

BiGR: Harnessing Binary Latent Codes for Image Generation and Improved Visual Representation Capabilities

Shaozhe Hao¹, Xuantong Liu², Xianbiao Qi³, Shihao Zhao¹, Bojia Zi⁴,
Rong Xiao³, Kai Han¹, Kwan-Yee K. Wong¹

¹The University of Hong Kong

²Hong Kong University of Science and Technology

³Intellifusion

⁴The Chinese University of Hong Kong



香港大學

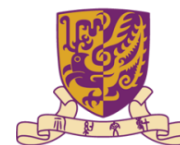
THE UNIVERSITY OF HONG KONG



香港科技大學

THE HONG KONG
UNIVERSITY OF SCIENCE
AND TECHNOLOGY

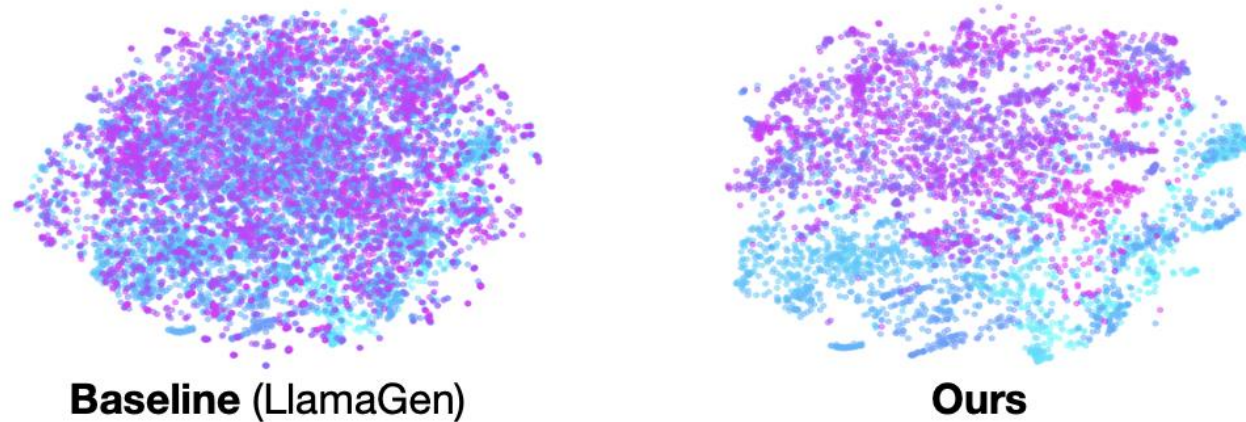
intell**if**usion
云天励飞



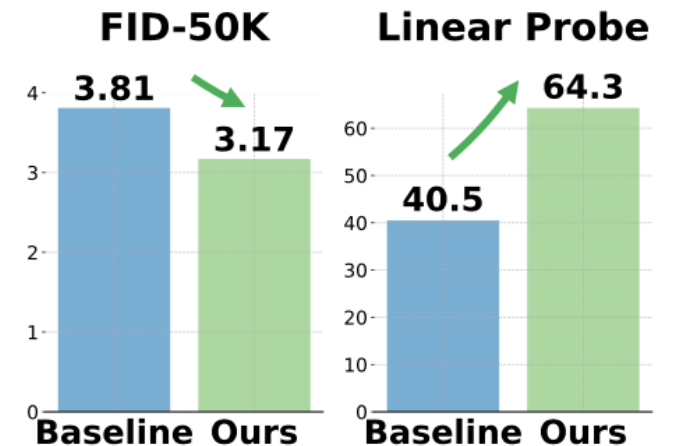
香港中文大學

The Chinese University of Hong Kong

Can image generative models produce discriminative visual representations?

