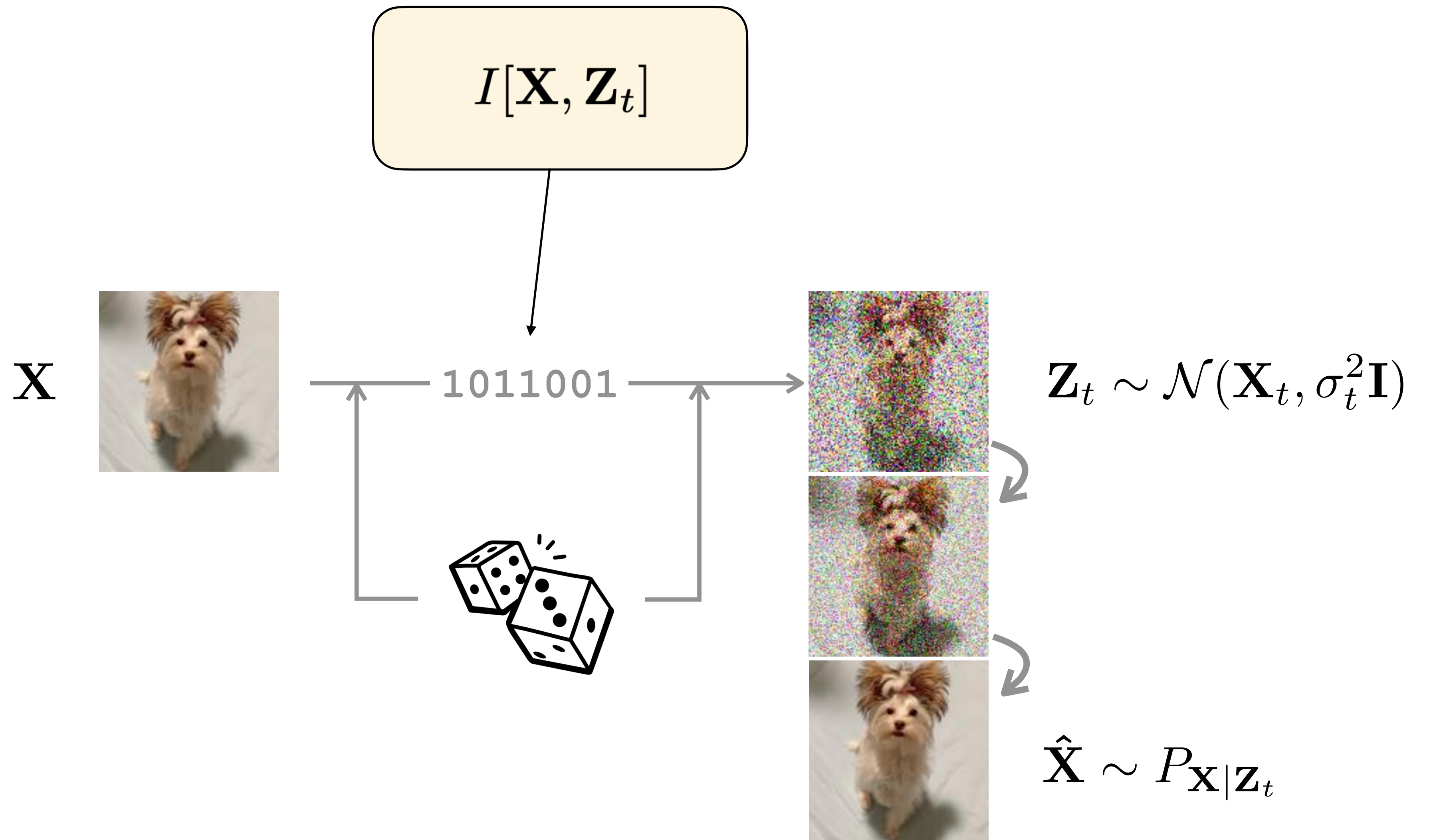
# Encoder randomness



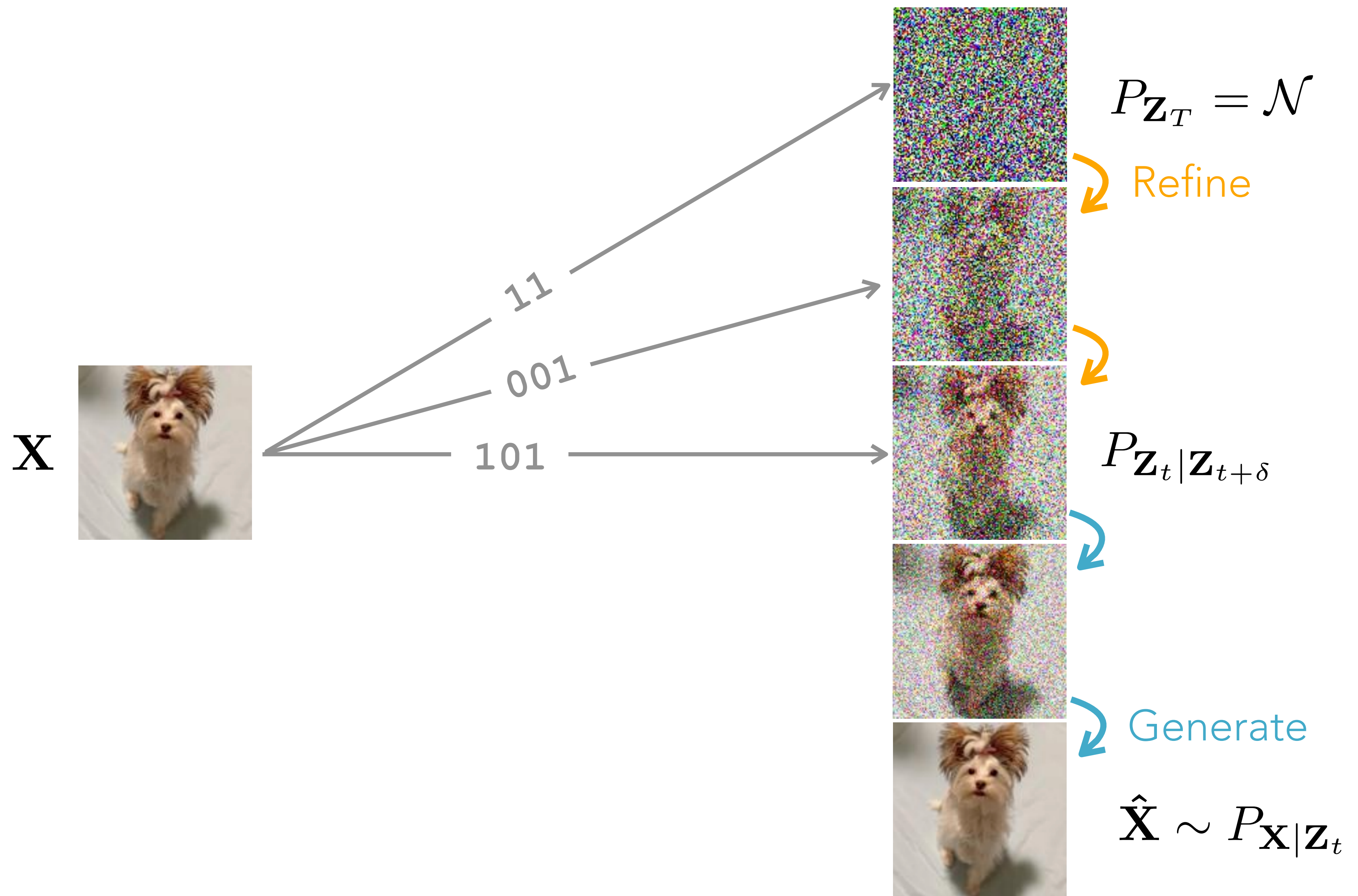
# Shared randomness



# DiffC



# DiffC



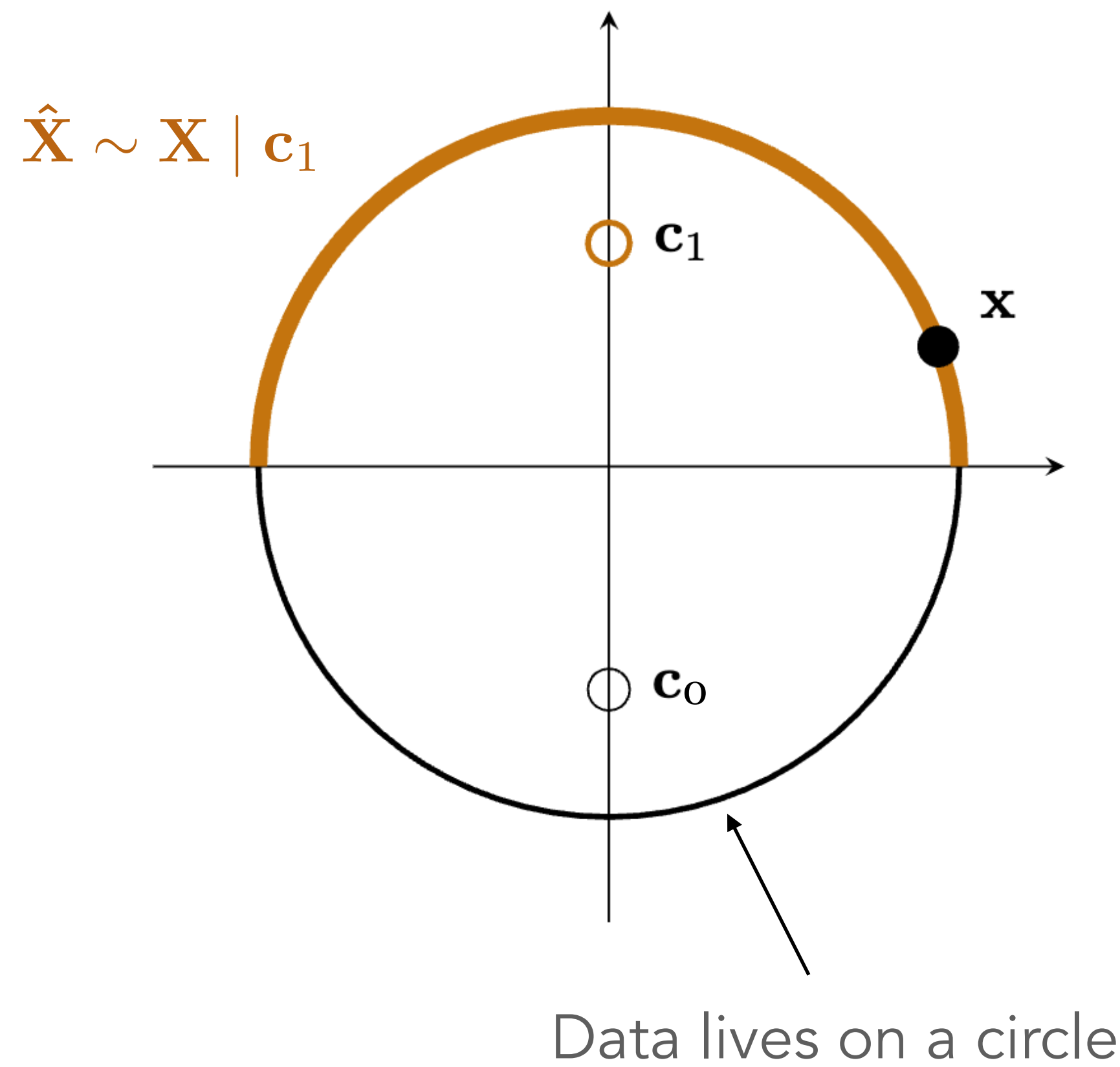
