

Глубинное обучение

Из NLP в CV

Ирина Сапарина

Vision Transformer

Vision Transformer

Идея: применить архитектуру Transformer для CV (image classification)

- нужна последовательность на вход

Как из картинки сделать последовательность?

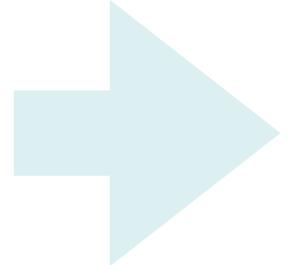


[Image credit](#)

Vision Transformer

Как из картинки сделать последовательность?

- разделим на 2D патчи



Vision Transformer

Как из картинки сделать последовательность?

- разделим на 2D патчи
- вытянем в последовательность 2D картинок

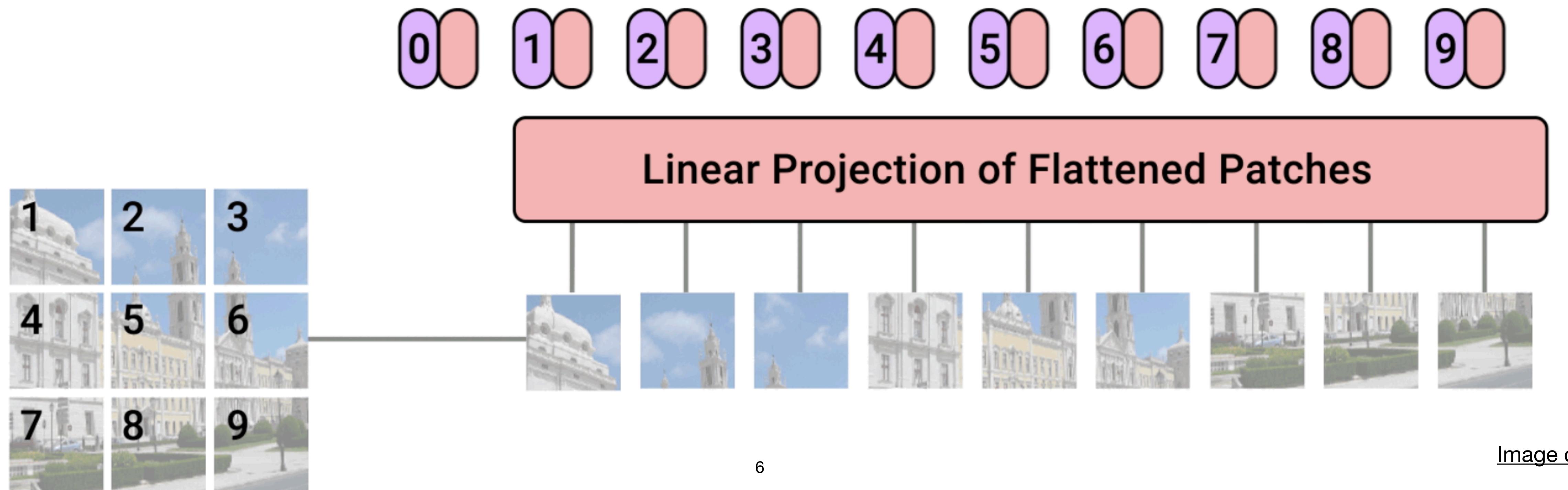


[Image credit](#)

Vision Transformer

Как из картинки сделать последовательность?

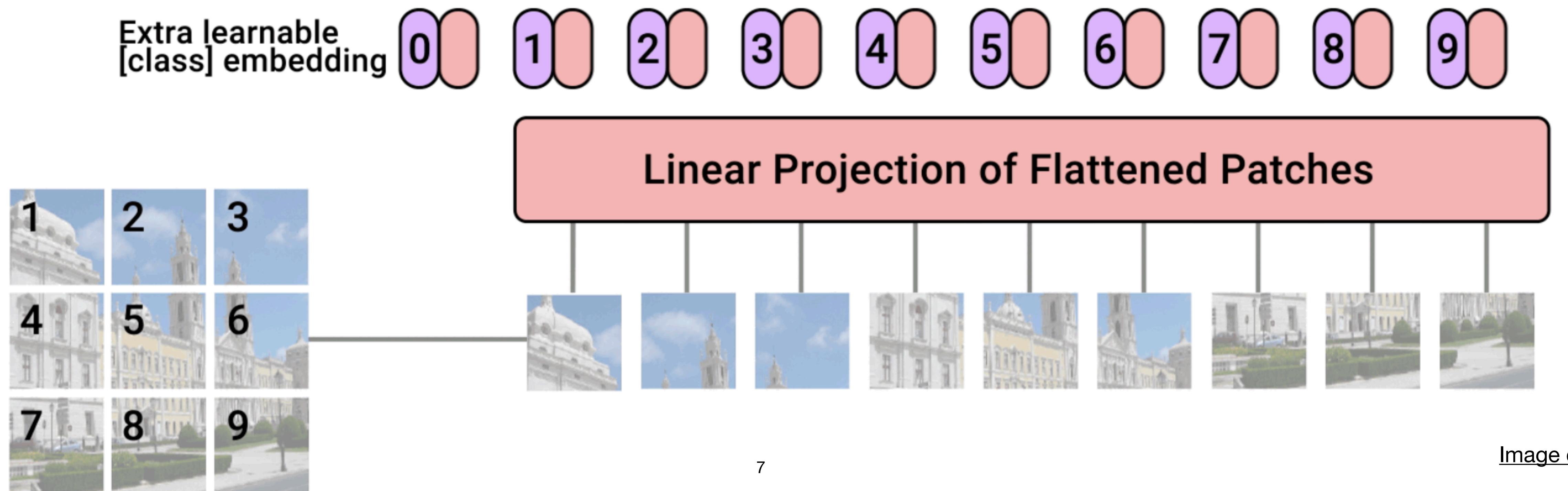
- разделим на 2D патчи
- вытянем в последовательность 2D картинок
- Linear mapping - в меньшую размерность



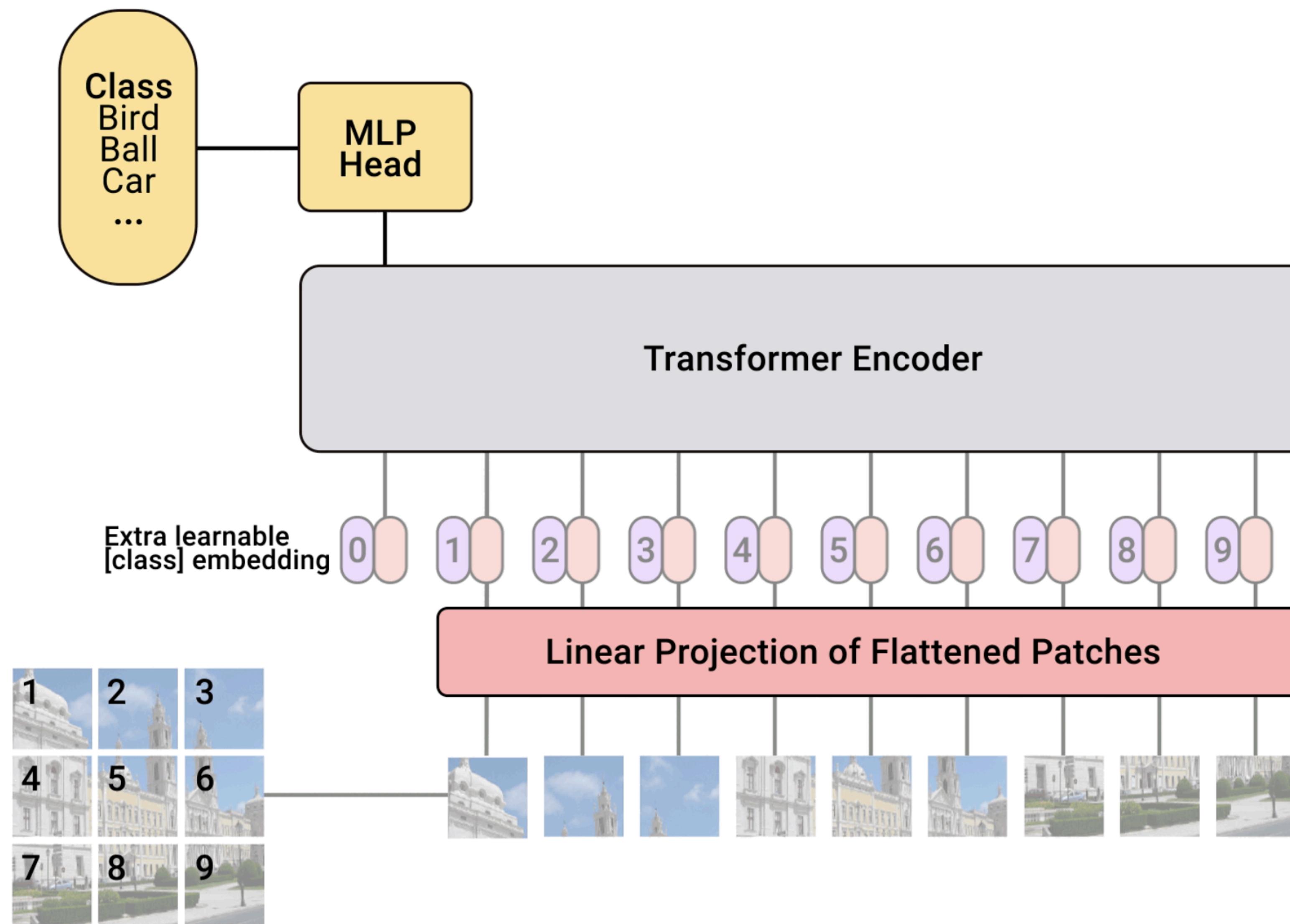
Vision Transformer

Как из картинки сделать последовательность?