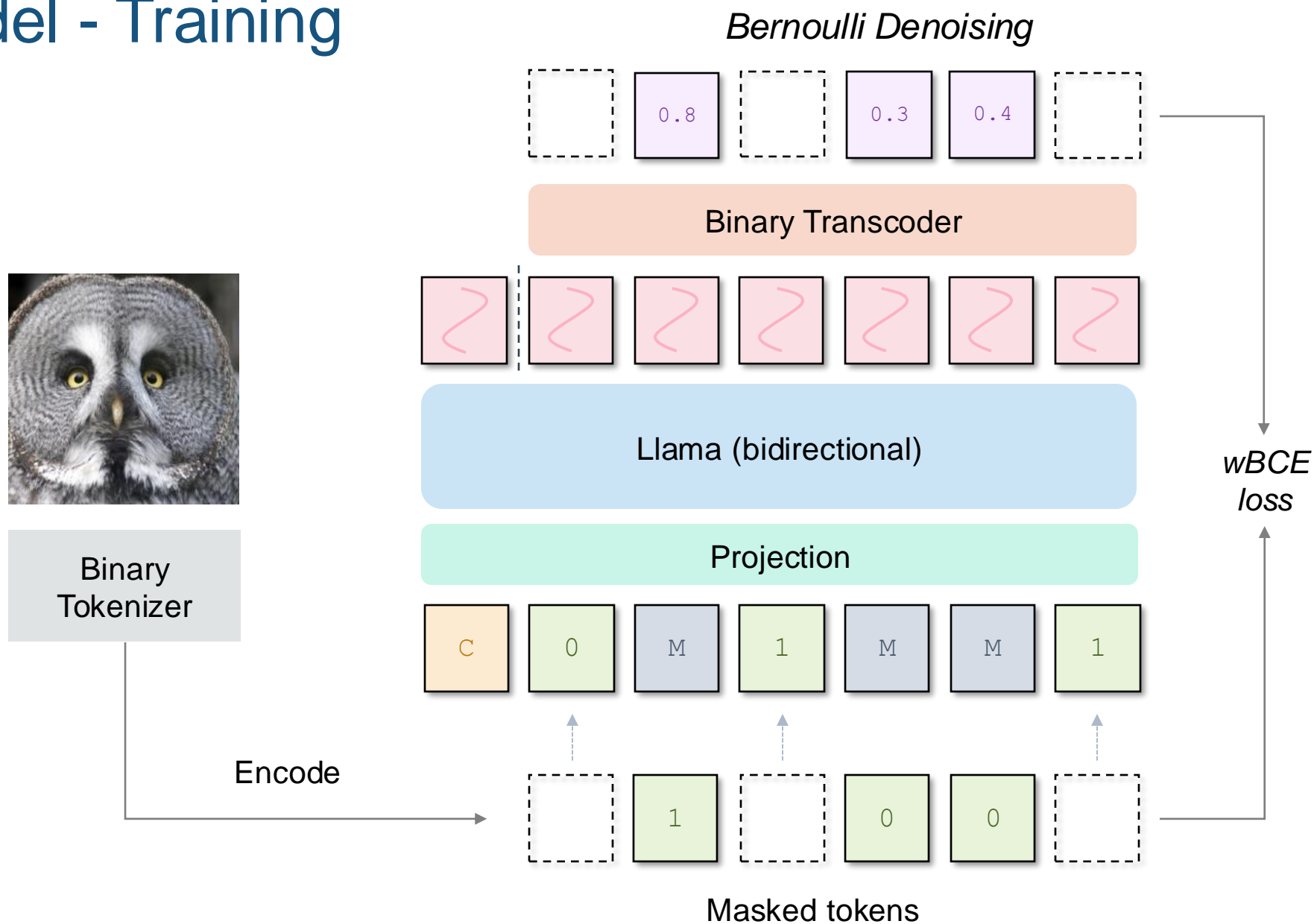


Samples from the same class are in the same color.



We introduce a novel conditional image generation model, which unifies **generative** and **discriminative** tasks.

Our model - Training



Binary transcoder - Bernoulli diffusion

➤ Forward diffusion process:

$$q(z^t|z^{t-1}) = \mathcal{B}(z^t; z^{t-1}(1 - \beta^t) + 0.5\beta^t) \quad t = 1, 2, \dots, T.$$

z^t : z at timestep t

β^t : predefined coefficient at timestep t

➤ Backward denoising process:

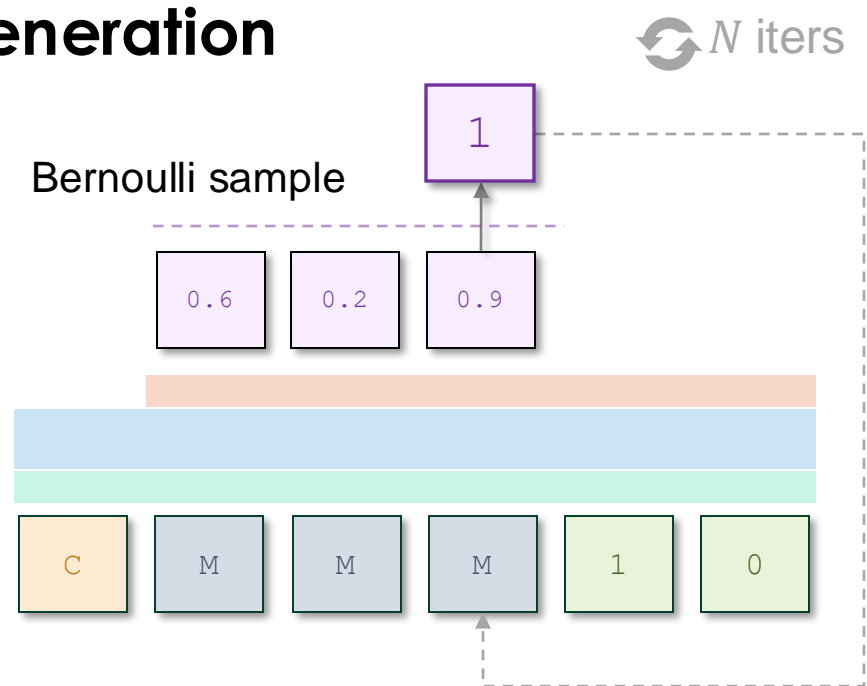
$$p_\phi(z^{t-1}|z^t) = \mathcal{B}(z^{t-1}; S(g_\phi(z^t, t, h)))$$

S : sigmoid function

g_θ : denoising network (MLPs)

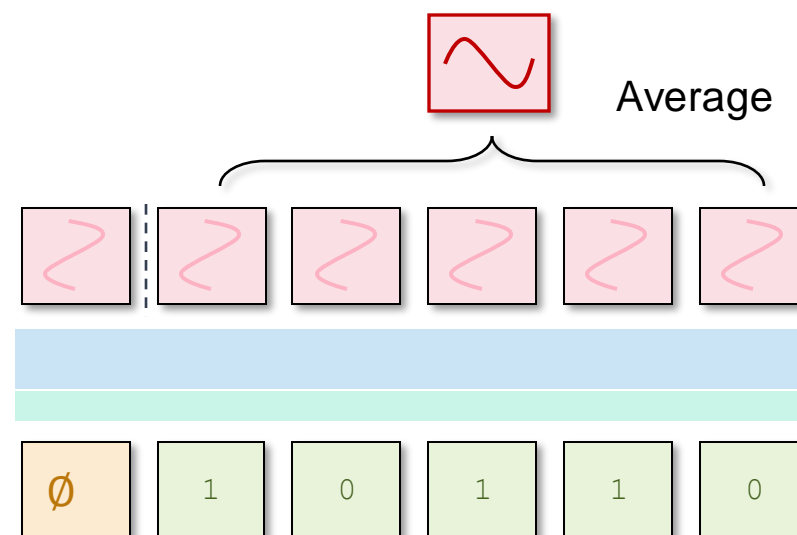
Inference

Generation



$$\mathcal{H} = -\frac{1}{K} \sum_{k=1}^K p_k \log_2 p_k + (1 - p_k) \log_2 (1 - p_k),$$

Representation



Evaluation

Table 1: **Uniformity comparison.** We compare the generative and discriminative performance of our model against LlamaGen (Sun et al., 2024) and three other settings, varying by tokenizers, training objectives, and modeling types. We use KV cache to accelerate all AR models.

Model	Tokenizer	Objective	Type	Time↓	Generative					Discriminative	
					FID↓	IS↑	sFID↓	Pre.↑	Rec.↑	ACC1	ACC5
LlamaGen	VQGAN	Cat.	AR	0.13	3.81	248.28	8.49	0.83	0.52	40.5	64.4
S0	B-AE	Cat.	AR	0.15	3.21	239.17	5.38	0.83	0.54	23.8	44.2
S1	B-AE	Cat.	Mask	0.10	3.85	261.81	6.10	0.85	0.47	61.1	83.2
S2	B-AE	Bin.	AR	1.04	7.50	164.31	6.56	0.85	0.41	45.2	69.3
S3 (Ours)	B-AE	Bin.	Mask	0.69	3.17	262.14	5.59	0.86	0.50	64.3	85.4