# Universal quantization

$$Z = \underbrace{[Y - U]}_K + U \sim Y + U'$$

Encoder      Decoder

$U, U' \sim \text{Uniform}([-1/2, 1/2])$

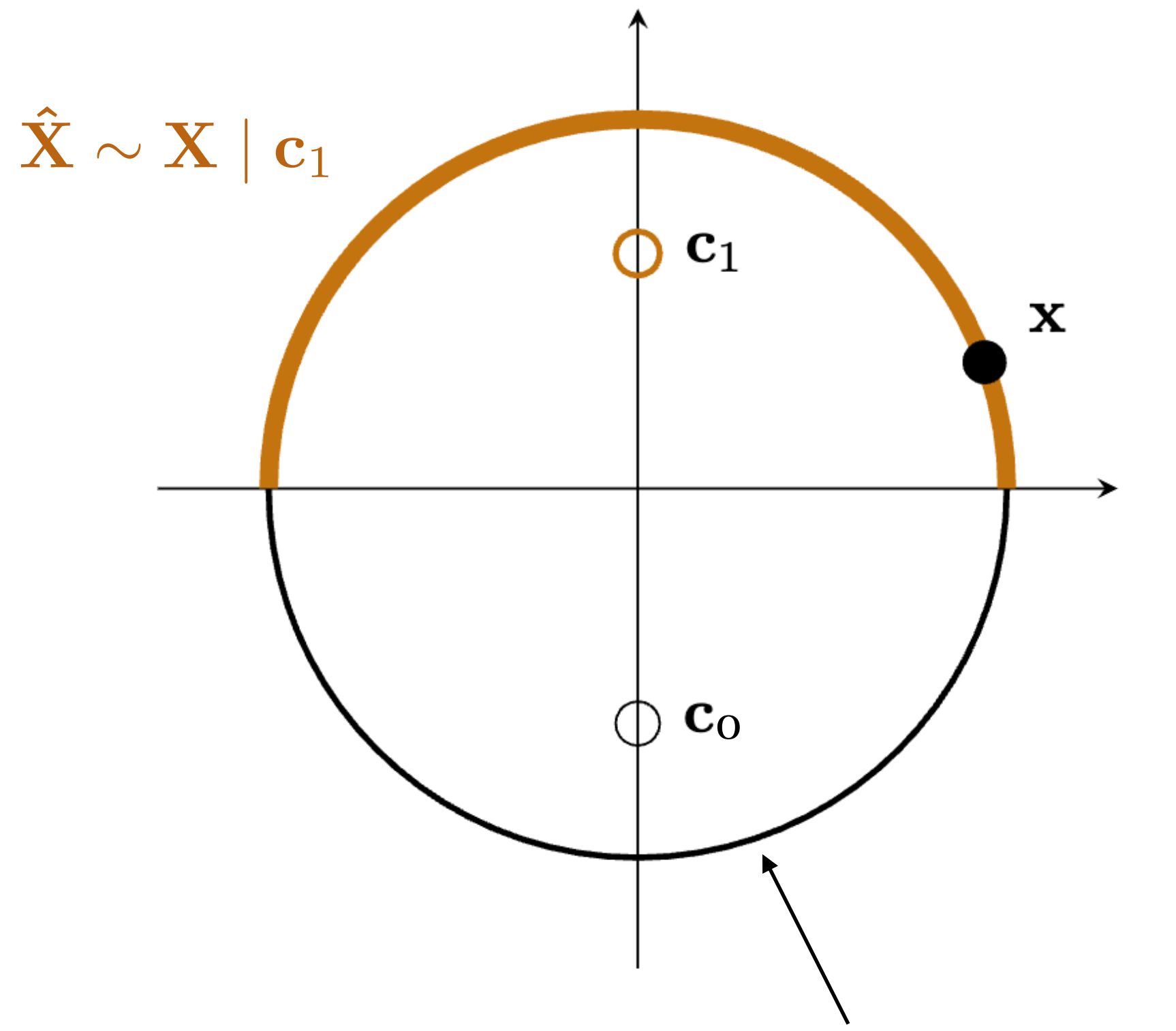
# Toy example



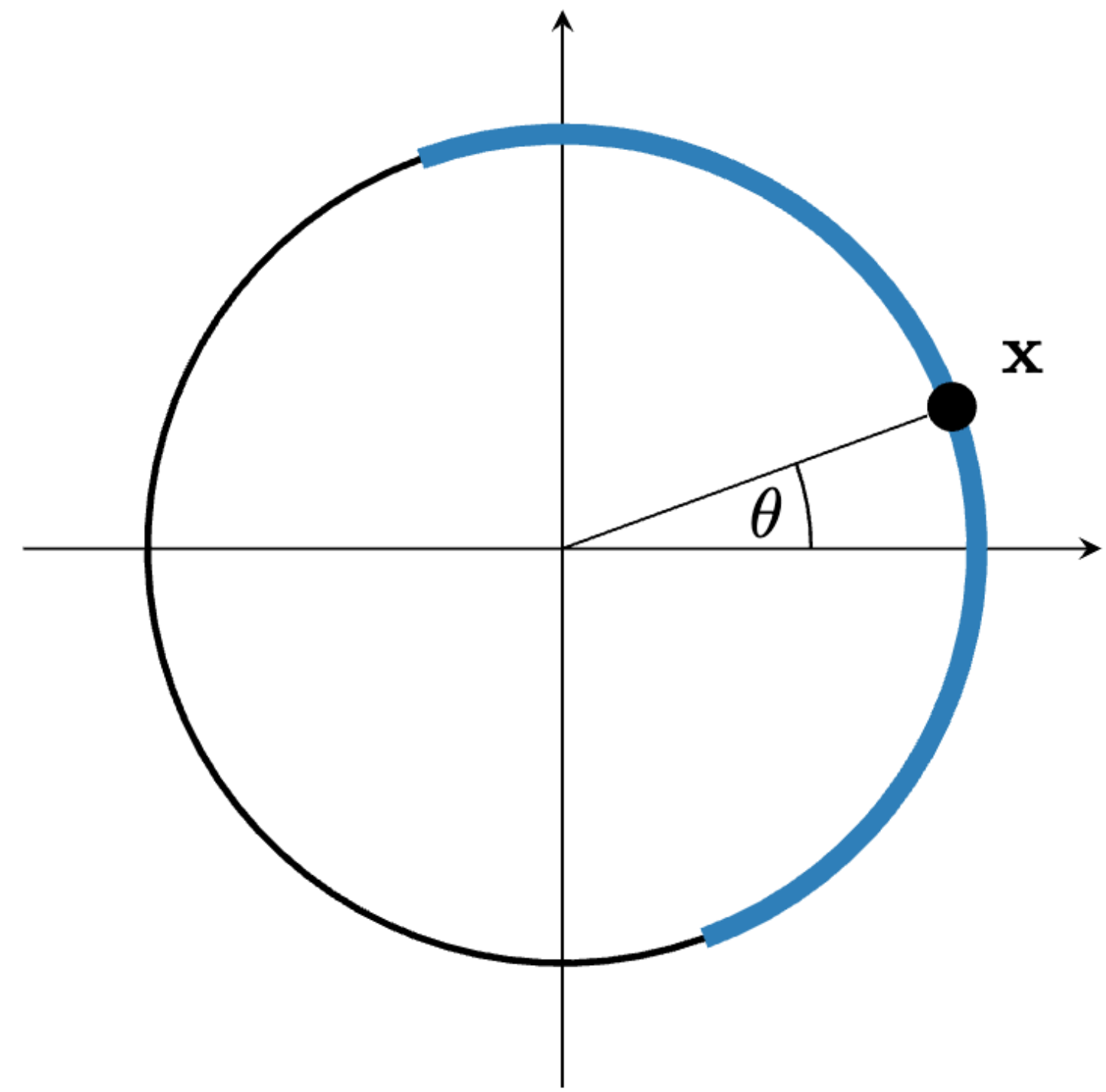