- разделим на 2D патчи
- вытянем в последовательность 2D картинок
- Linear mapping - в меньшую размерность



Vision Transformer



Vision Transformer

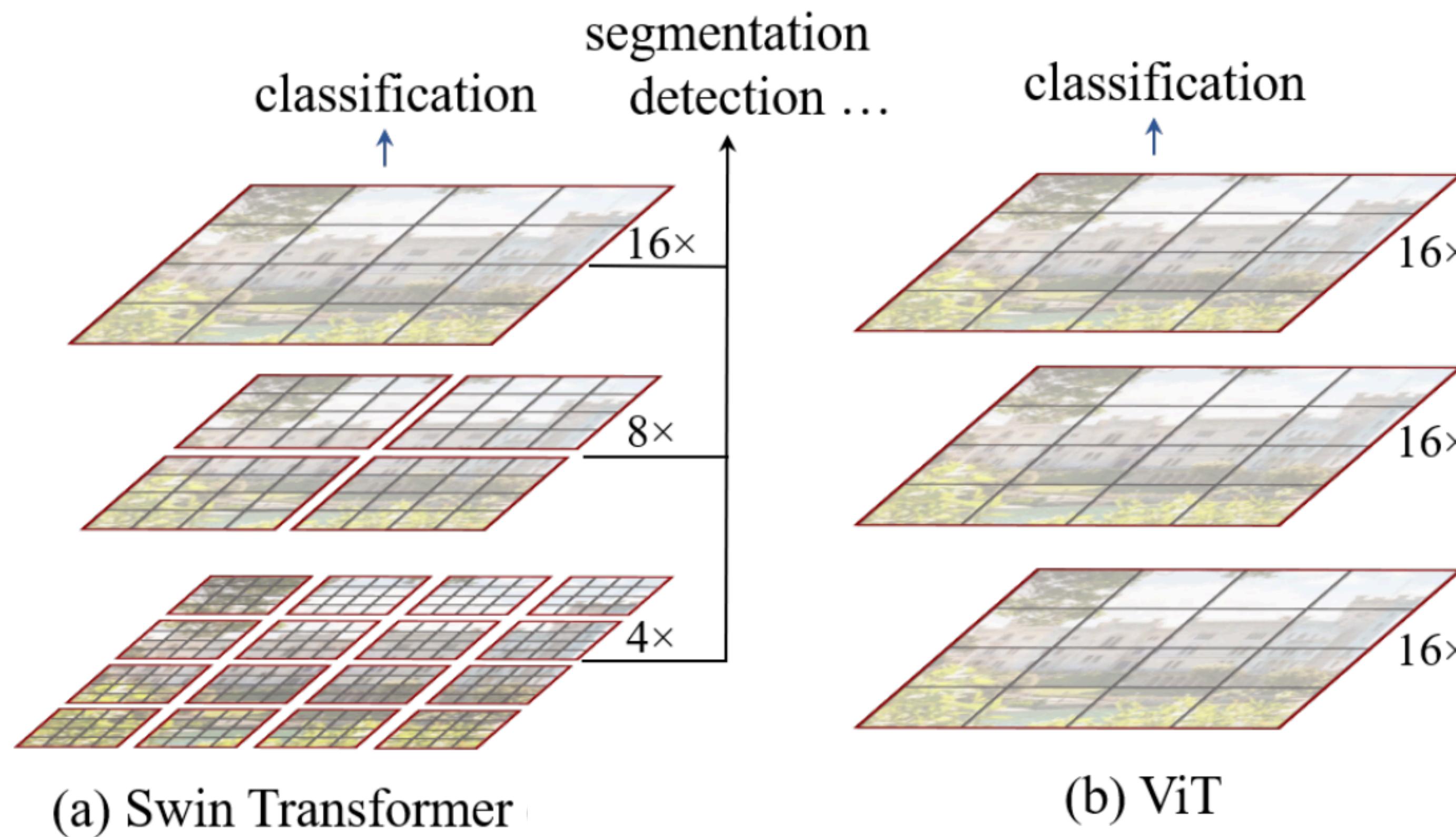
	ViT-H	Previous SOTA
ImageNet	88.55	88.5
ImageNet-ReaL	90.72	90.55
Cifar-10	99.50	99.37
Cifar-100	94.55	93.51
Pets	97.56	96.62
Flowers	99.68	99.63



[Image credit](#)

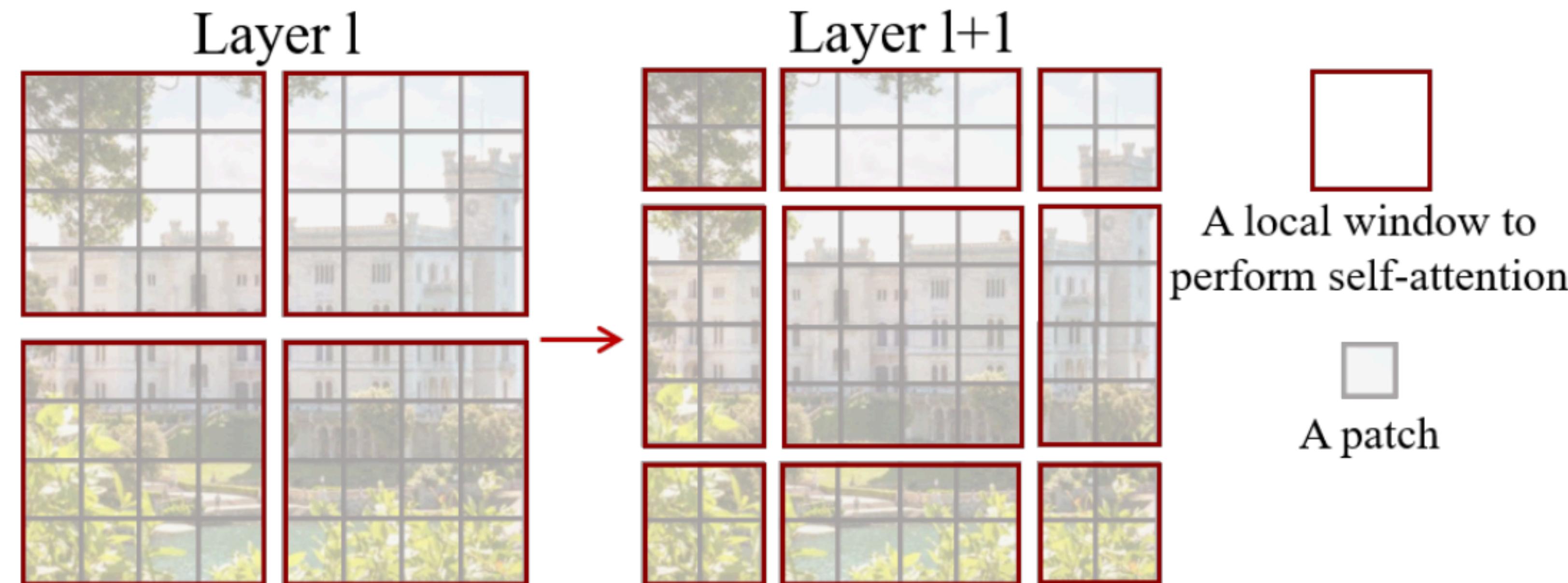
Swin Transformer

Для других задач CV (сегментация, детекторы) фиксированный патчи работают плохо - это исправляет Swin Transformer



Swin Transformer

Shifted Window: self-attention между блоками патчей,
сдвигаем окна на следующем слое



CLIP: Contrastive Language– Image Pre-training

CLIP (Contrastive Language–Image Pre-training)

BERT-like модели предобучаются **на неразмеченных данных** (тексты)

Image classifiers учат **на размеченных людьми данных** (ImageNet, etc.)