Model scaling

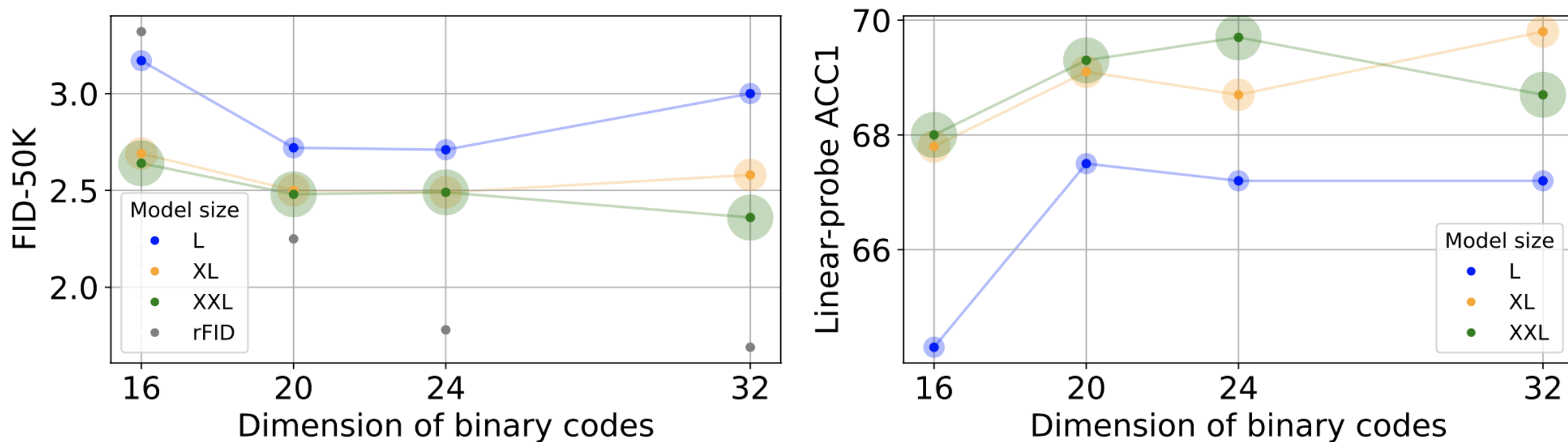