# Toy example

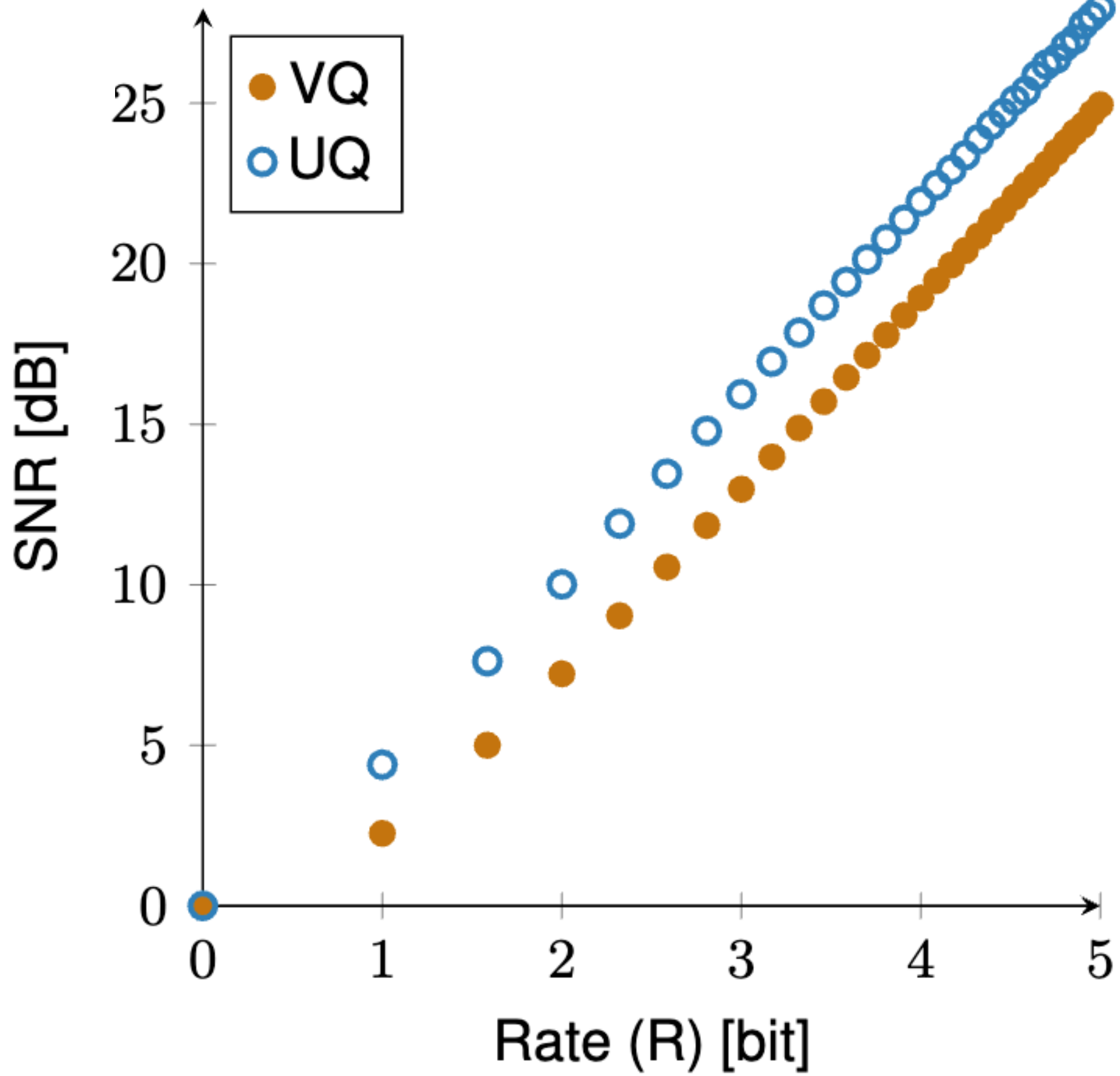
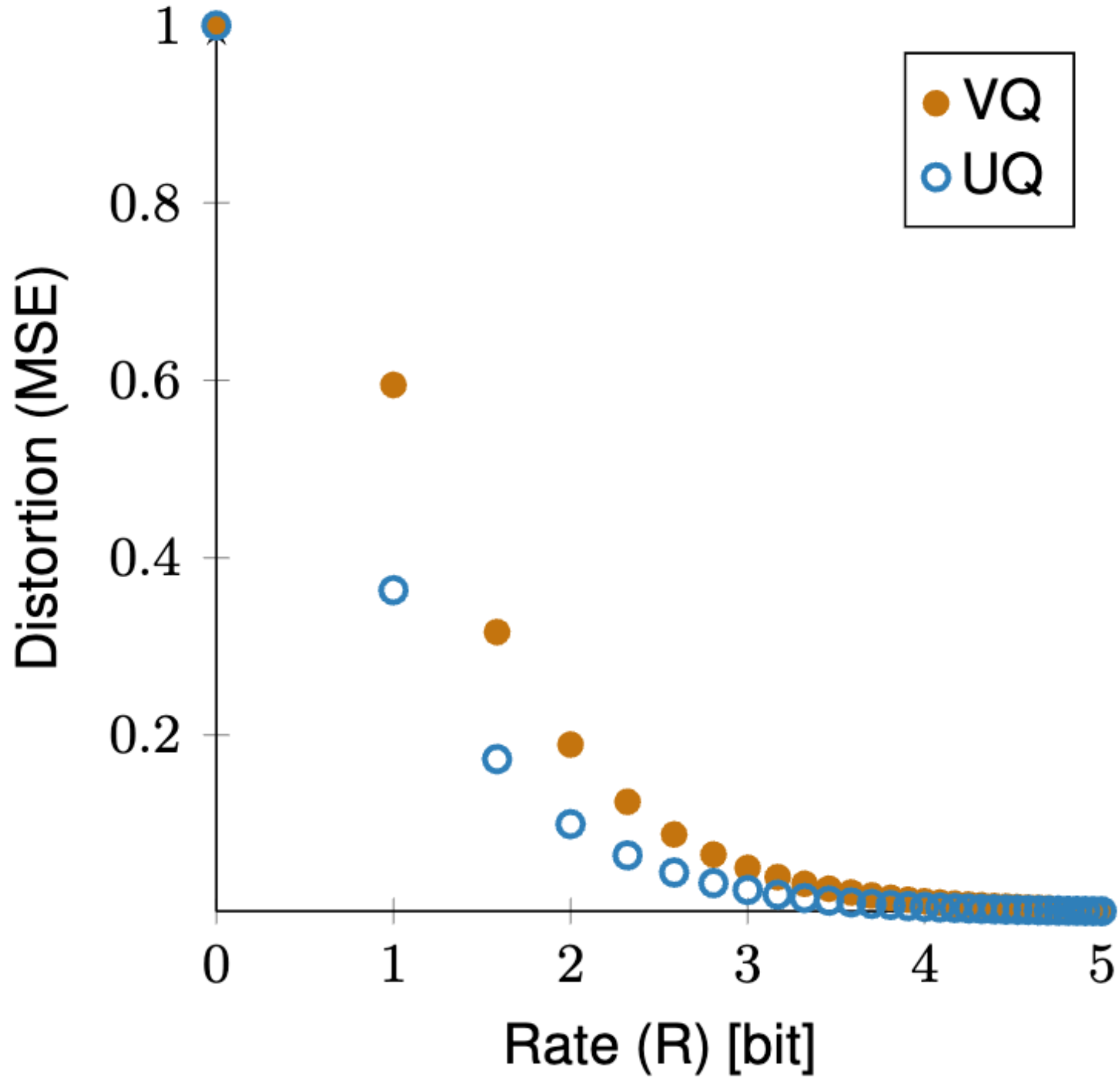
$$K = \left\lfloor \frac{\Theta}{\pi} - U \right\rfloor \text{ mod } 2$$
$$\hat{\Theta} = \pi(K + U)$$