Идея: предобучить image classifier на больших данных без использования аннотации (краудсорсинга)

CLIP (Contrastive Language–Image Pre-training)

- Собрали 400 миллионов пар (картишка, подпись) - из интернета

**a train traveling down a track
next to a forest.**

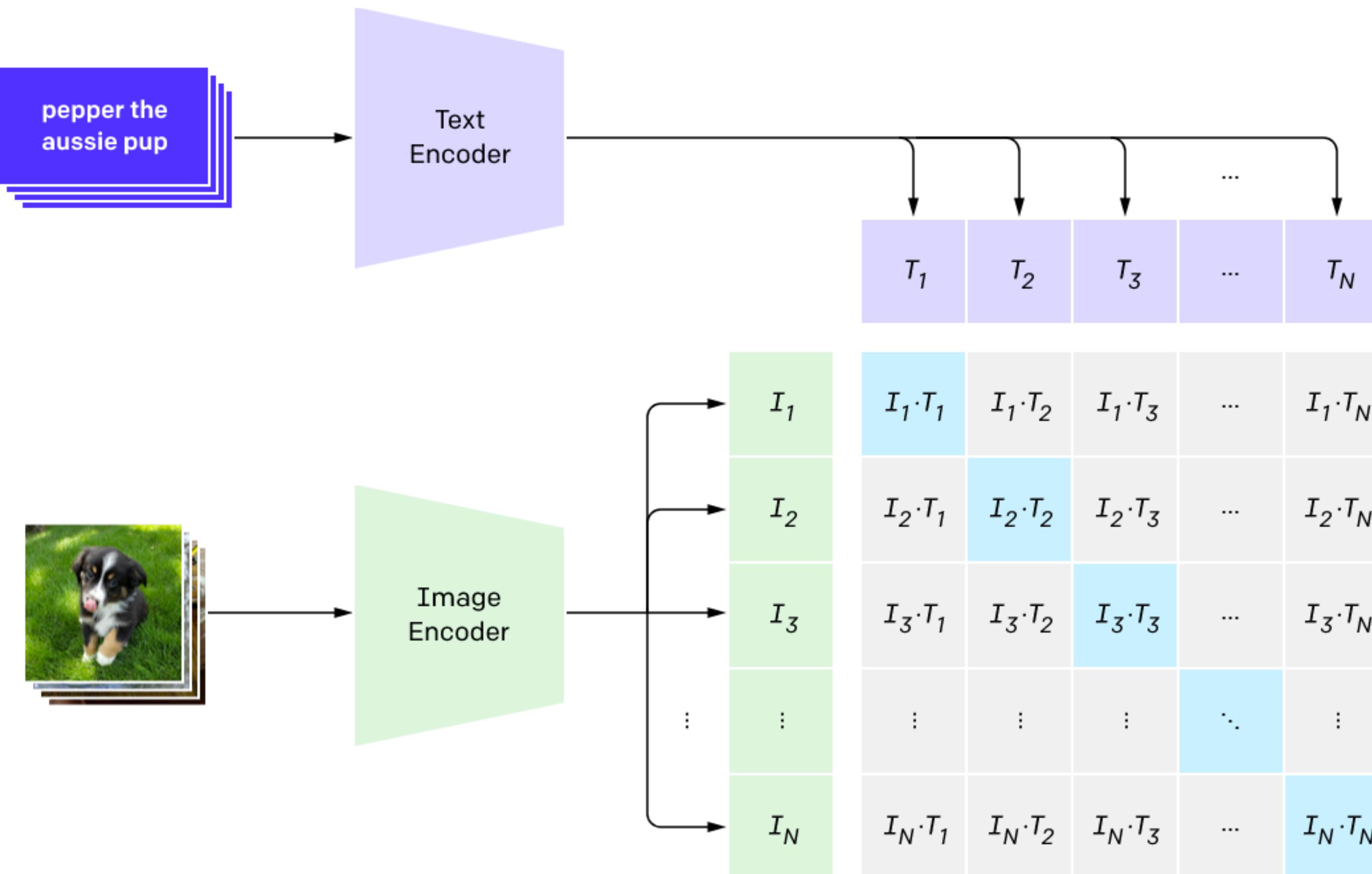


**a group of young boys playing
soccer on a field.**



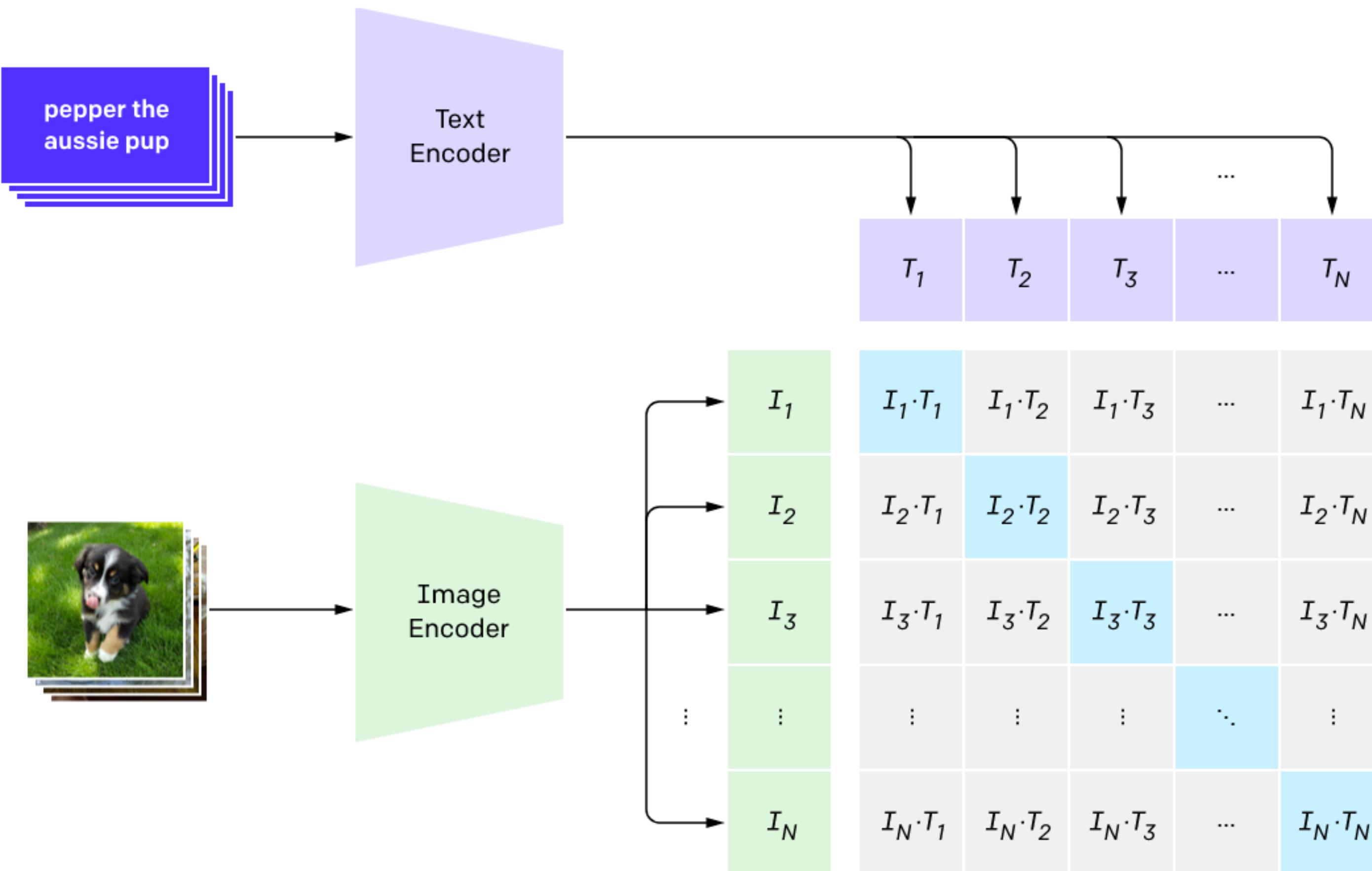
CLIP (Contrastive Language–Image Pre-training)

- Предобучение: энкодим N картинок (ViT) и N подписей к ним (Transformer)



CLIP (Contrastive Language–Image Pre-training)

- Предобучение: энкодим N картинок (ViT) и N подписей к ним (Transformer)



Лосс:

```
# symmetric loss function
labels = np.arange(n)
loss_i = cross_entropy_loss(logits, labels, axis=0)
loss_t = cross_entropy_loss(logits, labels, axis=1)
loss  = (loss_i + loss_t)/2
```

CLIP (Contrastive Language–Image Pre-training)

- Применение для Image Classification: меняем метку на текст “the photo of ...”

