# A. Effectiveness of TexTok on Discrete Tokens

In the main paper, we demonstrate that TexTok works well with continuous tokens. In this section, we further validate the effectiveness of TexTok with Vector-Quantized (VQ) discrete tokens. We use a codebook size of 4096. As shown in Table 6, TexTok consistently delivers significant improvements over Baseline (w/o text) in reconstruction performance. As the number of tokens decreases, the performance gains are more pronounced. These results verify the effectiveness of TexTok on discrete tokens, highlighting its versatility as a universal tokenization framework. This inherent compatibility with both continuous and discrete tokens allows TexTok to seamlessly integrate with a wide range of generative models, including diffusion models, autoregressive models, and others.

tokenizer	# tokens	$r\text{FID}\downarrow$	rIS↑	PSNR↑	SSIM $\uparrow$	LPIPS↓
Baseline-32	32	7.71	84.0	15.52	0.3822	0.4524
TexTok-32		<b>4.11</b>	<b>141.4</b>	<b>16.52</b>	<b>0.4040</b>	<b>0.3855</b>
Baseline-64	64	4.34	110.1	17.11	0.4200	0.3470
TexTok-64		<b>2.50</b>	<b>161.8</b>	<b>18.06</b>	<b>0.4462</b>	<b>0.2933</b>
Baseline-128	128	2.34	139.8	18.91	0.4737	0.2476
TexTok-128		<b>1.76</b>	<b>167.9</b>	<b>19.96</b>	<b>0.4926</b>	<b>0.2166</b>
Baseline-256	256	1.45	159.4	20.67	0.5371	0.1848
TexTok-256		<b>1.17</b>	<b>180.3</b>	<b>21.56</b>	<b>0.5526</b>	<b>0.1594</b>

Table 6. Reconstruction performance comparison of TexTok with Baseline (w/o text) using *discrete* tokens on ImageNet  $256 \times 256$ . TexTok works well with discrete tokens. It consistently outperforms Baseline (w/o text) by a large margin in image reconstruction quality, achieving more pronounced gains as the number of tokens decreases.



Figure 7. Training reconstruction FID comparison of TexTok-32 v.s. Baseline-32 (w/o text) on ImageNet  $512 \times 512$ . TexTok training is more efficient and effective, achieving faster convergence and better reconstruction quality.

#### **B.** Additional Training Analysis

In Figure 7, we further provide the training reconstruction FID comparison (evaluated on 10K samples) of TexTok-32 v.s. Baseline-32 (w/o text) on ImageNet 512×512. From

the figure, it is clear that TexTok training is more efficient and effective, achieving faster convergence and better reconstruction quality.

#### **C. Additional Implementation Details**

We provide detailed default training hyperparameters for TexTok-*N* as listed below:

- ViT encoder/decoder hidden size: 768.
- ViT encoder/decoder number of layers: 12.
- ViT encoder/decoder number of heads: 12.
- ViT encoder/decoder MLP dimensions: 3072.
- ViT patch size: 8 for 256 × 256 image resolution and 16 for 512 × 512.
- Discriminator base channels: 128.
- Discriminator channel multipliers: 1, 2, 4, 4, 4, 4 for  $256 \times 256$  image resolution, and 0.5, 1, 2, 4, 4, 4, 4 for  $512 \times 512$ .
- Discriminator starting iterations: 80,000.
- Latent shape:  $N \times 8$ .
- Reconstruction loss weight: 1.0.
- Generator loss type: Non-saturating.
- Generator adversarial loss weight: 0.1.
- Discriminator gradient penalty: r1 with cost 10.
- Perceptual loss weight: 0.1.
- LeCAM weight: 0.0001.
- Peak learning rate:  $10^{-4}$ .
- Learning rate schedule: linear warm up and cosine decay.
- Optimizer: Adam with  $\beta_1 = 0$  and  $\beta_2 = 0.99$ .
- EMA model decay rate: 0.999.
- Training epochs: 270.
- Batch size: 256.

We provide detailed default training and evaluation hyperparameters for DiT as listed below:

- Patch size: 1.
- Peak learning rate:  $5 \times 10^{-4}$ .
- Learning rate schedule: linear warm up and cosine decay.
- Optimizer: AdamW with  $\beta_1 = 0.9$ ,  $\beta_2 = 0.99$ , and 0.01 weight decay.
- Diffusion steps: 1000.
- Noise schedule: Linear.
- Diffusion  $\beta_0 = 0.0001$ ,  $\beta_{1000} = 0.02$ .
- Training objective: v-prediction.
- Sampler: DDIM.
- Sampling steps: 250.
- Training epochs: 350.
- Batch size: 1024.

For classifier-free guidance, we adopt the guidance schedule from [12] following prior arts [13, 49]. For the model configurations and other implementation details, please refer to the original DiT paper [32].