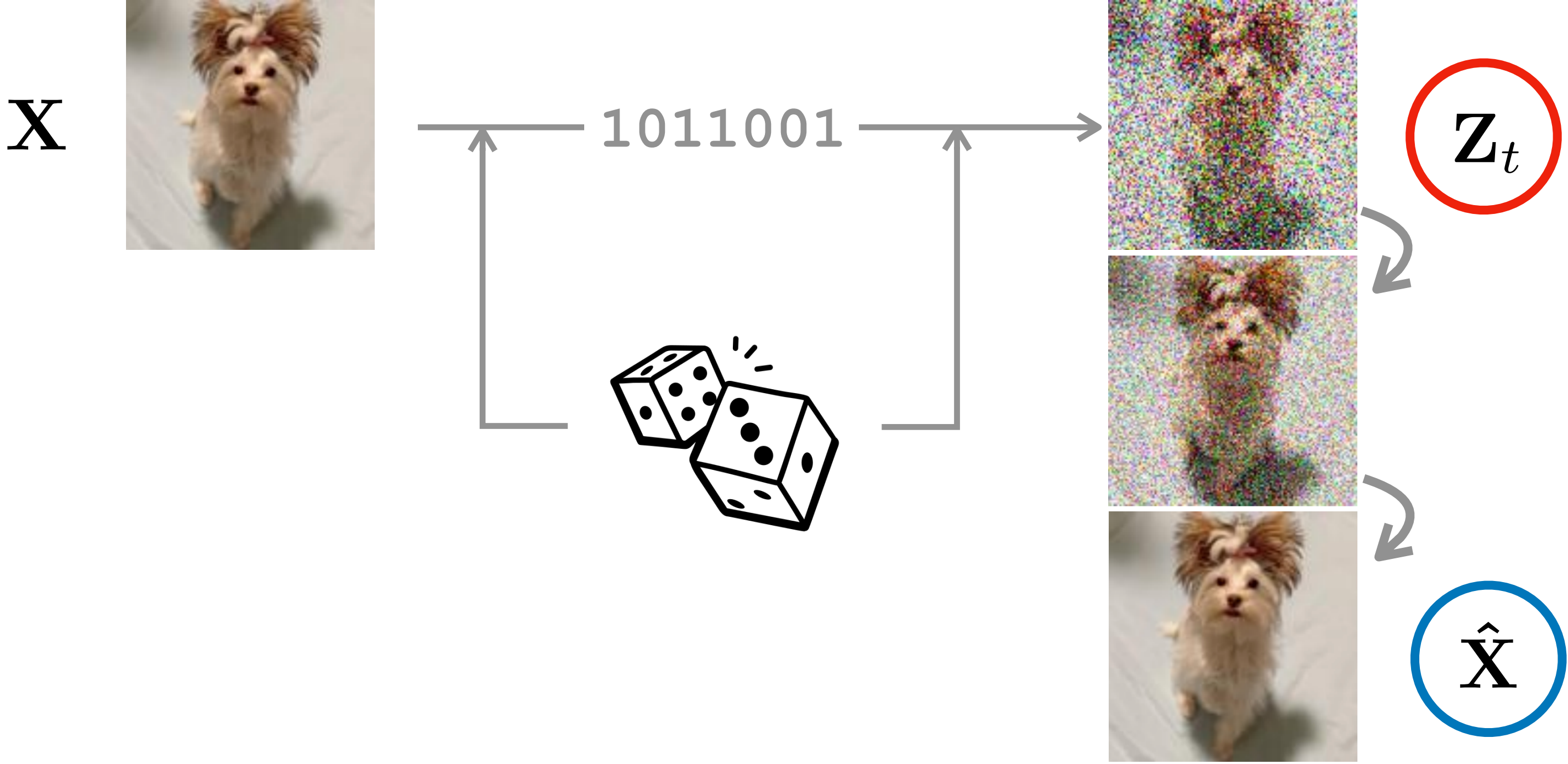
Data lives on a circle



# Toy example



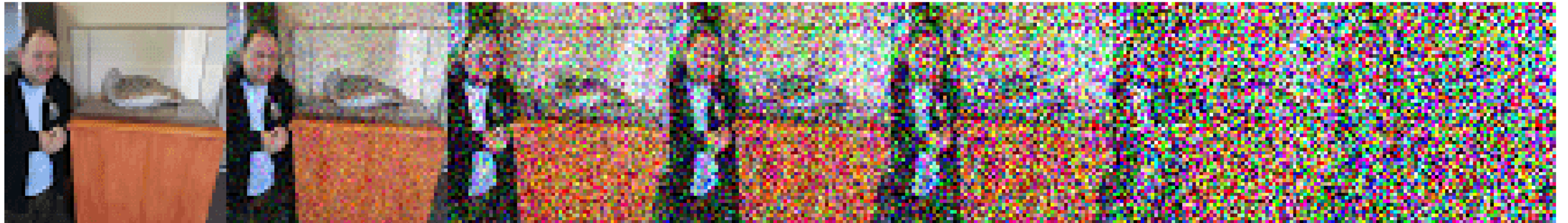
# DiffC



# DiffC

9.1719    0.5654    0.2421    0.1916    0.1297    0.0538    0.0256

$Z_t$



$\hat{X}$

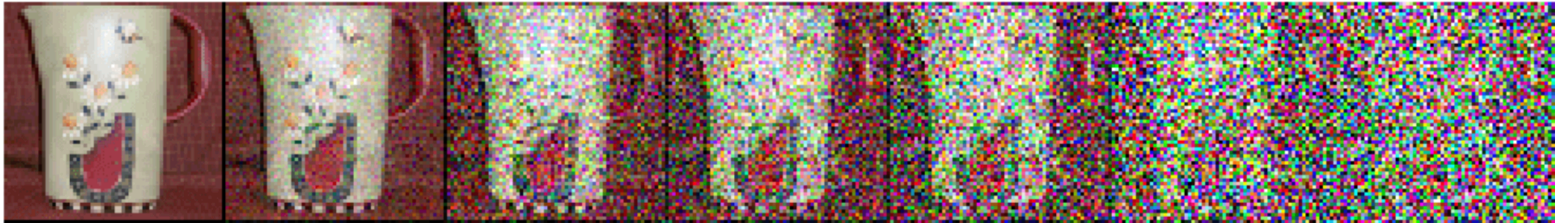


