

---

# DALLE 2

---

A Preprint

## 1 Main

### 1.1 General

Training dataset consists of pairs  $(x; y)$  where  $x$  is image and  $y$  is a corresponding captions. Given an image  $x$ ,  $z_i, z_t$  - CLIP image and text embeddings respectively. Authors suggest method of generating image from captions using 2 stages:

- a prior  $P(z_i|y)$  that produce CLIP image embeddings  $z_i$  conditioned on captions  $y$
- A decoder  $P(x|z_i, y)$  that produces images  $x$  conditioned on CLIP image embeddings  $z_i$  (optionally captions  $y$ )

The decoder allows us to invert images given their CLIP image embeddings, while the prior allows us to learn a generative model of the image embeddings themselves.

### 1.2 Decoder

They generate images using a diffusion model, where the generation is guided by a CLIP image embedding — a 512-dimensional vector that represents either a real image or a text prompt (via CLIP’s joint embedding space). Sometimes, they also use the original text as an additional input.

$$\epsilon_{\theta}(x_t, t | z_{\text{CLIP}}, \text{text})$$

where:

- $x_t$  - noisy image at timestep  $t$
- $z_{\text{CLIP}} \in \mathbb{R}^{512}$  is the CLIP embedding
- $\text{text}$  - optional caption
- $\epsilon_{\theta}$  predicts noise to remove

They take the timestep embedding (which encodes the current noise step  $t$ ) and add to it a version of the CLIP embedding that’s been linearly transformed to the same size.

$$e_t = e_t + W_{z_{\text{CLIP}}}$$

They take the CLIP embedding, project it into 4 separate “tokens” (vectors), and attach them to the sequence of text tokens generated by GLIDE’s text encoder.

Let the projected CLIP embedding be reshaped into 4 tokens:

$$\mathbf{Z}_{\text{CLIP-tokens}} = \text{reshape}(W' \mathbf{z}_{\text{CLIP}}) \in \mathbb{R}^{4 \times d}$$

Let the text encoder output be:

$$\mathbf{H} = (\mathbf{h}_1, \dots, \mathbf{h}_n) \in \mathbb{R}^{n \times d}$$

Then the full conditioning sequence is the concatenation:

$$(\mathbf{z}_1, \dots, \mathbf{z}_4, \mathbf{h}_1, \dots, \mathbf{h}_n) \in \mathbb{R}^{(n+4) \times d}$$

To prepare the model for guidance at inference time, they simulate “unconditional” cases during training. Specifically:

10% of the time, CLIP embedding is replaced with 0 (or a special vector),

50% of the time, text is removed.

This allows the model to learn both conditional and unconditional versions of the task.

They use two more diffusion models:

First upsampler: from  $64 \times 64 \rightarrow 256 \times 256$

Second upsampler: from  $256 \times 256 \rightarrow 1024 \times 1024$

Each one is trained to increase image resolution while keeping visual consistency.

They deliberately damage the low-res images during training, so the upsamplers learn to fix flaws and don’t overfit.

First upsampler: uses Gaussian blur

Second upsampler: uses BSR degradation (a more realistic corruption pipeline: blur, noise, resize, etc.)

To save compute, they don’t train the upsampler on full-size  $1024 \times 1024$  images. Instead, they crop out smaller parts (e.g.,  $256 \times 256$ ) and train on those.

Why it works: The model still learns how to upsample locally, and generalizes to full resolution during inference.

### 1.3 Prior

---

Algorithm 1 DALL·E 2 Decoder Pipeline: From CLIP Embedding to 1024×1024 Image

---

Require: CLIP embedding  $\mathbf{z}_{\text{CLIP}} \in \mathbb{R}^{512}$ , optional text caption

Ensure: Final image  $x_0^{\text{hi-res}} \in \mathbb{R}^{1024 \times 1024 \times 3}$

- 1: (Optional) Encode caption to token sequence:  $\mathbf{H} = (\mathbf{h}_1, \dots, \mathbf{h}_n)$
- 2: Project  $\mathbf{z}_{\text{CLIP}}$  into:
- 3: (a) Token embeddings:  $\mathbf{Z}_{\text{CLIP-tokens}} = \text{reshape}(W' \mathbf{z}_{\text{CLIP}}) \in \mathbb{R}^{4 \times d}$
- 4: (b) Timestep bias:  $\mathbf{z}_{\text{proj}} = W \mathbf{z}_{\text{CLIP}}$
- 5: Form conditioning sequence:  $\mathbf{C} = (\mathbf{z}_1, \dots, \mathbf{z}_4, \mathbf{h}_1, \dots, \mathbf{h}_n)$

