

Stable Diffusion finetune

Links

<https://dreambooth.github.io/>

<https://huggingface.co/blog/dreambooth>

<https://github.com/huggingface/diffusers/tree/main/examples/dreambooth>

<https://arxiv.org/pdf/2208.01618.pdf>

<https://arxiv.org/abs/2106.09685>

<https://stable-diffusion-art.com/lora/>

Методы дообучения модели

- Dreambooth
- Textual Inversion
- LoRA
- Hypernetworks

DreamBooth



Input images



in the Acropolis



swimming



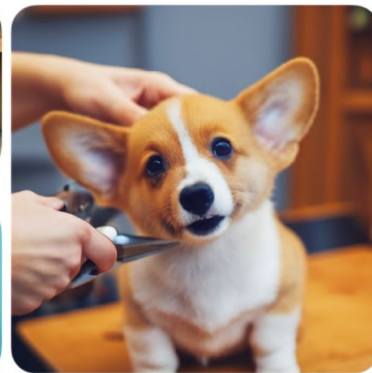
sleeping



in a doghouse



in a bucket



getting a haircut



Input images



A [V] teapot floating
in milk



A transparent [V] teapot
with milk inside

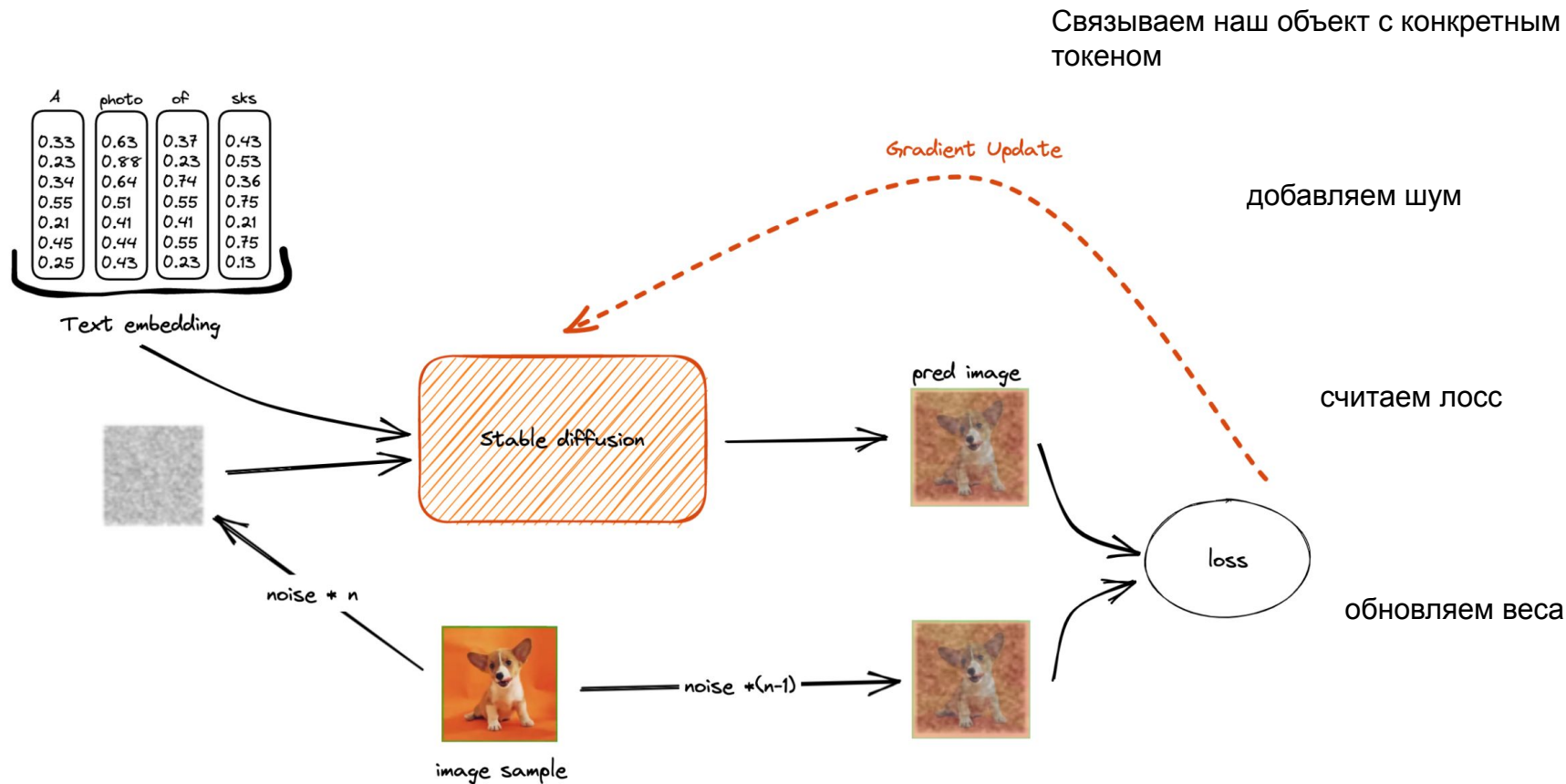


A [V] teapot
pouring tea



A [V] teapot floating
in the sea

DreamBooth



DreamBooth

$$\mathbf{x}_{\text{gen}} = \hat{\mathbf{x}}_{\theta}(\epsilon, \mathbf{c})$$

ϵ - шум; \mathbf{c} - conditional вектор; \mathbf{x}_{θ} - модель

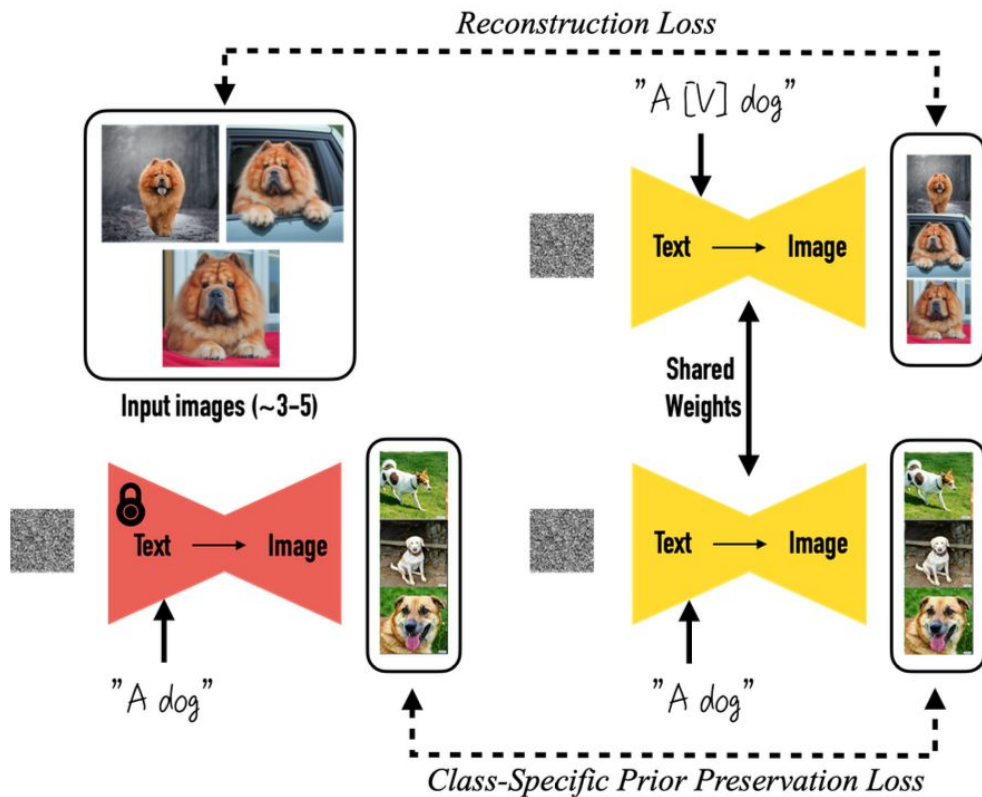
Считаем MSE loss для шума из разницы изображений

функция потерь

$$\mathbb{E}_{\mathbf{x}, \mathbf{c}, \epsilon, t} [w_t \|\hat{\mathbf{x}}_{\theta}(\alpha_t \mathbf{x} + \sigma_t \epsilon, \mathbf{c}) - \mathbf{x}\|_2^2]$$