25 Words: Black and white husky, tongue out, sits snowy mountains, chained to a heavy anchor, thick fur coat.

75 Words: A medium-sized husky with a thick, predominantly white coat speckled with black patches sits in a snowy landscape against a backdrop of distant, snow-capped mountains. tongue lolls out slightly, suggesting a relaxed state, and it is tethered to a heavy anchor with a thick metal chain. The husky's fur appears soft and dense, its black markings accentuating its white body, and its expression is calm and gentle.

**25 Words:** Fluffy silver-gray Persian cat lounges relaxed on a tan fabric couch, showcasing a plush, long-haired coat and a fluffy tail.

75 Words: A fluffy, silver-and-white Persian cat

lounges comfortably on a beige, textured sofa, its long, luxurious fur creating a soft cloud-like

texture. The cat's round body and plush tail are relaxed, its paws tucked gently beneath it, and its green eyes are partially visible in a somewhat pensive expression. It appears content and at ease

in its domestic environment.





25 Words: Sunlit, orange-brown acoustic guitar rests wall near a woven fabric, showcasing smooth wood and strings.

75 Words: A sun-kissed, medium-sized acoustic guitar with a warm, sunburst finish rests against a tan stucco wall and a rustic woven fence. Its smooth stucco wall and a rustic woven fence. Its smooth, glossy body contrasts with the rough texture of the wall and the earthy tones of the fence. The guitar's dark brown fretboard and strings stand out against the vibrant orange-brown body, and a simple, elegant soundhole is centrally positioned, capturing the light. A soft fabric partially obscures the ground near the guitar.

**25 Words:** A close-up of a metallic mesh microphone head, gleaming against a pink backdrop. Its textured, silver grid contrasts with the smooth, dark mic body.

75 Words: A silver dynamic microphone dominates the frame against a soft pink backdrop. Its metallic mesh grille, a tightly woven texture reflecting light, grine, a tightly word texture reflecting light, curves over the spherical head. A dark, matte handle angles out of view, suggesting its handheld size. The close-up perspective emphasizes the microphone's readiness to capture sound.



25 Words: Snow-capped mountains frame a tranquil valley, reflecting in a still, dark-blue lake under a vibrant blue sky.

75 Words: A serene valley nestled between snow capped mountains reflects perfectly in a tranquil, Capped mountains reflects perfectly in a tranquil, deep blue lake under a vibrant azure sky. Sparse, verdant vegetation lines the valley floor, contrasting with the rugged, dark brown and gray mountain textures. Small clusters of houses dot the landscape near the water's edge, creating a picturesque scene of quiet, natural beauty.



25 Words: Brick red church with tall spire, nestled amidst green trees and a landscaped park at sunset 75 Words: A vibrant red brick church with a tall,

pointed spire stands amidst lush green trees and pullicus spire stands millus Lush green cites and manicured lawns under a soft, dusky sky. Its dark gray roof contrasts with the warm brick, and arch windows line the facade, hinting at a traditional architectural style. The church is nestled in a tranquil urban setting, with a hint of city skyline in the distance.

Figure 8. ImageNet captioning examples. Both 25-word and 75-word captions are displayed to the right of the corresponding images.

## **D. ImageNet Captioning Details**

We provide the prompt fed into Gemini that we used to generate the captions for images in ImageNet training and validation sets:

long,

Describe an image of a {class\_name} from the ImageNet dataset in a single continuous detailed paragraph, using no more than {word\_count} words. Include descriptions of the appearance, colors, textures, size, environment, and any notable features or actions. Make sure to provide a vivid and engaging description that captures the essence of the {class\_name}. Output only the description without any additional words or commentary and the description must not exceed {word\_count} words. <image\_bytes>

Figure 8 shows several captioning examples, including both 25-word and 75-word captions.

### **E.** Additional Qualitative Results

We present additional qualitative class-conditional image generation results on ImageNet  $256 \times 256$  in Figure 9. We observe that TexTok generates high-quality, semantically meaningful images with intricate fine-grained details.

We also present additional text-to-image generation results on ImageNet  $256 \times 256$  in Figure 10. Note that TexTok generates photo-realistic images that accurately align with the given prompts and even share many visual details with the reference images that the model has never seen, demonstrating both high fidelity and the capability to follow the text prompts.

### F. Additional Discussion on Related Work

There have been some recent efforts in image tokenization methods [29, 30, 48, 53] that aim to align image tokens with textual semantics. Approaches like LQAE [30], SPAE [48], and V2L [53] map images into tokens derived from the codebooks of large language models (LLMs), while LG-VO [29] focuses on aligning the decoder features of image tokens with text representations. These methods are designed primarily for bridging visual and textual modalities to improve multi-modal understanding tasks. However, by aligning image tokens directly to textual semantic spaces, they often suffer from limited image reconstruction quality due to the inherent divergence between vision and language representations. Consequently, these approaches fail to achieve reasonable performance, or even do not report results, on standard image generation benchmarks, such as ImageNet class-conditional generation.

