# Zero-Shot Text-to-Image Generation(Dall-e)

Aditya Ramesh 1 Mikhail Pavlov 1 Gabriel Goh 1 Scott Gray 1 Chelsea Voss 1 Alec Radford 1 Mark Chen 1 Ilya Sutskever 1

OpenAI

July, 1, 2021
구재원

# Introduction

**>> DALL-E: Zero-shot Image Generation by text prompt**

- Input: Text tokens + image tokens (250 millions of text-image pairs)
- Output: image tokens
- Models: VQ-VAE + Transformer

https://openai.com/blog/dall-e/#fn2
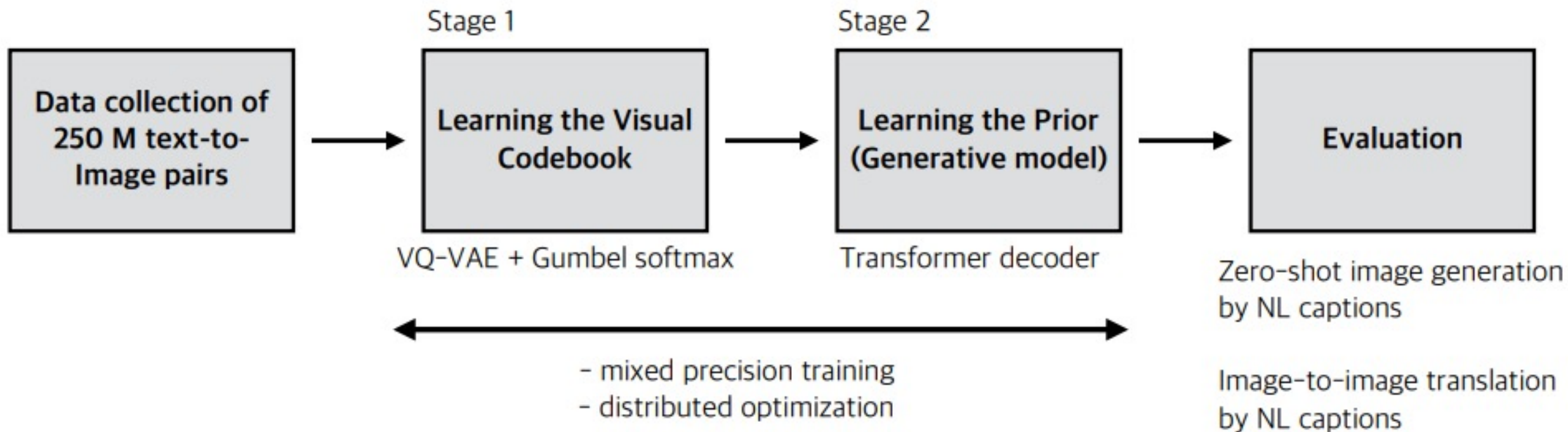
# Large-scale Generative Models

>> Problems: Still ==suffer from severe artifacts== such as object distortion, illogical object placement, or unnatural blending of foreground and background elements

>> Recent advances of ==autoregressive transformers== have achieved impressive results in several domains by large-scale of compute, model, and data

>> Text-to-image generation has typically been evaluated on ==relatively small datasets== such as MS-COCO and CUB-200

Could dataset size and model size be the limiting factor of current approaches?

# Contributions

>> Training a 12-billion params of an autoregressive transformer on 250 M image-text pairs

>> Flexible & high fidelity generative model of images controllable by natural language.

>> Zero-shot image generation on MS-COCO dataset.

>> 90 % of people prefer the DALL-E's images than those of prior work.

>> Image-to-image translation

# Overall Method



Stage 1
Stage 2

| Data collection of 250 M text-to-Image pairs | Learning the Visual Codebook | Learning the Prior (Generative model) | Evaluation |

VQ-VAE + Gumbel softmax

Transformer decoder

Zero-shot image generation by NL captions

Image-to-image translation by NL captions

- mixed precision training
- distributed optimization

# Training Objective

>> overall procedure can be views and <mark>maximizing the evidence of lower bound</mark> by the joint distribution of an Images x, captions y, and the image tokens z.