функция потерь с регуляризацией

$$\mathbb{E}_{\mathbf{x}, \mathbf{c}, \epsilon, \epsilon', t} [w_t \|\hat{\mathbf{x}}_{\theta}(\alpha_t \mathbf{x} + \sigma_t \epsilon, \mathbf{c}) - \mathbf{x}\|_2^2 + \lambda w_{t'} \|\hat{\mathbf{x}}_{\theta}(\alpha_{t'} \mathbf{x}_{\text{pr}} + \sigma_{t'} \epsilon', \mathbf{c}_{\text{pr}}) - \mathbf{x}_{\text{pr}}\|_2^2],$$



DreamBooth

Цель: добавить в модель конкретный объект при этом не потерять накопленные знания

1. расширяем словарь и соотносим объект с непопулярным токеном[sks] чтобы не было *language drift*
2. Плохая идея выбрать случайный набор символов “хху5syt00” токенайзер может их разбить на разные токены
3. внедряем пару “a [identifier] [class noun]” (e.g. cat, dog, watch, etc.)
4. На вход достаточно 5 картинок на разном фоне и с разных ракурсов (максимально разнообразные)
5. тренируем модель целиком, приходится сохранять веса целиком

Inference



Textual inversion



Textual inversion

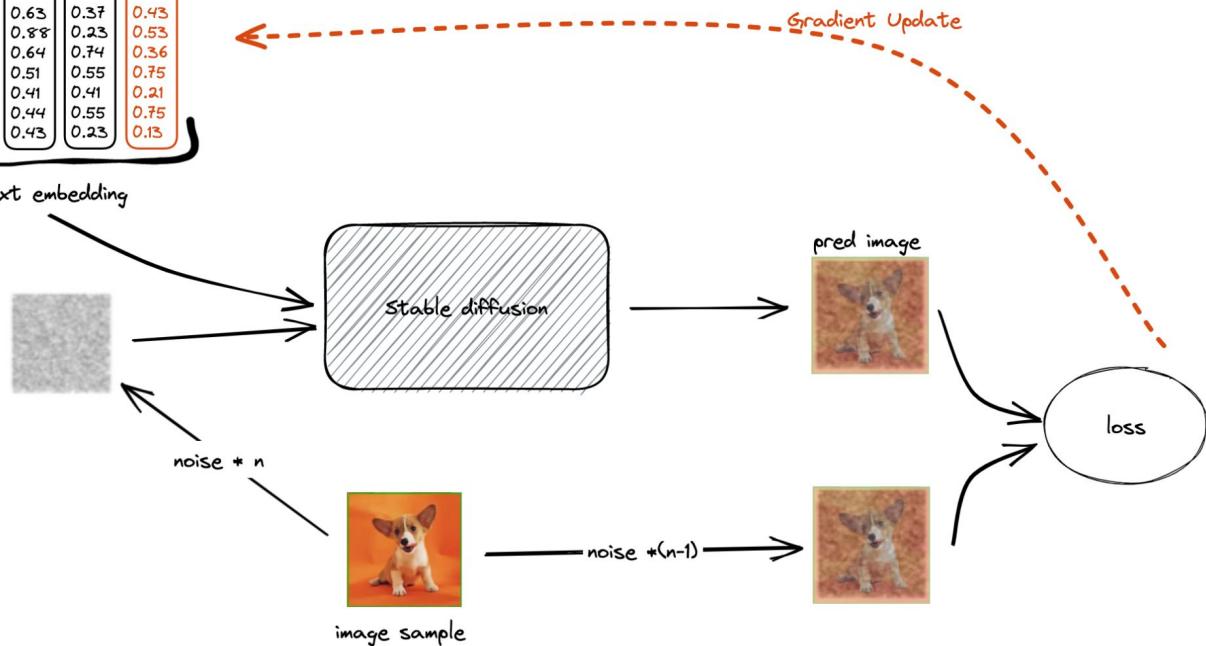
Вместо целой модели
“обучаем” наш текстовый
embedding