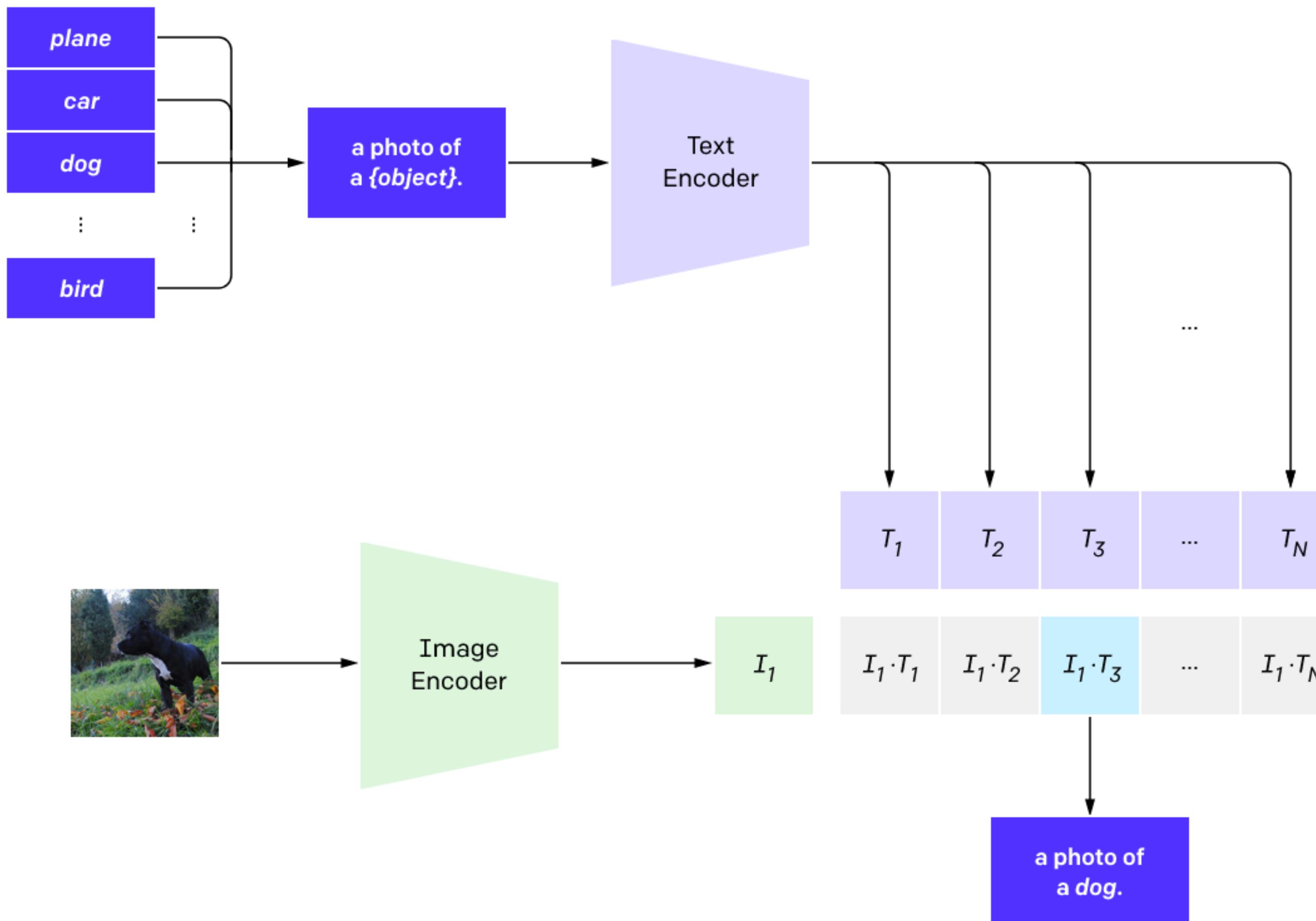


Image credit

CLIP (Contrastive Language–Image Pre-training)

FOOD101

guacamole (90.1%) Ranked 1 out of 101 labels



- ✓ a photo of **guacamole**, a type of food.
- ✗ a photo of **ceviche**, a type of food.
- ✗ a photo of **edamame**, a type of food.
- ✗ a photo of **tuna tartare**, a type of food.
- ✗ a photo of **hummus**, a type of food.

SUN397

television studio (90.2%) Ranked 1 out of 397



- ✓ a photo of a **television studio**.
- ✗ a photo of a **podium indoor**.
- ✗ a photo of a **conference room**.
- ✗ a photo of a **lecture room**.
- ✗ a photo of a **control room**.

YOUTUBE-BB

airplane, person (89.0%) Ranked 1 out of 23



- ✓ a photo of a **airplane**.
- ✗ a photo of a **bird**.
- ✗ a photo of a **bear**.
- ✗ a photo of a **giraffe**.
- ✗ a photo of a **car**.