$$p_{\theta,\psi}(x, y, z) = p_\theta(x \mid y, z)p_\psi(y, z)$$

**VQ-VAE decoder**          **Transformer decoder**

$$\ln p_{\theta,\psi}(x, y) \geqslant \mathbb{E}_{z \sim q_\phi(z \mid x)} \Big( \ln p_\theta(x \mid y, z) -$$

$$\beta\, D_{\mathrm{KL}}(q_\phi(y, z \mid x), p_\psi(y, z))\Big),$$

# Step1: Learning the Visual Codebook

>> using pixels directly as image tokens would require an <mark>inordinate amount of memory</mark> for high-resolution images

>> low frequency < high frequency

>> <mark>Downsampling</mark> of an input image is required for the scalable feasibility
 - ex) 256 x 256 -> 32 x 32

>> VQ-VAE is used to downsample

# Step1: Learning the Visual Codebook

>> VQ-VAE : The <mark>code book</mark> is updated to contain latent representations for the reconstruction of training samples.
>> <mark>gumbel softmax</mark> (Jang et al., ICLR'17) is used to backpropagate the reconstruction loss, beyond the stochastic layer with categorical random variables.
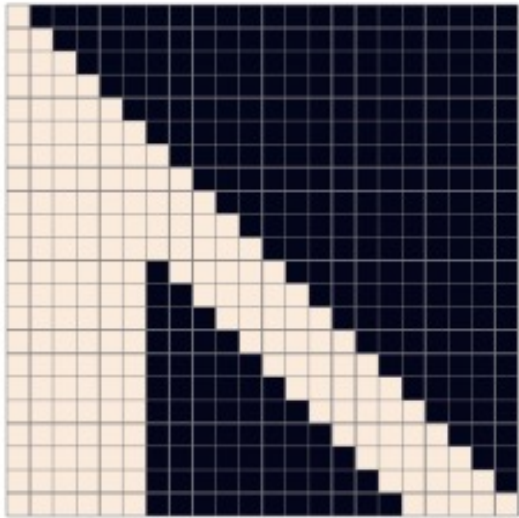
# Overall Method



Stage 1

| Data collection of 250 M text-to-Image pairs | Learning the Visual Codebook | Learning the Prior (Generative model) | Evaluation |

Stage 2

Data collection of 250 M text-to-Image pairs

Learning the Visual Codebook
VQ-VAE + Gumbel softmax

Learning the Prior (Generative model)
Transformer decoder

Evaluation
Zero-shot image generation by NL captions

Image-to-image translation by NL captions

- mixed precision training
- distributed optimization

# Step2: Learning the Prior

>> Text: BPE-encode using at most 256 tokens / Image: encode 32x32=1024 tokens
>> The text and image tokens are ==concatenated== and used for the next token prediction.

| start of text | text embed 0 | text embed 1 | text embed 2 | pad embd 0 | pad embd 1 | start of image | image embd 0 | image embd 0 |
|---|---|---|---|---|---|---|---|---|
| text pos embd 0 | text pos embd 1 | text pos embd 2 | text pos embd 3 | | | row embd 0 | row embd 0 | row embd 0 |
| | | | | | | col embd 0 | col embd 1 | col embd 2 |

+

| state 0 | state 1 | state 2 | state 3 | state 4 | state 5 | state 6 | state 7 | state 8 |
|---|---|---|---|---|---|---|---|---|

# Step2: Learning the Prior

>> model uses three kinds of sparse attention masks



(a) Row attention mask.  (b) Column attention mask.  (c) Column attention mask with transposed image states.  (d) Convolutional attention mask.