Stage 1: 64×64 Base Image Generation

- 6: for  $t = T \rightarrow 1$  do
- 7:   Sample noise image  $x_t \sim \mathcal{N}(0, I)$
- 8:   Compute denoised prediction:

$$\epsilon_{\text{guided}} = (1 + w) \cdot \epsilon_{\theta}(x_t, t, \mathbf{C}) - w \cdot \epsilon_{\theta}(x_t, t, \emptyset)$$

- 9:   Sample  $x_{t-1}$  using reverse DDPM step
- 10: end for
- 11:  $x_0^{\text{low}} \leftarrow x_0$

▷ Generated 64×64 image

Stage 2: 256×256 Upsampling

- 12: Corrupt  $x_0^{\text{low}}$  with Gaussian blur:  $\tilde{x}_0^{\text{low}}$
- 13: for  $t = T \rightarrow 1$  do
- 14:   Sample noise  $x_t$
- 15:   Predict noise:  $\epsilon_{\theta}(x_t, t \mid \tilde{x}_0^{\text{low}})$
- 16:   Sample  $x_{t-1}$
- 17: end for
- 18:  $x_0^{\text{mid}} \leftarrow x_0$

▷ Generated 256×256 image

Stage 3: 1024×1024 Upsampling

- 19: Corrupt  $x_0^{\text{mid}}$  with BSR degradation:  $\tilde{x}_0^{\text{mid}}$
- 20: Train on 256×256 crops of targets to reduce compute
- 21: for  $t = T \rightarrow 1$  do
- 22:   Sample noise  $x_t$
- 23:   Predict noise:  $\epsilon_{\theta}(x_t, t \mid \tilde{x}_0^{\text{mid}})$
- 24:   Sample  $x_{t-1}$
- 25: end for
- 26:  $x_0^{\text{hi-res}} \leftarrow x_0$

▷ Generated 1024×1024 image

---

---

Algorithm 2 DALL·E 2 Prior: Generating CLIP Image Embedding from Caption

---

Require: Caption  $y$ , CLIP text encoder  $\text{CLIP}_{\text{text}}$

Ensure: Predicted image embedding  $\hat{\mathbf{z}}_i \in \mathbb{R}^{512}$

1: Compute CLIP text embedding:  $\mathbf{z}_t \leftarrow \text{CLIP}_{\text{text}}(y)$

Option 1: Autoregressive (AR) Prior

2: Compute target image embedding  $\mathbf{z}_i$

3: Apply PCA:  $\mathbf{z}_i^{\text{PCA}} \in \mathbb{R}^{319}$

4: Quantize each component into 1024 bins  $\Rightarrow$  sequence  $\{q_1, \dots, q_{319}\}$

5: Compute dot product:  $s = \mathbf{z}_i \cdot \mathbf{z}_t$

6: Quantize  $s \Rightarrow$  discrete token  $q_s$

7: Prepare prefix tokens from  $y$ ,  $\mathbf{z}_t$ , and  $q_s$

8: Train Transformer with causal mask to predict  $\{q_1, \dots, q_{319}\}$

Sampling from AR Prior

9: Encode prefix as above

10: Predict  $q_1, \dots, q_{319}$  autoregressively

11: Dequantize and apply inverse PCA:  $\hat{\mathbf{z}}_i$

Option 2: Diffusion Prior

12: Train diffusion model on  $\mathbf{z}_i$  using noising schedule  $q_t$

13: for  $t = T \rightarrow 1$  do

14: Sample noise:  $\mathbf{z}_i^{(t)} \sim \mathcal{N}(0, I)$

15: Form input sequence:

16: Text tokens from  $y$

17: CLIP text embedding  $\mathbf{z}_t$

18: Timestep embedding  $\mathbf{e}_t$

19: Noised  $\mathbf{z}_i^{(t)}$

20: Placeholder token for prediction

21: Use Transformer output to predict  $\hat{\mathbf{z}}_i^{(0)}$

22: end for

Inference with Dot Product Re-ranking

23: Generate 2 samples:  $\hat{\mathbf{z}}_i^{(1)}, \hat{\mathbf{z}}_i^{(2)}$

24: Compute  $\hat{s}_1 = \hat{\mathbf{z}}_i^{(1)} \cdot \mathbf{z}_t$ ,  $\hat{s}_2 = \hat{\mathbf{z}}_i^{(2)} \cdot \mathbf{z}_t$

25: Return the one with higher score:  $\hat{\mathbf{z}}_i = \arg \max_{\hat{\mathbf{z}}} (\hat{s}_1, \hat{s}_2)$

---