32.4dB    26.7dB    25.4dB    23.3dB    18.6dB    15.7dB

DiffC

10.4572 0.6486 0.2554 0.1974 0.1240 0.0429 0.0198

$Z_t$

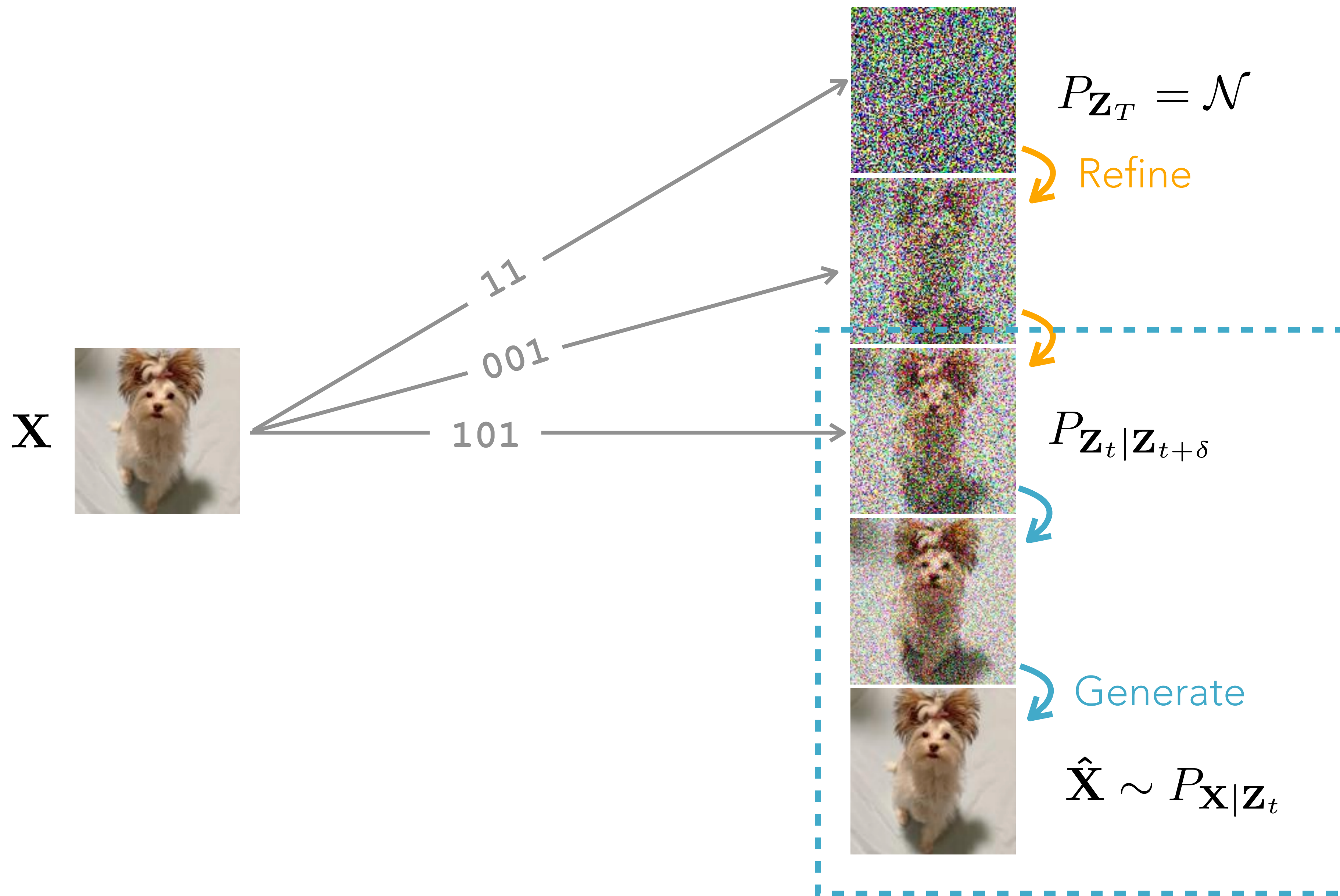


$\hat{X}$

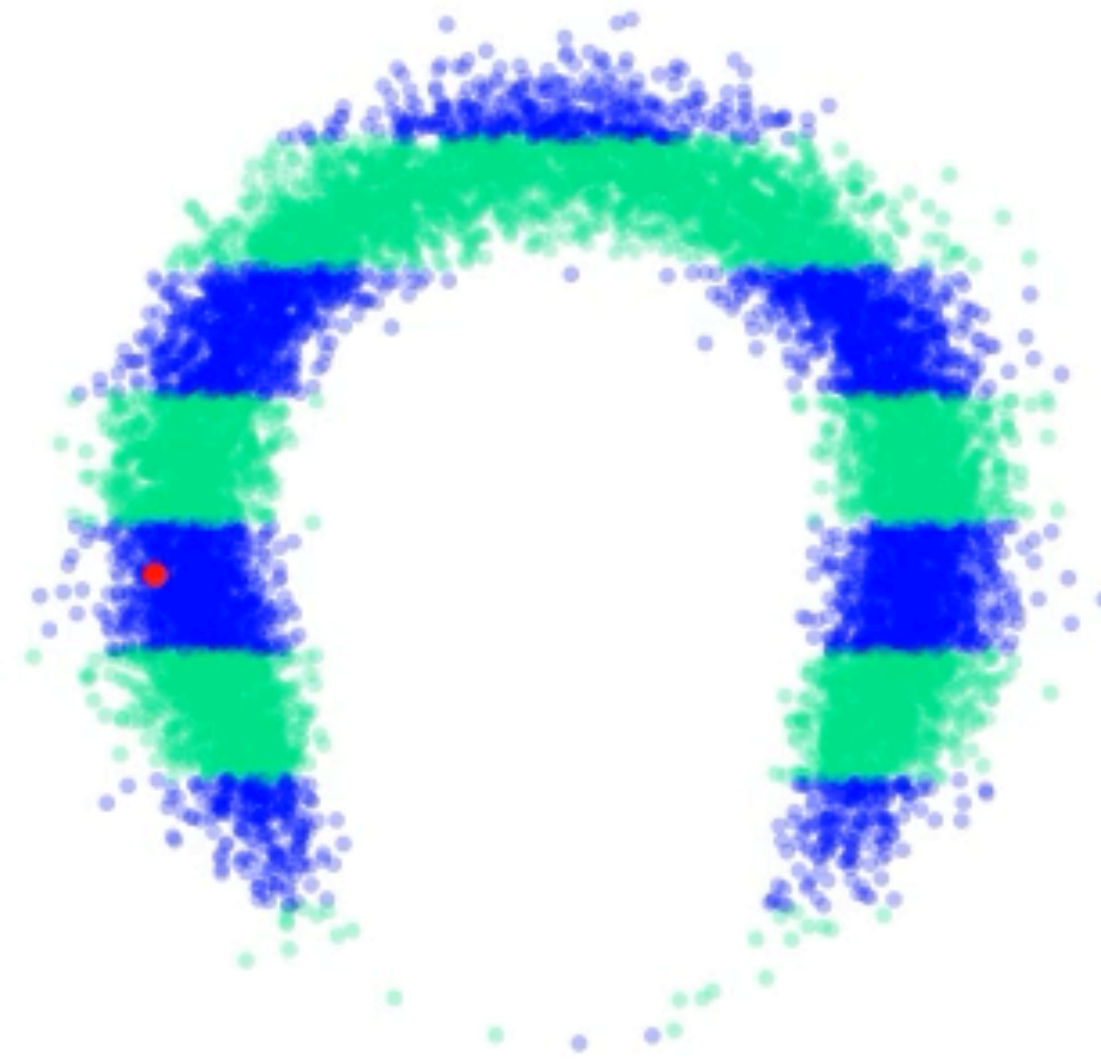


31.7dB 25.8dB 24.7dB 22.4dB 19.3dB 16.3dB