# Step2: Learning the Prior

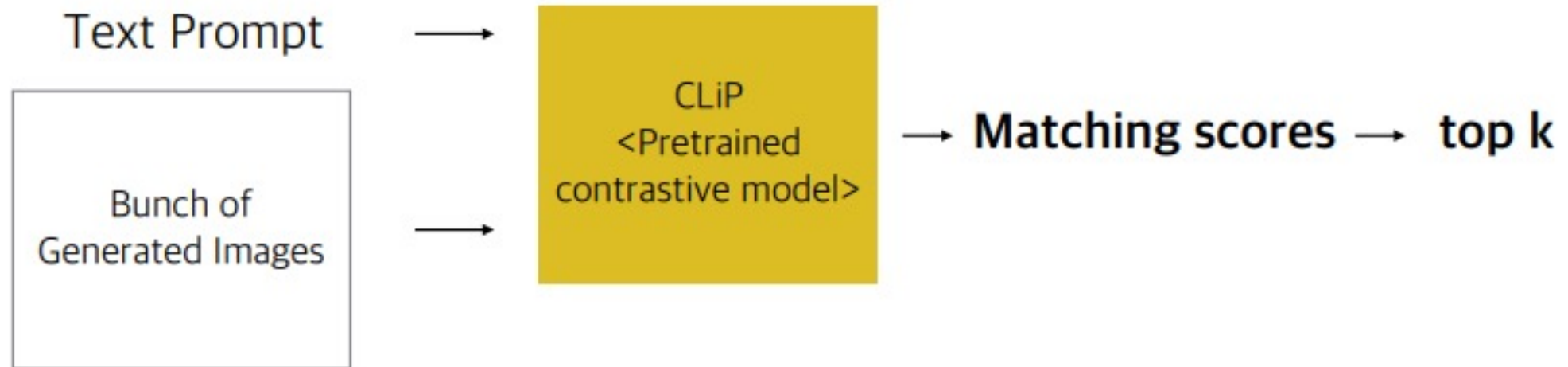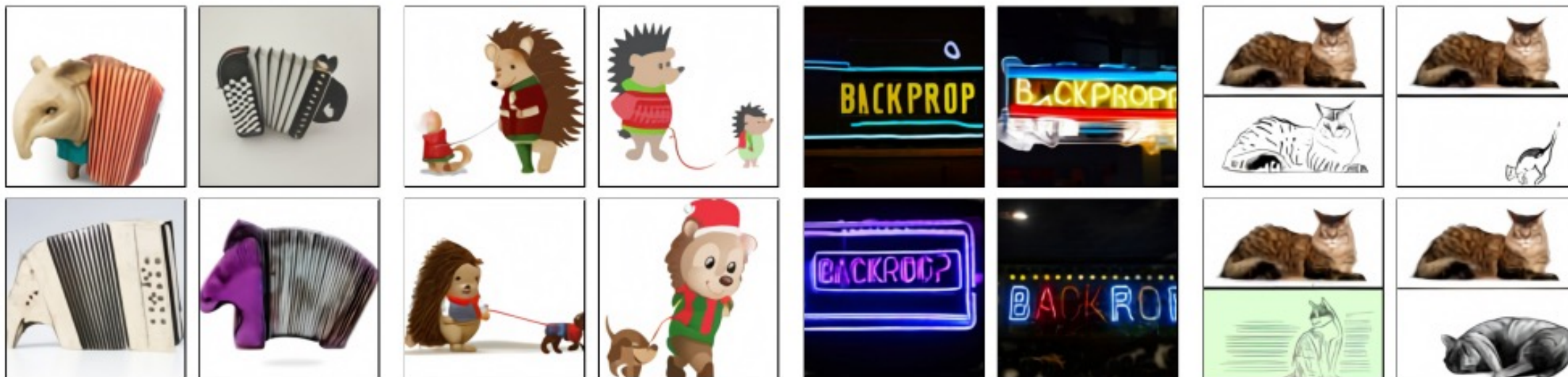>> The transformer decoder is trained to predict next tokens



(a) Autoregressive

Target

# Overall Method

[출처] https://drive.google.com/file/d/1QQH1Dyg-r_wlmZHw17sdQPunthDxg6MS/view

# Results : Sample Generation

>> At the test time, an image is generated by the text prompts with/without the part of given images

>> The samples drawn from the transformer are reranked by a retrained contrastive model (CLIP).

# Results



(a) a tapir made of accordion. a tapir with the texture of an accordion.

(b) an illustration of a baby hedgehog in a christmas sweater walking a dog

(c) a neon sign that reads "backprop". a neon sign that reads "backprop". backprop neon sign

(d) the exact same cat on the top as a sketch on the bottom

*Figure 2.* With varying degrees of reliability, our model appears to be able to combine distinct concepts in plausible ways, create anthropomorphized versions of animals, render text, and perform some types of image-to-image translation.

# Results



(a) "the exact same cat on the top as a sketch on the bottom"

(b) "the exact same photo on the top reflected upside-down on the bottom"

(c) "2 panel image of the exact same cat. on the top, a photo of the cat. on the bottom, an extreme close-up view of the cat in the photo."

(d) "the exact same cat on the top colored red on the bottom"

(e) "2 panel image of the exact same cat. on the top, a photo of the cat. on the bottom, the cat with sunglasses."