A	photo	of	sks
0.33	0.63	0.37	0.43
0.23	0.88	0.23	0.53
0.34	0.64	0.74	0.36
0.55	0.51	0.55	0.75
0.21	0.41	0.41	0.21
0.45	0.44	0.55	0.75
0.25	0.43	0.23	0.13

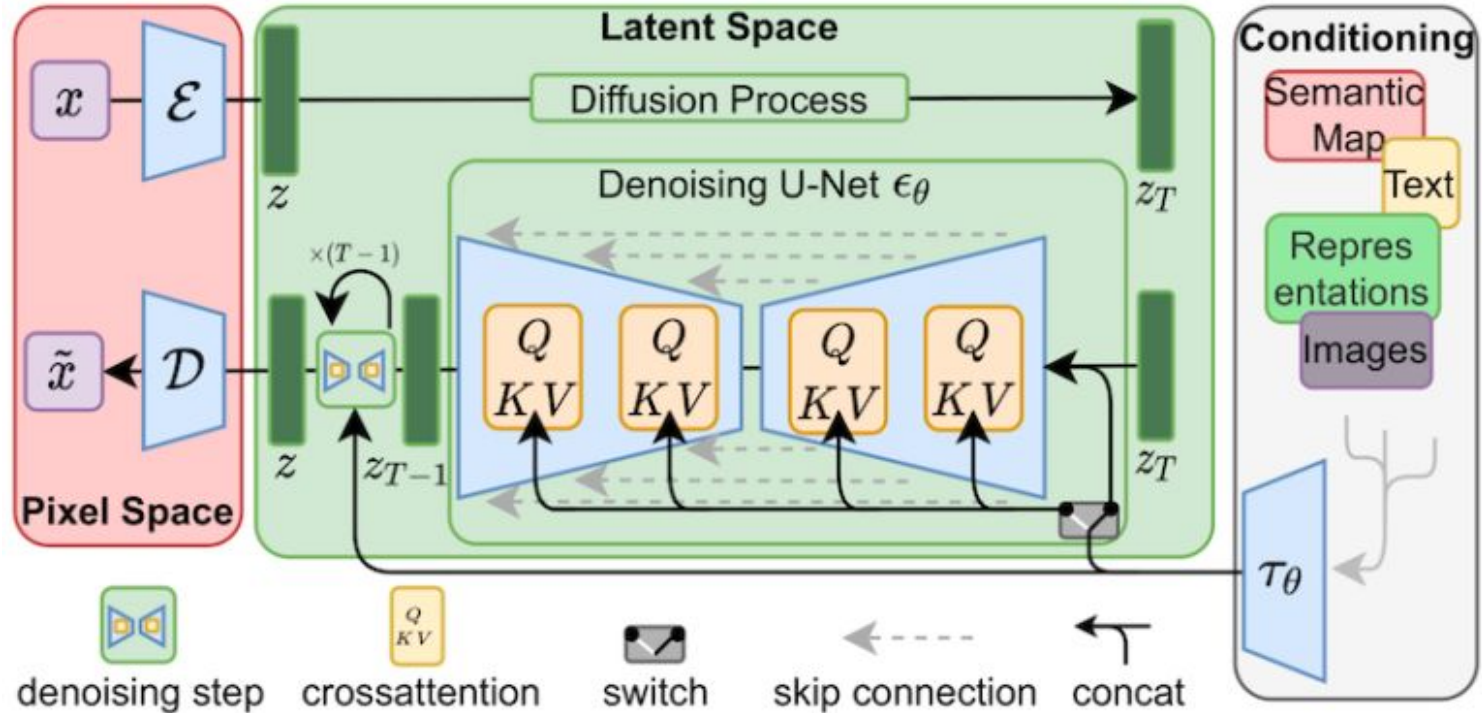
Text embedding

Шаблоны

- “a photo of a S_* .”,
- “a rendering of a S_* .”,
- “a cropped photo of the S_* .”,
- “the photo of a S_* .”,
- “a photo of a clean S_* .”,
- “a photo of a dirty S_* .”,
- “a dark photo of the S_* .”,
- “a photo of my S_* .”,