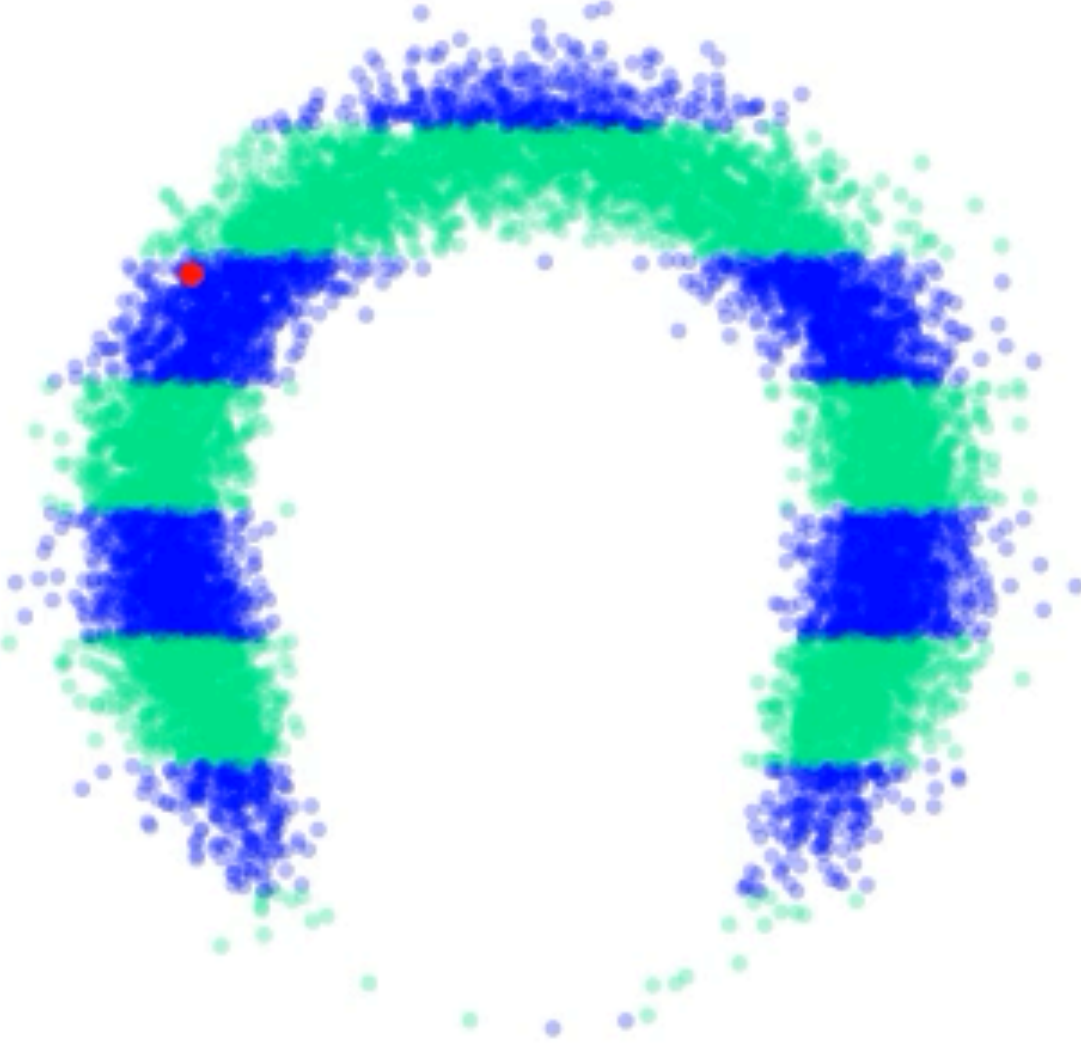
# DiffC



DiffC-A (SDE)



DiffC-F (ODE)





# DiffC-F vs DiffC-A

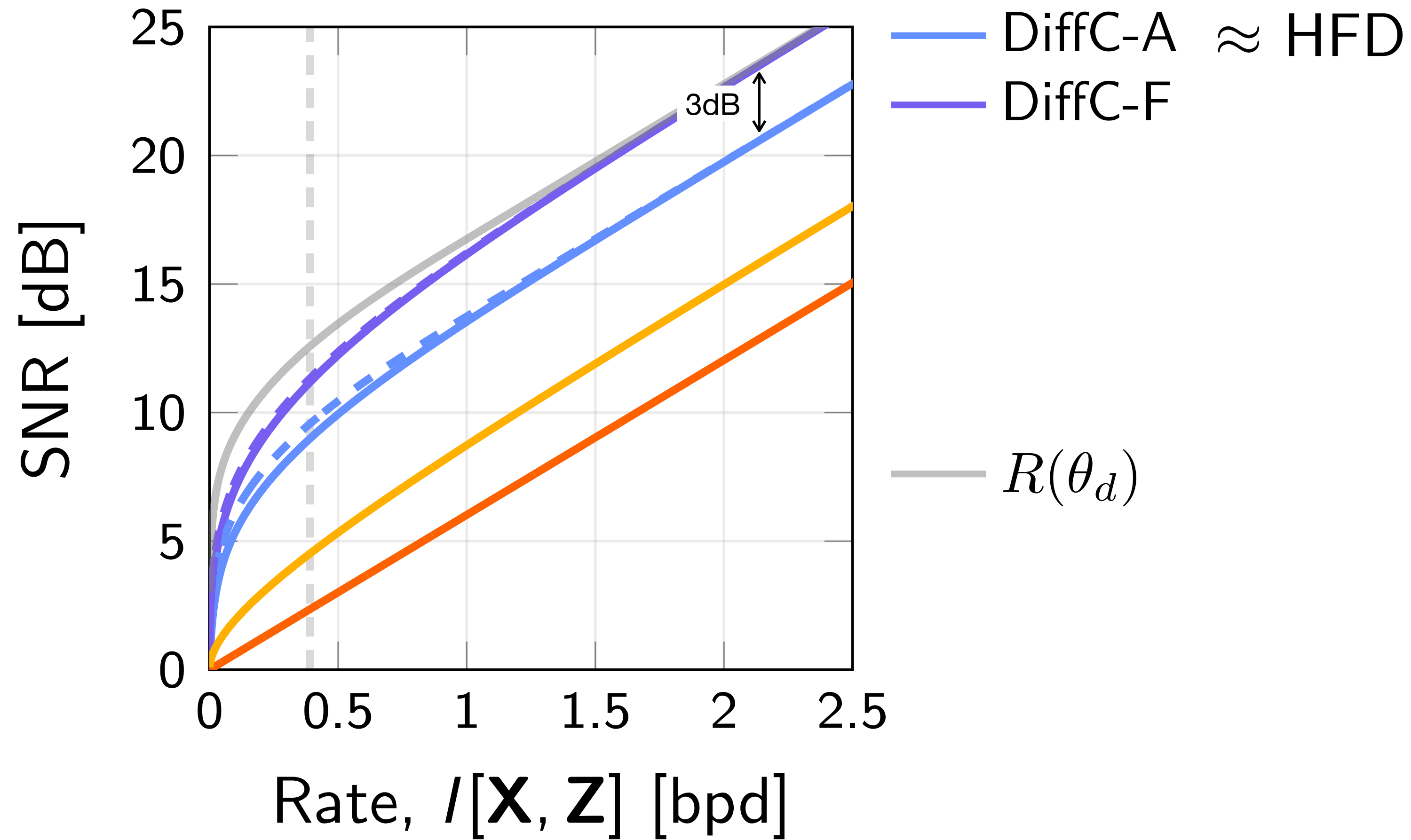
**Theorem.** Let  $\mathbf{X} : \Omega \rightarrow \mathbb{R}^M$  have a smooth density  $p$  with finite