In contrast, our work introduces a novel image tokenization framework specifically designed for generative tasks. We are the first work that proposes conditioning the tokenization process directly on the text captions of images, a strategy typically reserved for the generation phase. Rather than enforcing a strict alignment between image tokens and text captions, our method leverages descriptive captions to provide a compact semantic representation. This simplifies semantic learning, allowing more learning capacity and token space to be allocated to capture fine-grained visual details. This complementary approach significantly improves the reconstruction quality and compression rate of tokenization. Furthermore, we demonstrate that our approach achieves state-of-the-art performance on standard ImageNet conditional generation benchmarks on both  $256 \times 256$  and  $512 \times 512$  image resolutions, establishing our method as a distinct advancement over existing works.



Figure 9. **Qualitative class-conditional image generation results** on ImageNet  $256 \times 256$ . TexTok generates high-quality, semantically meaningful images with intricate fine-grained details. Results are generated with our class-conditional TexTok-256 + DiT-XL model.

#### **Reference Image**

Generated Image

**Reference Image** 

Generated Image

**Reference Image** 

Generated Image







A cheerful Pembroke Welsh Corgi, with a short, rust-colored and white coat, sits on a gray concrete step, possibly outdoors, with a silver chain leash attached to its collar. Its short legs and fluffy, erect ears are characteristic of the breed, and it displays a wide, happy grin, revealing a bright pink tongue. Its expression and compact size add to its endearing nature, captured against a neutral backdrop.



A vibrant pink flamingo stands gracefully on one leg in a dark, shallow pool of water, its soft, feathery plumage catching the light. Its long, slender pink legs and neck are elegantly curved as it preens its feathers. The bird is situated near a dark, rocky shore and a lush green palm frond, creating a striking contrast against its delicate, more here rosv hue.



A fork twirls strands of golden spaghetti coated in a creamy, eggy sauce, interspersed with bits of savory, pale pink ham and flecks of fresh parsley. The dish sits on a plain white plate, creating a visually appealing contrast, and exudes a rich, savory aroma. The spaghetti is long and thin, creating a luscious, glistening texture.



A majestic, snow-capped stratovolcano dominates the A majestic, snow-capped stratovolcano dominates the horizon, its dark, textured slopes contrasting against the bright white summit. It stands tall over a tranquil, deep blue lake, with a dark silhouette of a park bench and a few buoys in the foreground, adding a peaceful, serene touch to the scene under a clear, light blue sky.



In a small, dark-colored dish sits a scoop of creamy, off-In a small, dark-colored dish sits a scoop of creamy, of white vanilla ice cream, alongside two more, sprinkled with dark chocolate shavings. A vibrant red raspberry, drizzled with dark chocolate, crowns the dessert, restin atop the scoops. The ice cream appears smooth and cold, presented on a white plate with bold black and dark red . ting design elements, a silver spoon nearby.



A small, reddish-brown red panda, with a fluffy, dense coat, rests peacefully on a lichen-covered tree bra amidst a lush green forest. Its face is a striking branch blend of creamy white and rust-red, with dark, inquisitive eyes closed in slumber. The rich, dark brown of its back and legs contrasts with its lighter, russet-toned fur, creating a visually stunning image of tranquility nestled within the vibrant natural world.



A grand, ornate castle, primarily light blue and white with gold accents, stands tall against a twilight sky ablaze with celebratory fireworks. Its numerous towers and turrets are adorned with blue and gold banners, and intricate detailing covers its stone facade. Surrounded by dark green foliage, the castle's majestic size and enchanting glow create a magical scene in the evening atmosphere.



stands on a rough, dark gray concrete lakeshore, watching the calm, blue-gray water ripple gently. A smaller duck swins in the distance, and the scene is bathed in soft, muted light, creating a tranquil and serene lakeside environment



A medium-sized Golden Retriever with a luxurious, golden-A medium-sized volume Refriever with a lukarios, goluen-blonde coat and a broad, friendly smile sits in a vibrant green grassy yard speckled with tiny white flowers. The dog wears a red and gray athletic jersey, suggesting playful energy. Its soft, fluffy fur contrasts against the bright colors of the jersey and the lush green lawn, creating a cheerful and heartwarming scene.



A sleek, dark gray and black laptop with a slightly muted metallic sheen rests on a plain white backdrop. Its screen displays a rich, dark brown abstract image, possibly a screensaver or desktop background. The keyboard is dark gray and the touchpad is integrated into the palm rest. It's a standard size laptop with a closed lid, and appears to be in excellent condition, ready for use.

A medium-sized, white ceramic mug with a prominent Starbucks logo features a candy cane-striped handle in vibrant red and cream. Inside, a creamy liquid holds a bright red heart-shaped object, and the background is a muted gray with subtle snowflake designs. The cup's surface is smooth, and the handle's stripes are glossy, creating a festive and cozy ambiance.

Figure 10. Qualitative text-to-image generation results on ImageNet 256×256. For each sample, we present the reference image from the ImageNet validation set alongside the corresponding generated image, which is produced conditioned on the caption (displayed below) of the reference image. Results are generated with our text-to-image TexTok-128 + DiT-XL-T2I model. TexTok generates visually realistic images that accurately align with the given prompts, demonstrating both high fidelity and semantic relevance.