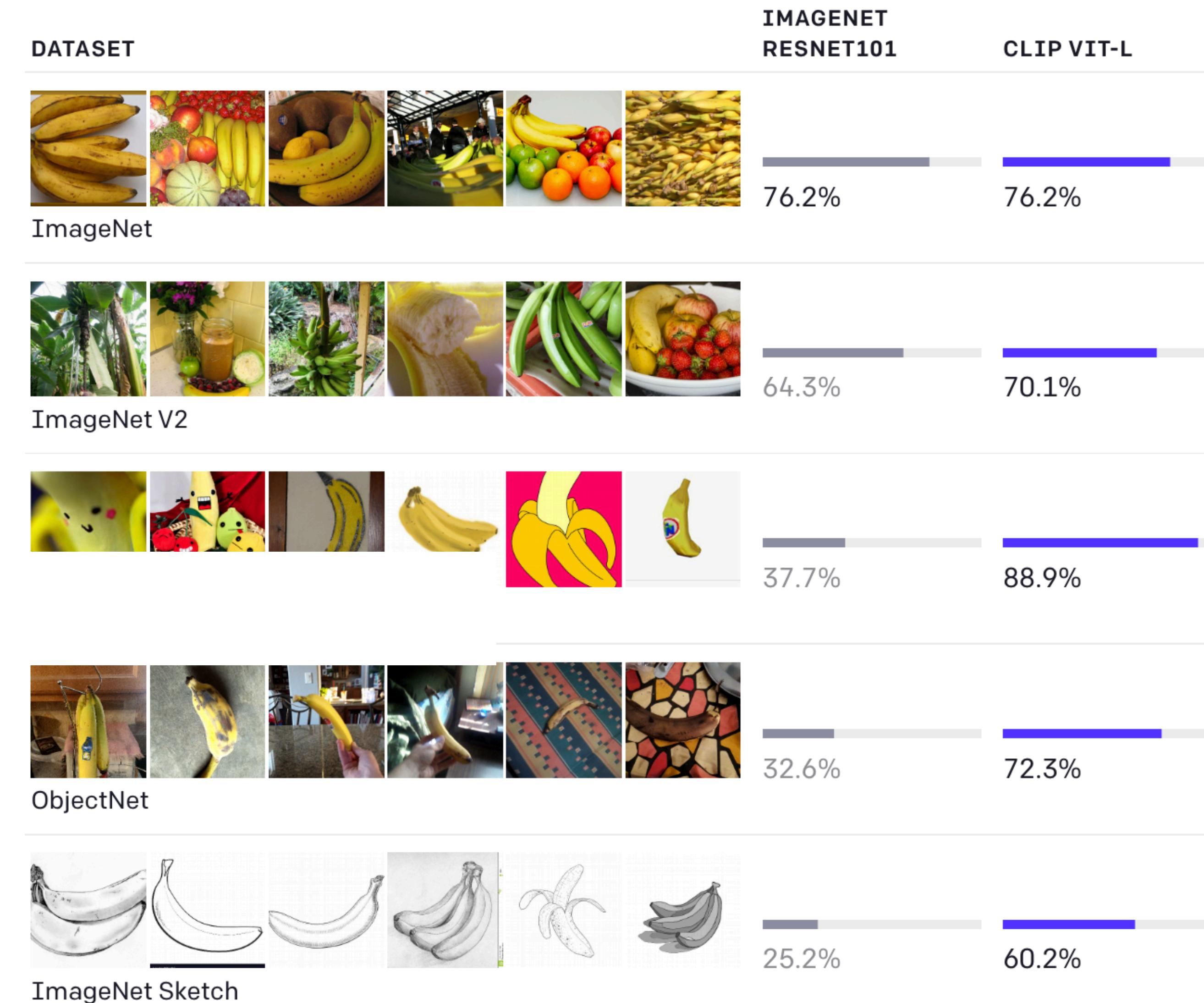
EUROSAT

annual crop land (12.9%) Ranked 4 out of 10



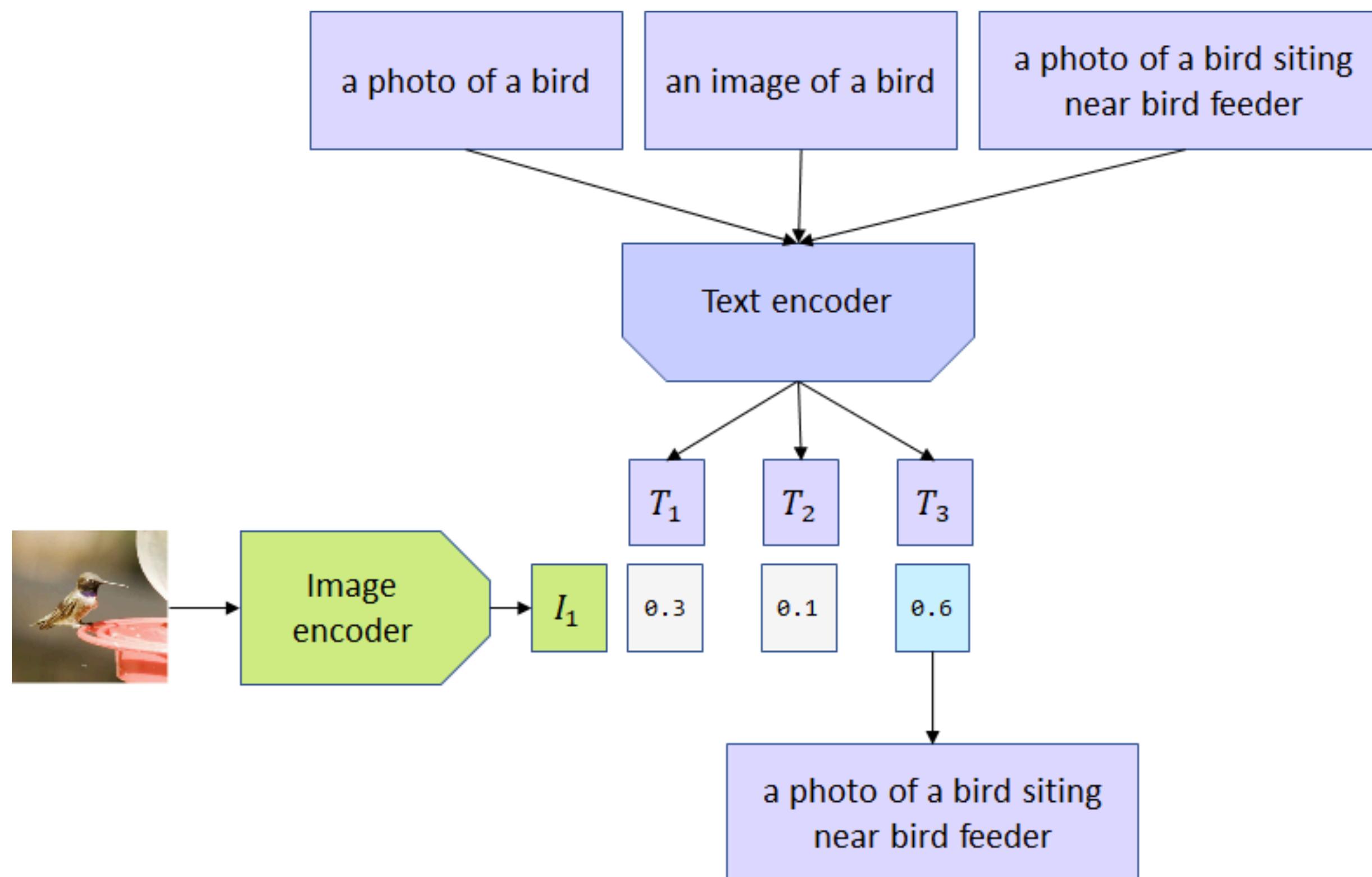
- ✗ a centered satellite photo of **permanent crop land**.
- ✗ a centered satellite photo of **pasture land**.
- ✗ a centered satellite photo of **highway or road**.
- ✓ a centered satellite photo of **annual crop land**.
- ✗ a centered satellite photo of **brushland or shrubland**.

CLIP (Contrastive Language–Image Pre-training)



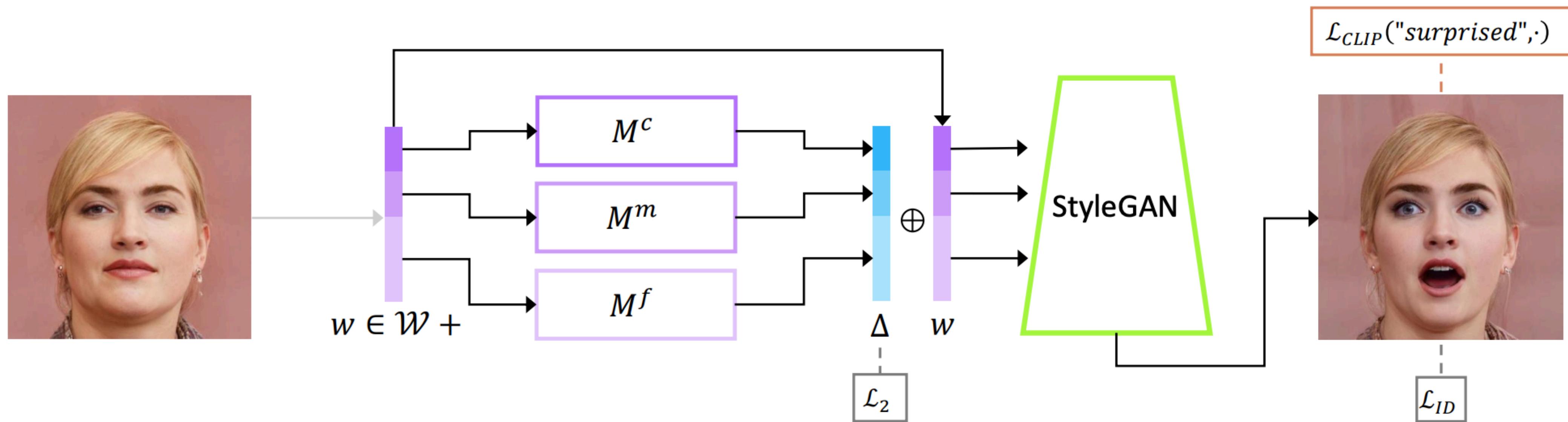
CLIP (Contrastive Language–Image Pre-training)

- Чувствителен к описанию

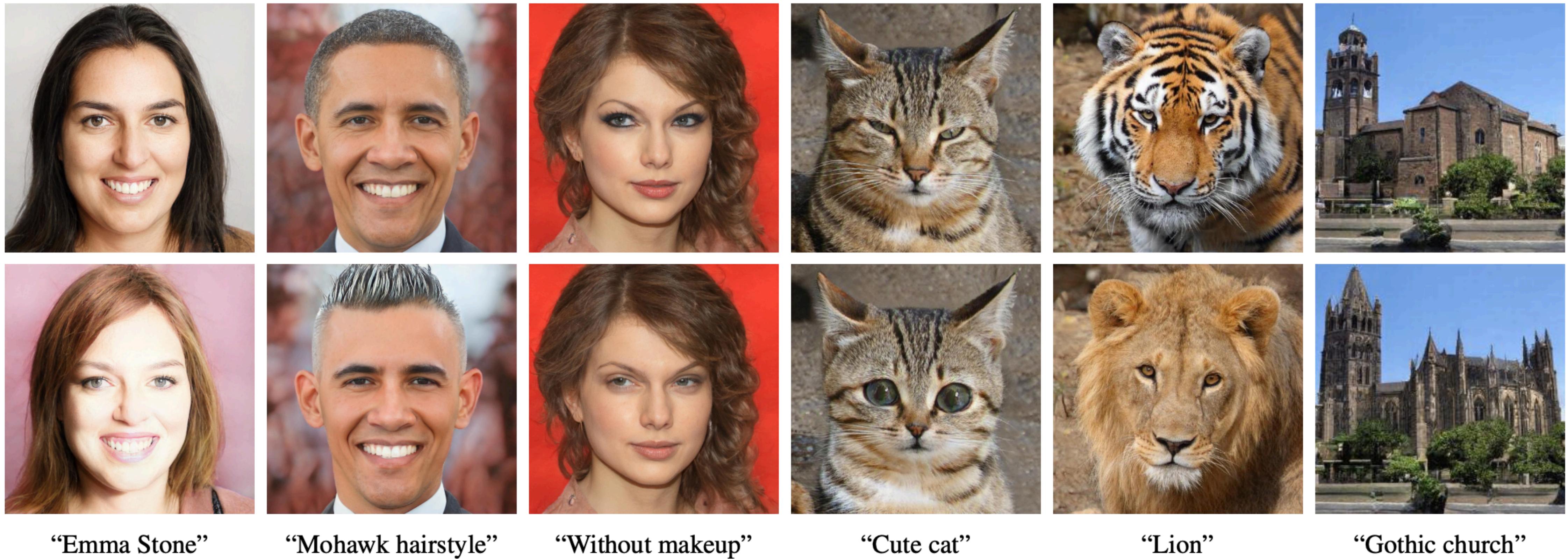


CLIP + StyleGAN

- Используют преодобученные StyleGAN и CLIP
- Цель - выучить трансформации над исходной картинкой, чтобы результат генерации StyleGAN был близок к текстовому описанию (эмбеддингам CLIP)



CLIP + StyleGAN



DALL-E

DALL-E

Задача: по текстовому описанию сгенерировать картинку

TEXT PROMPT an armchair in the shape of an avocado....