(f) "the exact same cat on the top as a postage stamp on the bottom"

Figure 14. Further examples of zero-shot image-to-image translation.
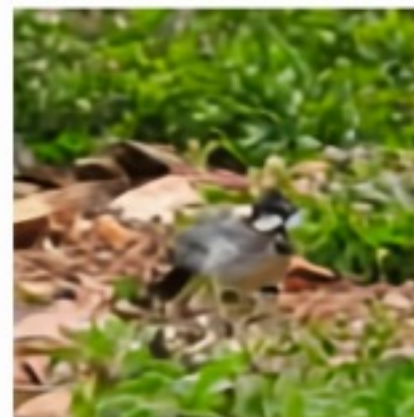
# Results



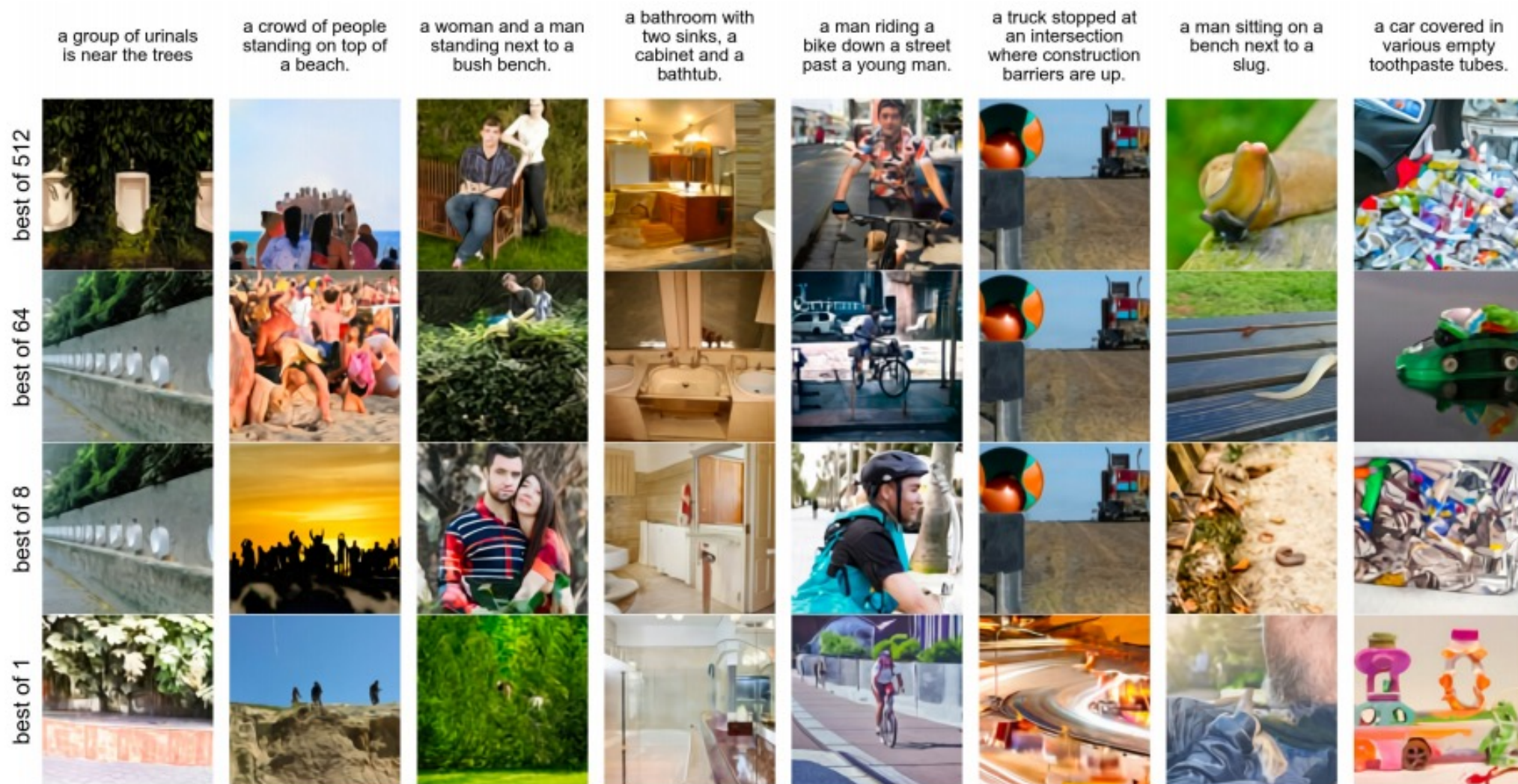*Figure 8.* Zero-shot samples from our model on the CUB dataset.

# Results



*Figure 6.* Effect of increasing the number of images for the contrastive reranking procedure on MS-COCO captions.