$$G = \mathbb{E}[\|\nabla \ln p(\mathbf{X})\|^2].$$

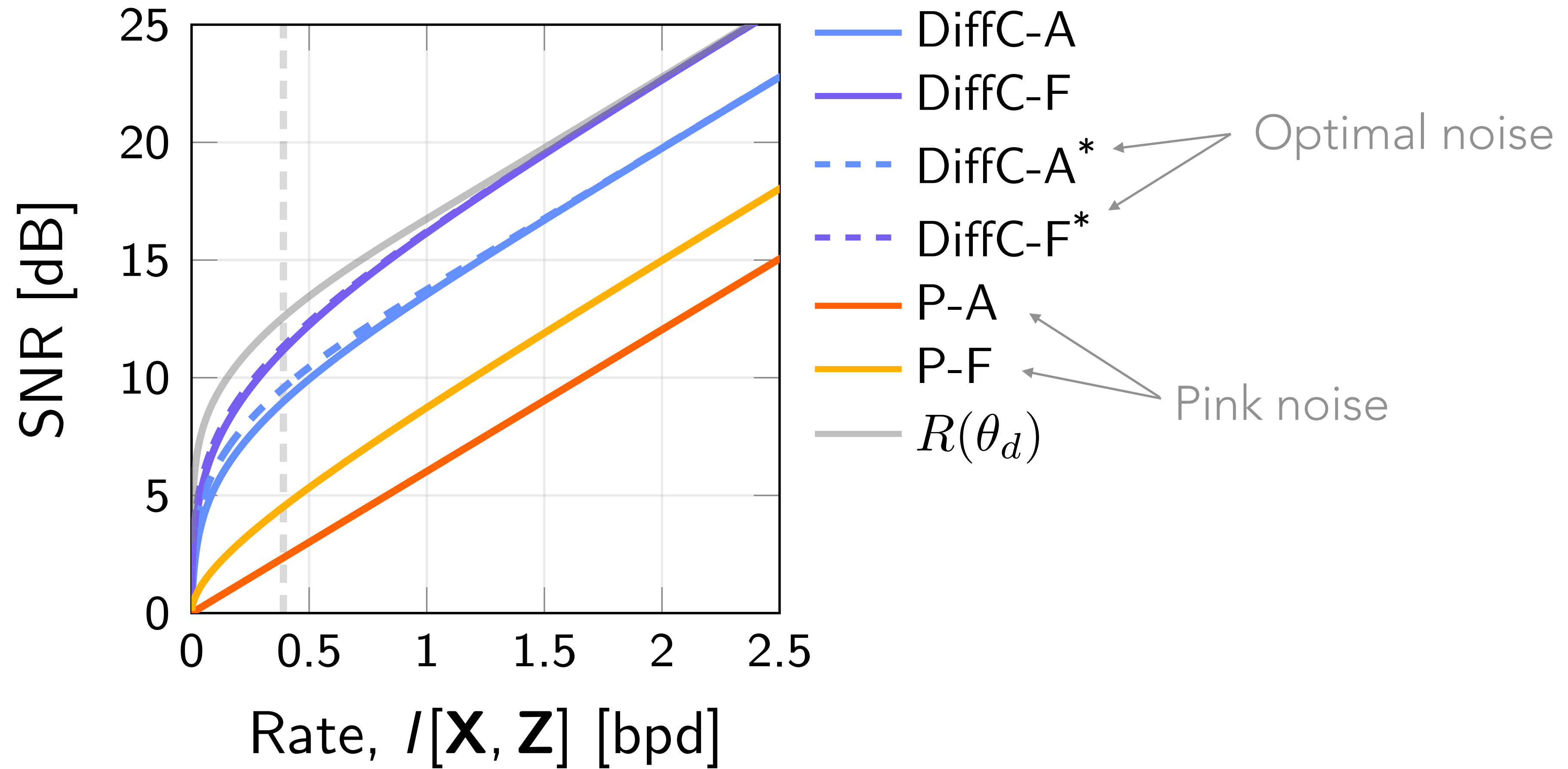
Let  $\mathbf{Z}_t = \sqrt{1 - \sigma_t^2} \mathbf{X} + \sigma_t \mathbf{U}$  with  $\mathbf{U} \sim \mathcal{N}(0, \mathbf{I})$ . Let  $\hat{\mathbf{X}}_A \sim P(\mathbf{X} \mid \mathbf{Z}_t)$  and let  $\hat{\mathbf{X}}_F = \mathbf{Z}_0$  be the solution to the ODE with  $\mathbf{Z}_t$  as initial condition. Then

$$\lim_{\sigma_t \rightarrow 0} \frac{\mathbb{E}[\|\hat{\mathbf{X}}_F - \mathbf{X}\|^2]}{\mathbb{E}[\|\hat{\mathbf{X}}_A - \mathbf{X}\|^2]} = \frac{1}{2}$$