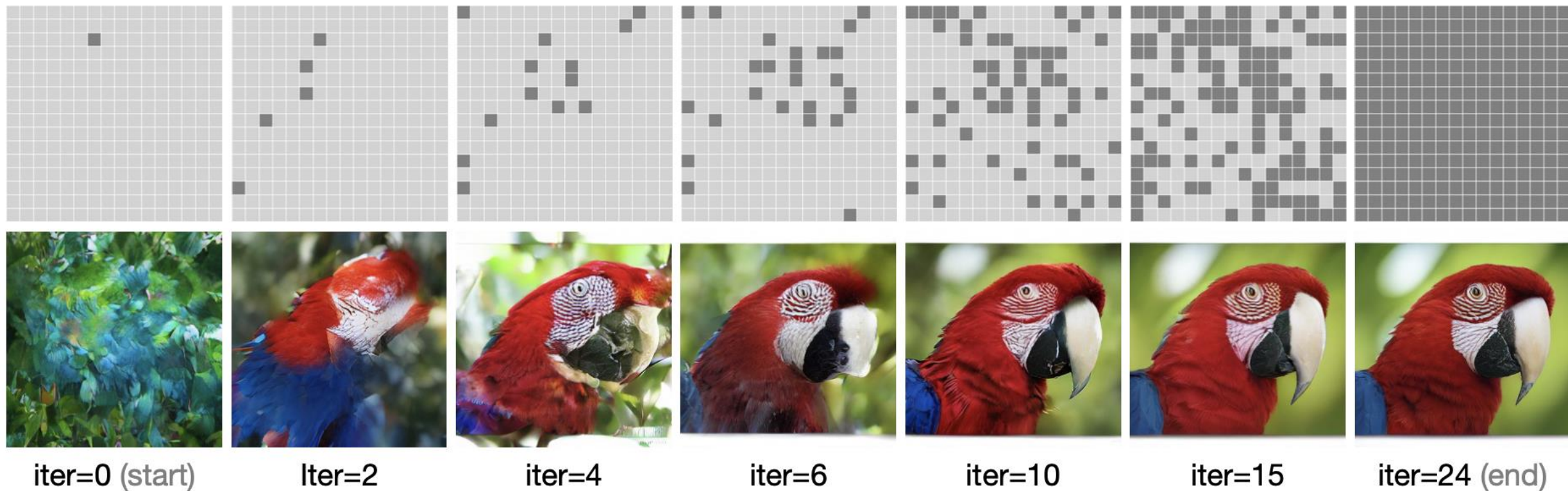
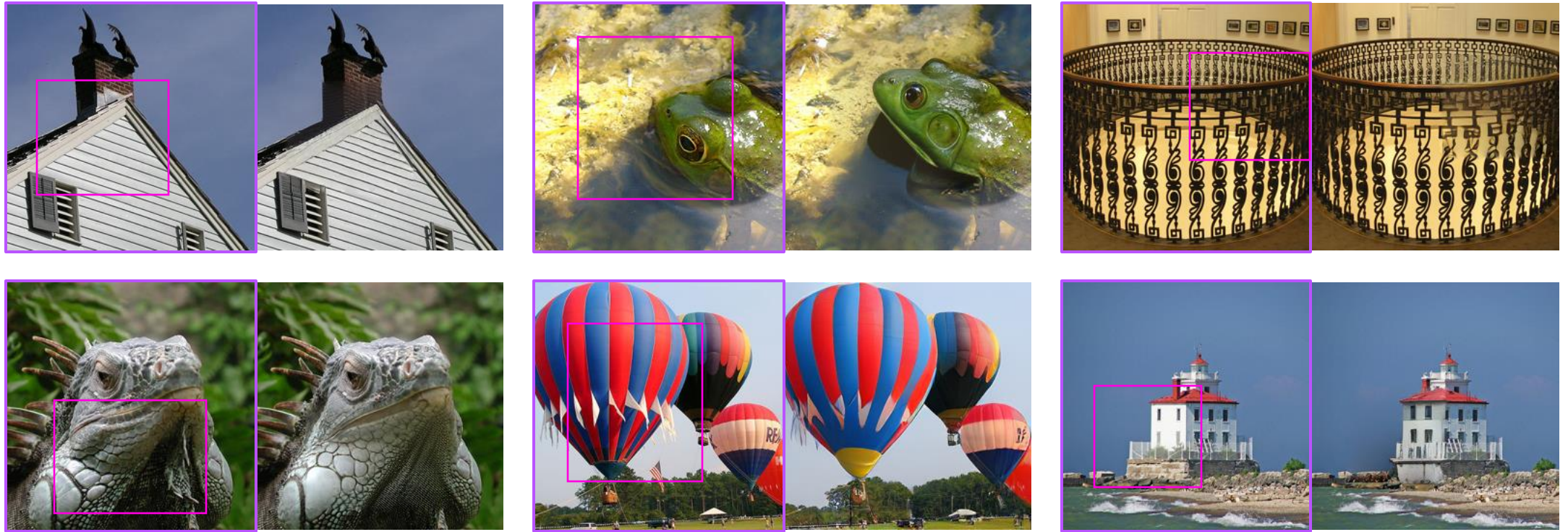