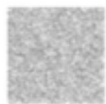
Lora Low-Rank Adaptation



Lora Low-Rank Adaptation

A	photo	of	skis
0.33	0.63	0.37	0.43
0.23	0.88	0.23	0.53
0.34	0.64	0.74	0.36
0.55	0.51	0.55	0.75
0.21	0.41	0.41	0.21
0.45	0.44	0.55	0.75
0.25	0.43	0.23	0.13

Text embedding

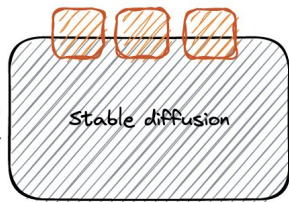


noise * n



image sample

noise *(n-1)



pred image



Метод для больших языковых моделей

экономит ресурсы и место

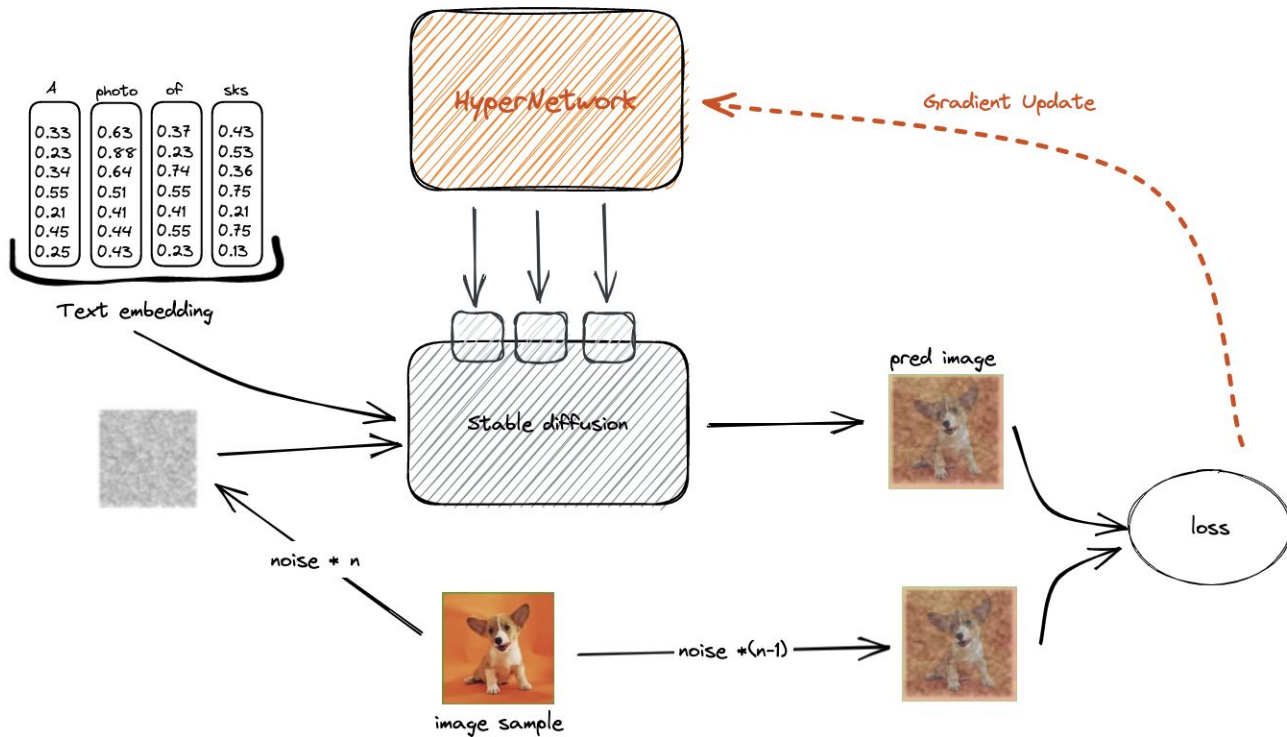


в матрице размером 1000*2000
2 000 000 параметров

в матрицах размером 1000*2 и 2000*2
6 000 параметров

$$C = A \cdot B = \begin{pmatrix} 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \end{pmatrix} \cdot \begin{pmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix} = \begin{pmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}$$

HeperNetwork





That's all Folks!