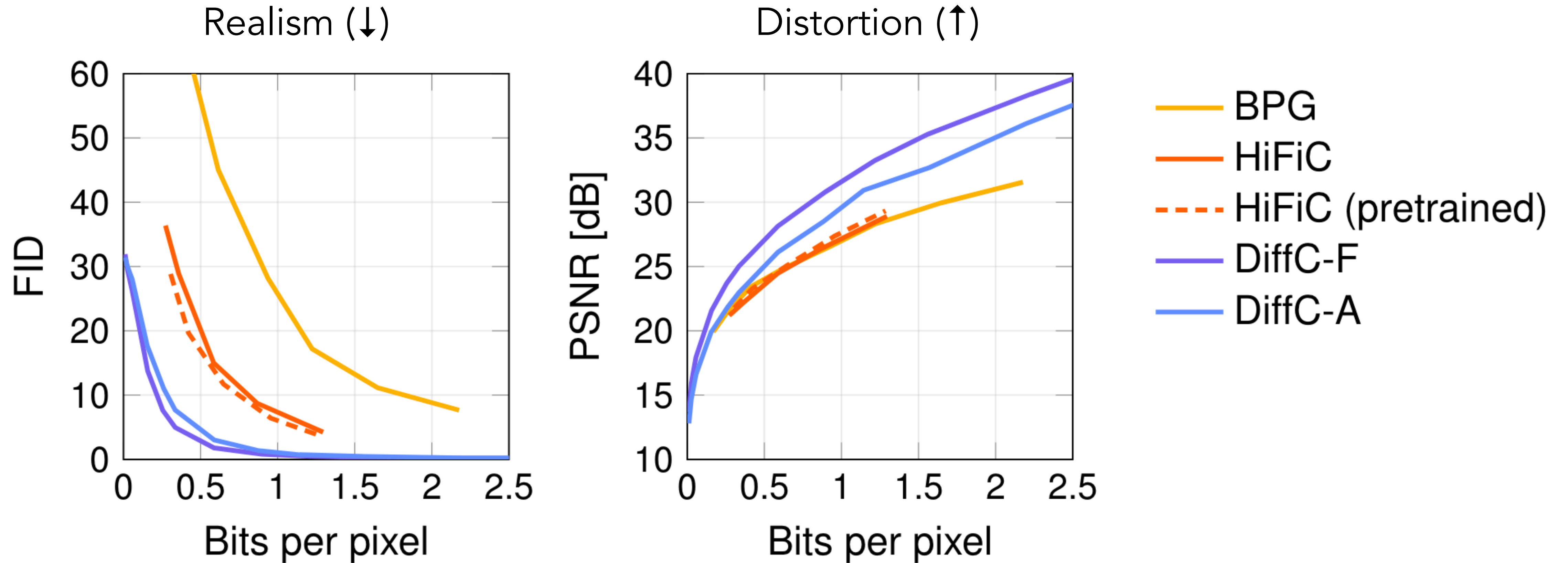
# Example: Multivariate Gaussian

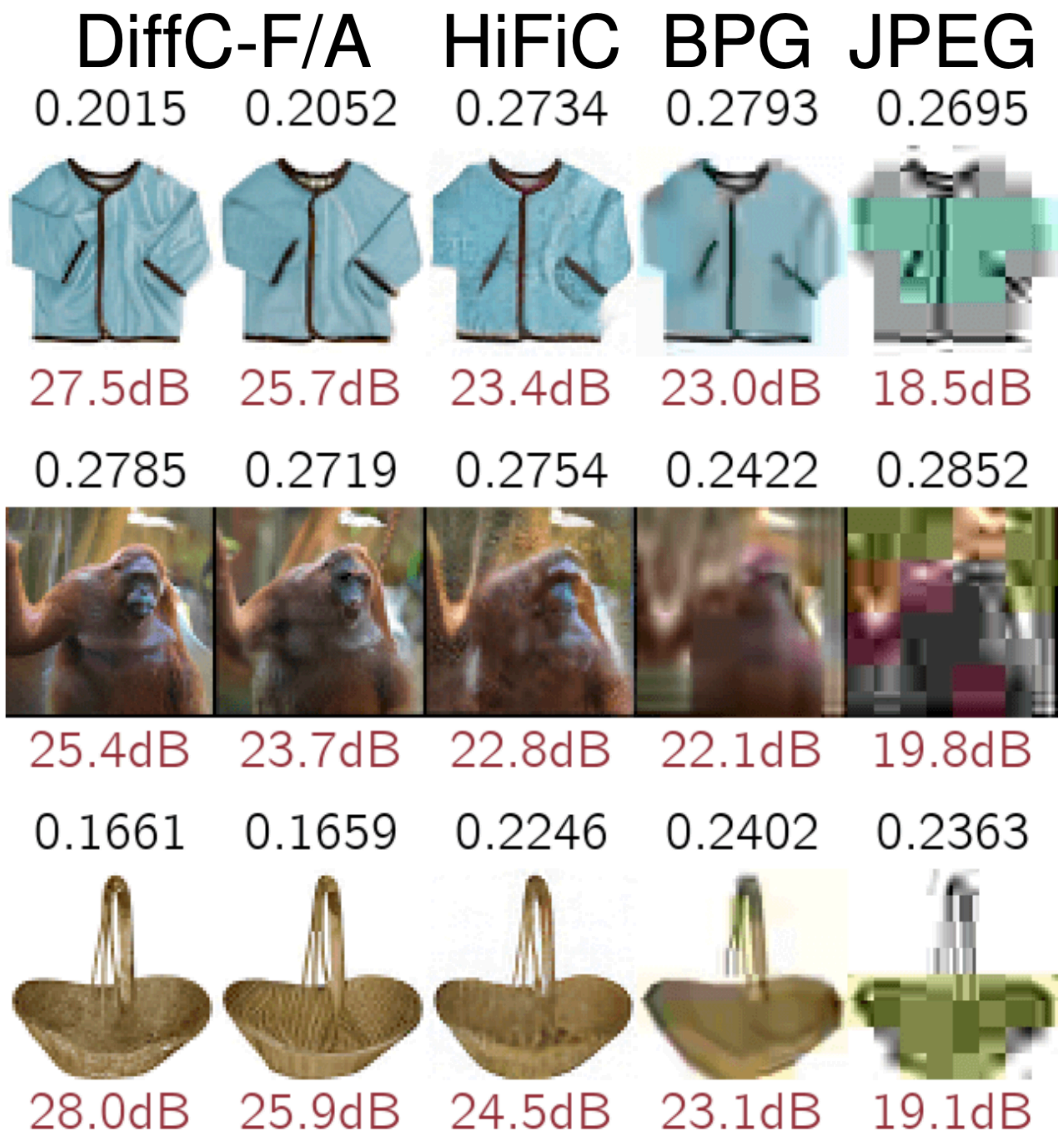


# Example: Multivariate Gaussian



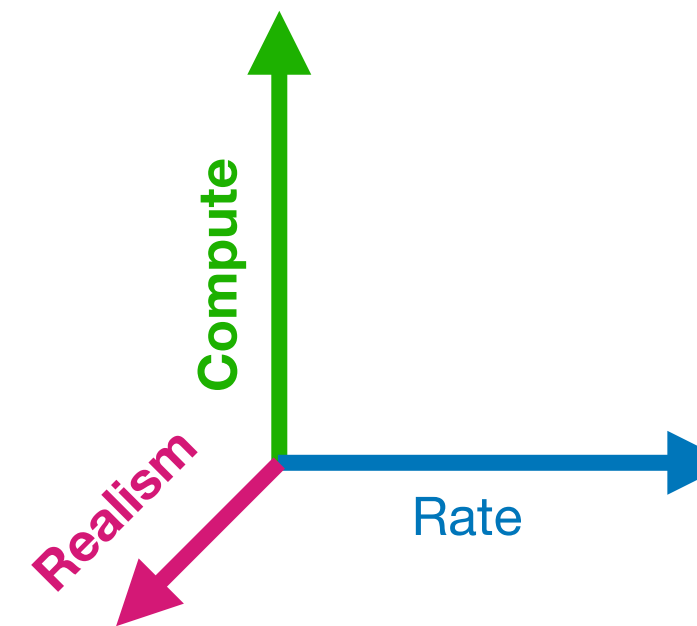
# ImageNet 64x64





# Future research

- **What do good analysis transforms and distortions look like?**
  - Assuming perfect realism, what should the distortion measure?
  - Assuming perfect realism, what does a 1dB change in PSNR mean perceptually?
- **Make diffusion-based compression practical**
  - Fast Gaussian reverse channel coding
  - Dithered quantization instead of Gaussian noise
  - Rectified flows (or “flow matching”)
  - Distillation methods
- What are the limits of **low rate, low complexity, & high realism?**



# Realism revisited

What makes an image realistic?

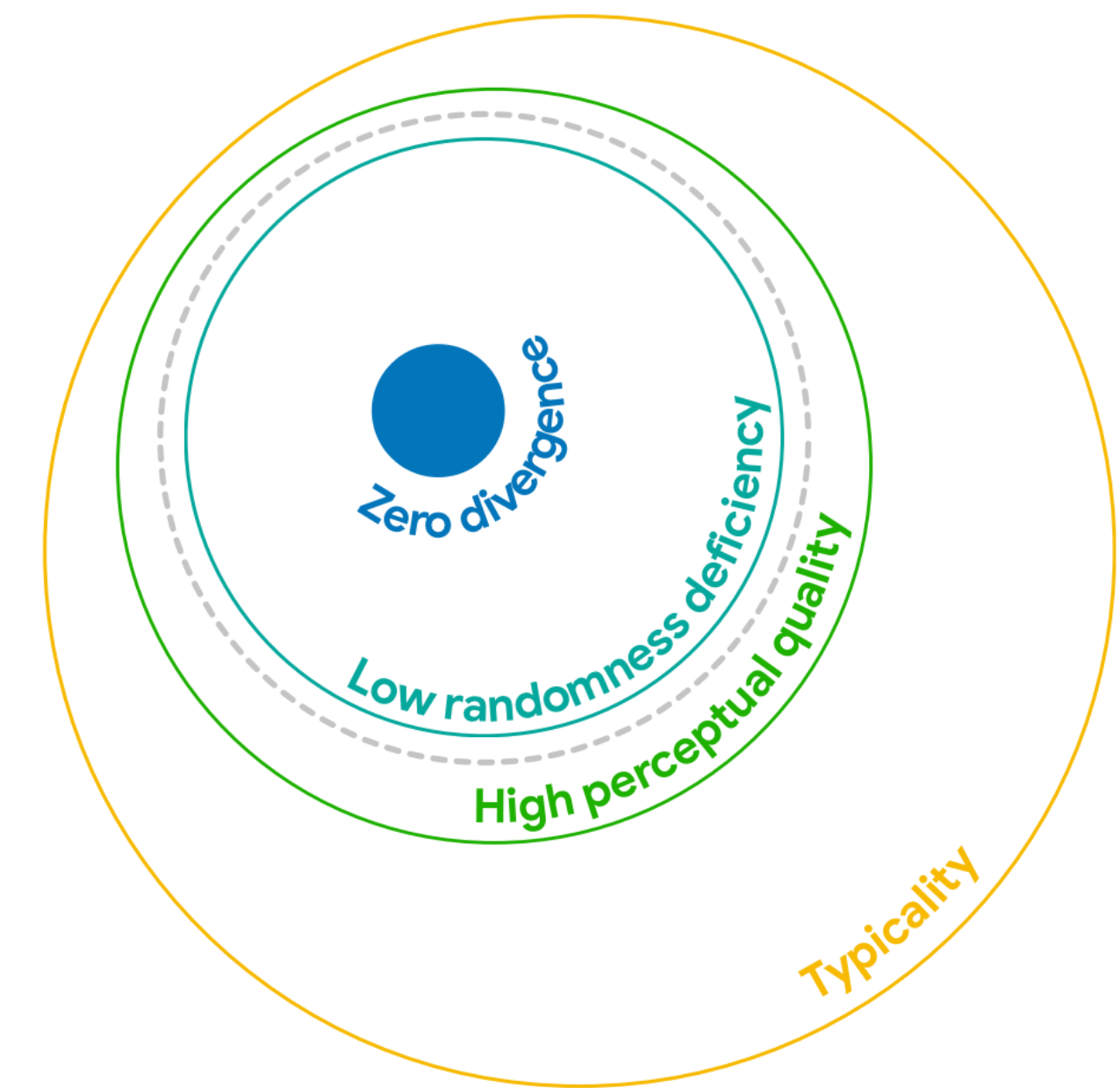
ICML 2024

“Universal critic”

$$U(\mathbf{x}) = \log \sum_k \pi_k Q_k(\mathbf{x}) - \log P(\mathbf{x})$$

$$U(\mathbf{x}_1, \dots, \mathbf{x}_N) = \log \sum_k \pi_k \prod_n Q_k(\mathbf{x}_n) - \log \prod_n P(\mathbf{x}_n)$$

$$\rightarrow ND_{\text{KL}}[Q||P] \quad \text{where} \quad \mathbf{x}_n \sim Q$$



# Thank you

Collaborators