Figure 4: **Evaluation of generative and discriminative performance across different model sizes.** We report results for all tested tokenizers across four different dimensions of binary codes. We include the reconstruction FID (rFID) for each binary tokenizer for reference (grey points).

Visualization of the entropy-ordered sampling process



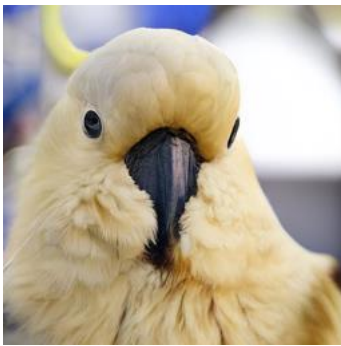
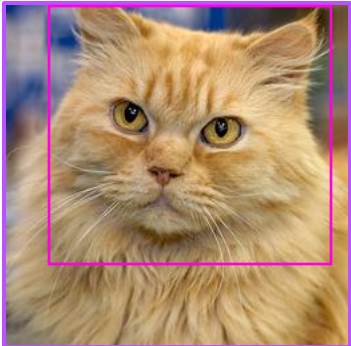
Zero-shot generalization – Inpainting



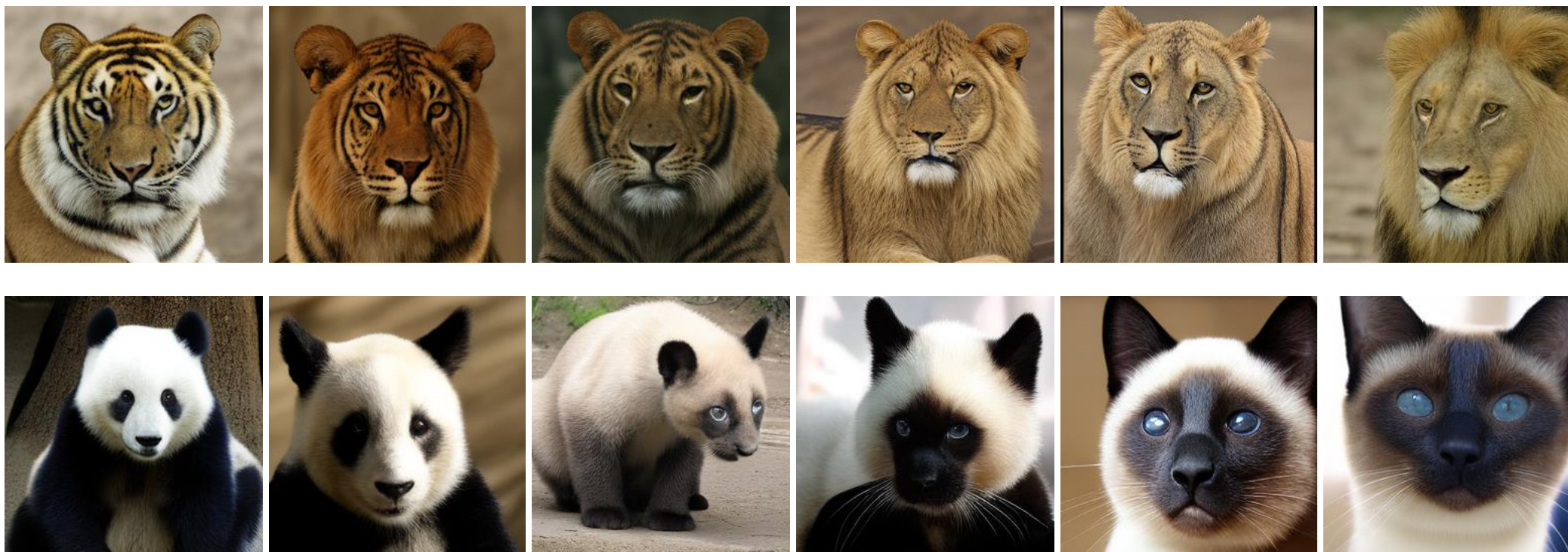
Zero-shot generalization – Outpainting



Zero-shot generalization – Editing



Zero-shot generalization – Interpolation



Text-to-image generation



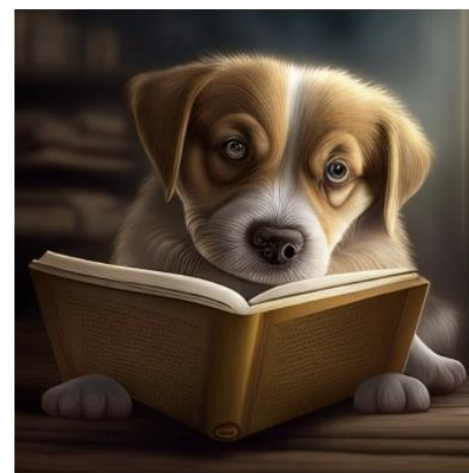
A man looks up at the starry sky, lonely.



An astronaut riding a horse.



Crocodile in a sweater.



A dog is reading a thick book.



the black hole in the space.



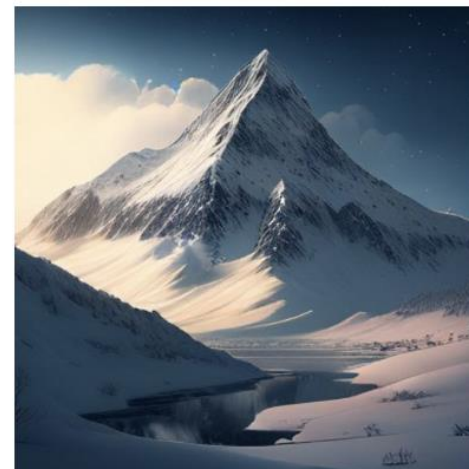
a table top with a vase of flowers on it.



an illustration of a teapot.



the Eiffel Tower in winter



A snowy mountain.



A jellyfish rocket.

Thank you for your listening!

Welcome to our Poster:

Thu 24 Apr 3 p.m. CST — 5:30 p.m. CST

Code & Model:

<https://github.com/haoosz/BiGR>