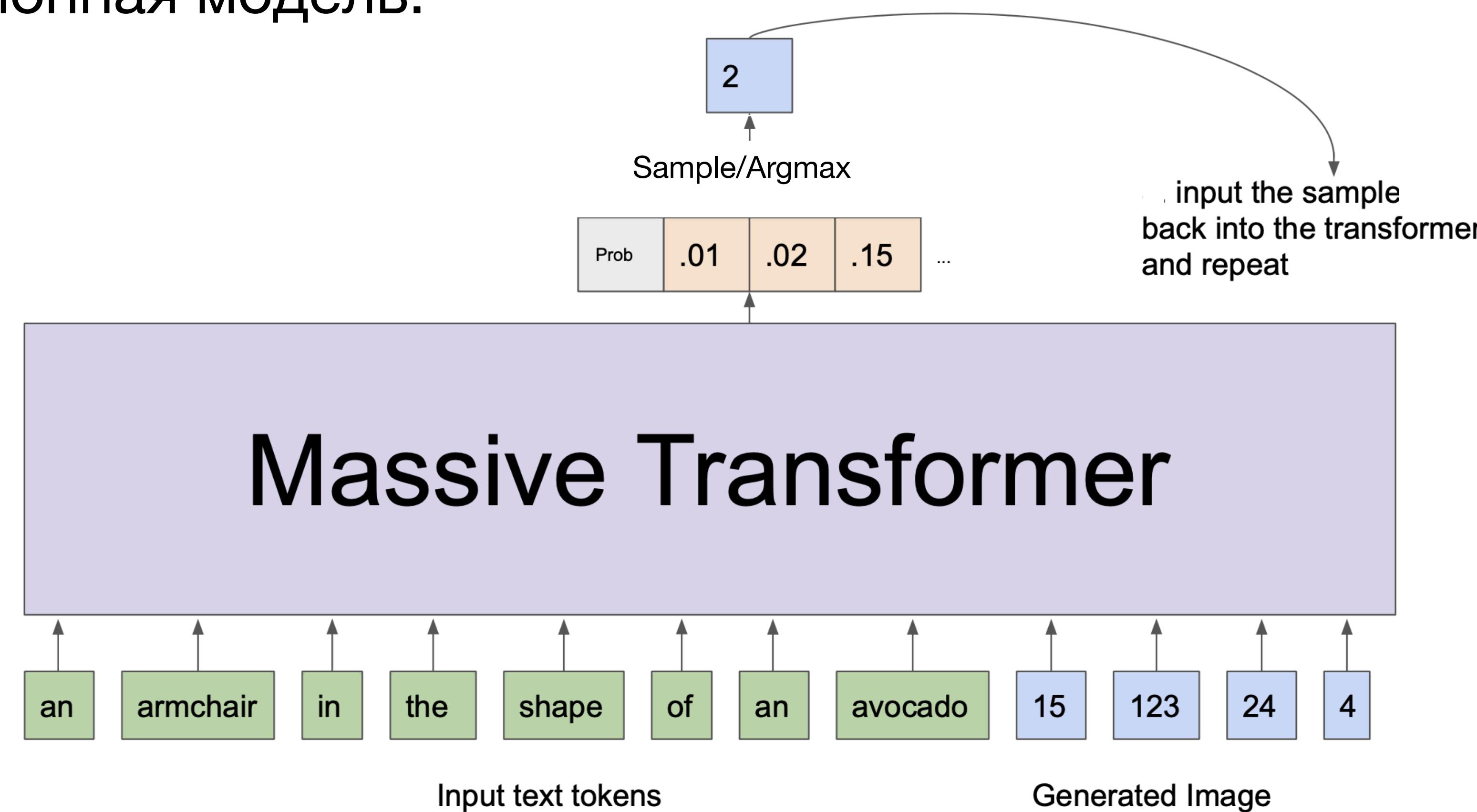
AI-GENERATED
IMAGES



DALL-E

Задача: по текстовому описанию сгенерировать картинку

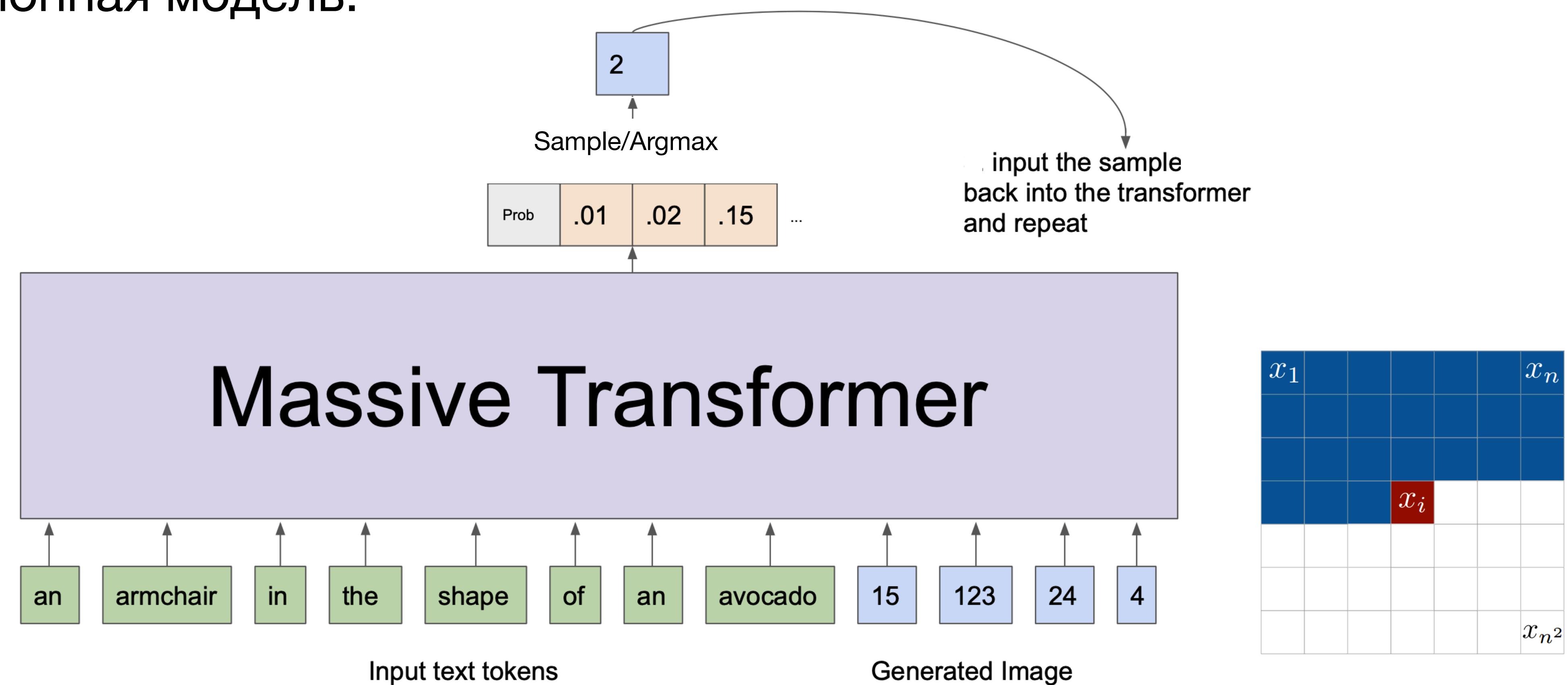
Авторегрессионная модель:



DALL-E

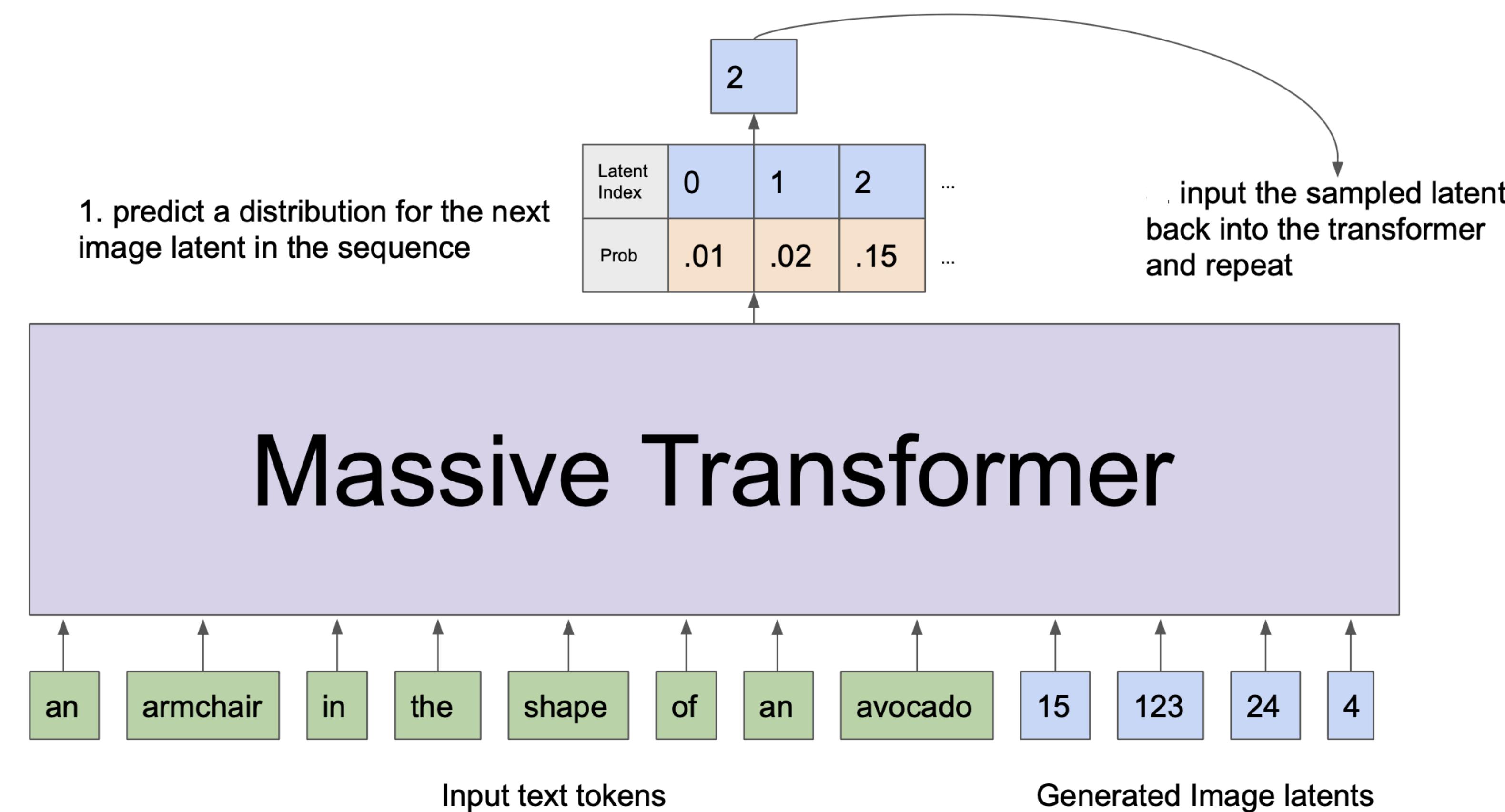
Задача: по текстовому описанию сгенерировать картинку

Авторегрессионная модель:

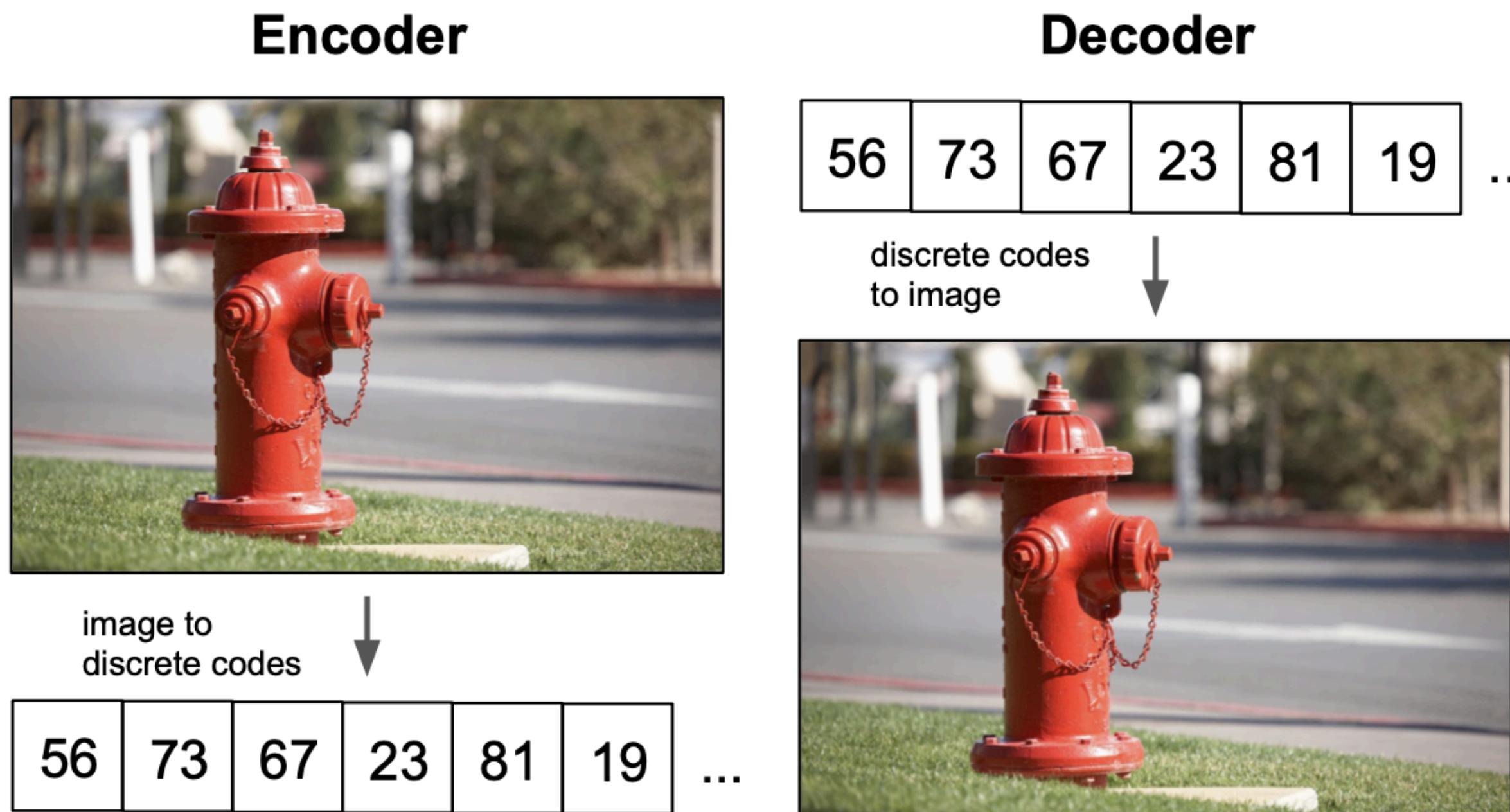


DALL-E

Генерировать по пикселям долго - переведем картинку в пространство меньшей размерности



VQ-VAE



VQ-VAE

Reconstructed Input:

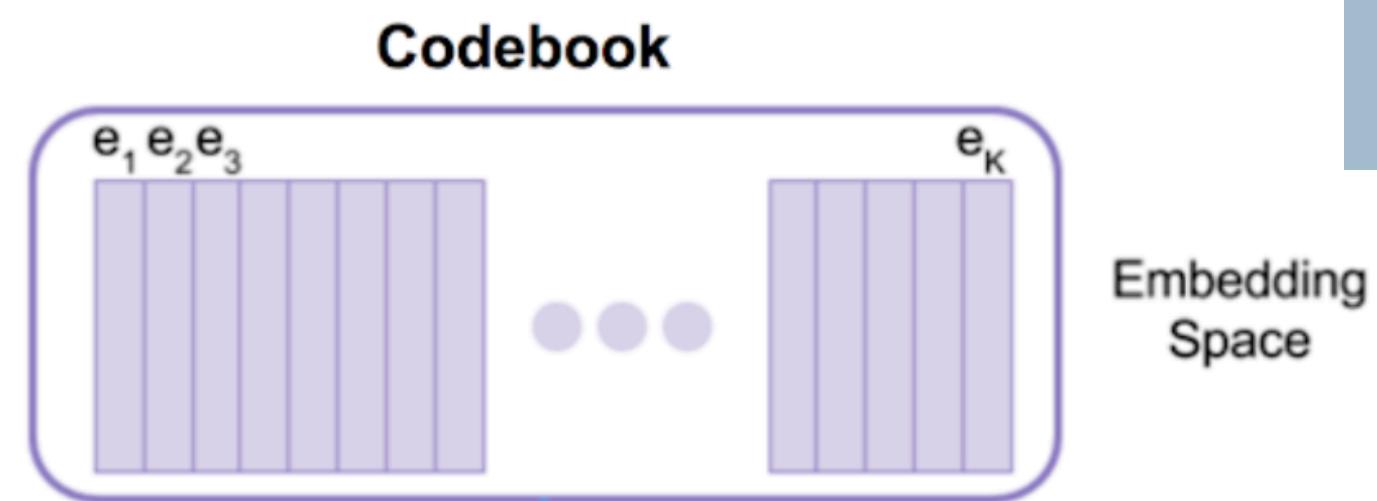
\hat{x}

Decoder

z_q

Vector Quantization Layer

$$z_q(x) = \operatorname{argmin}_i \|e_i - z_e(x)\|_2^2$$



z_e

Encoder

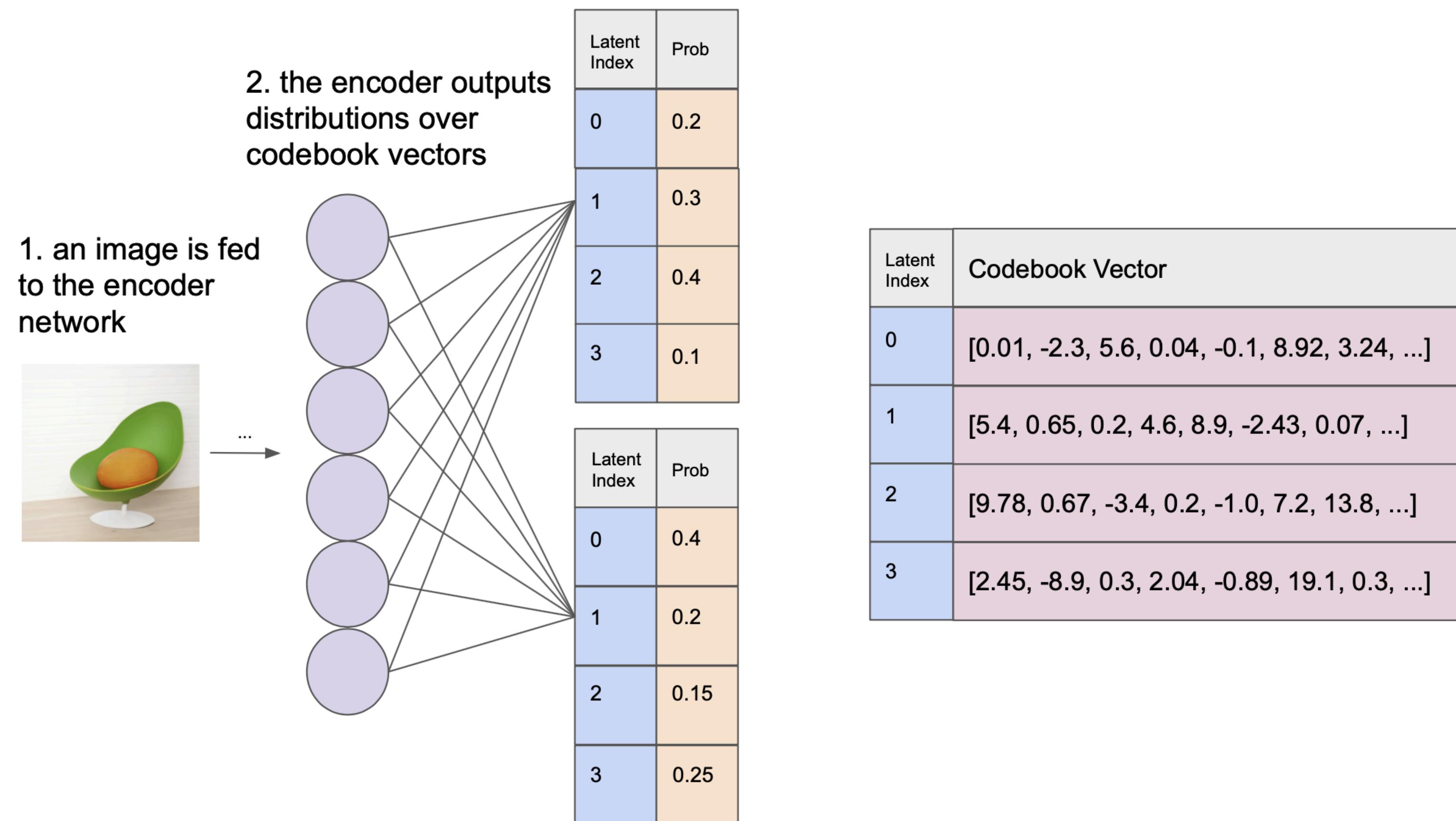
Input:

x

dVAE

Отличия от VQ-VAE:

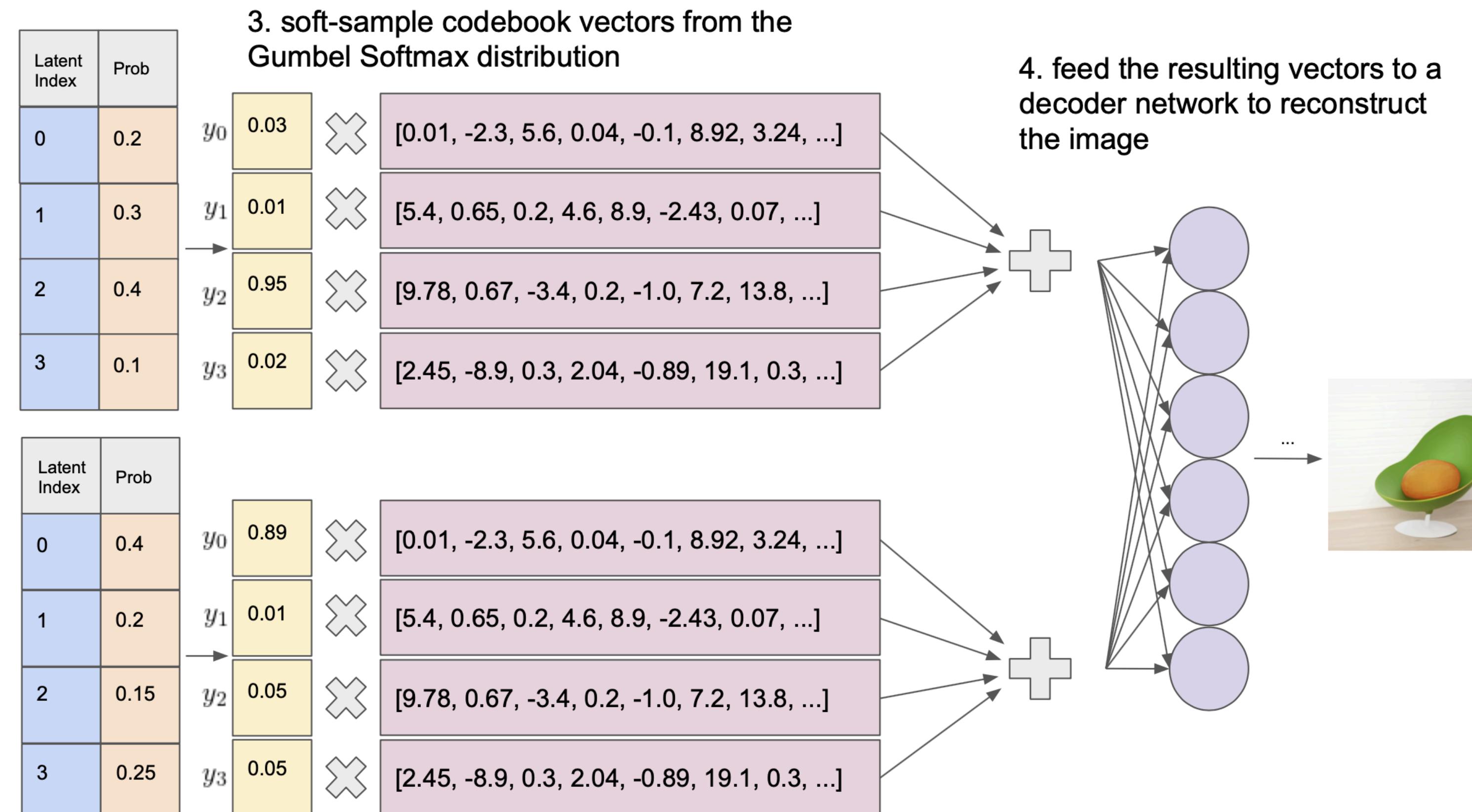
- Вместо выбора одного индекса для codebook предсказываем распределение



dVAE

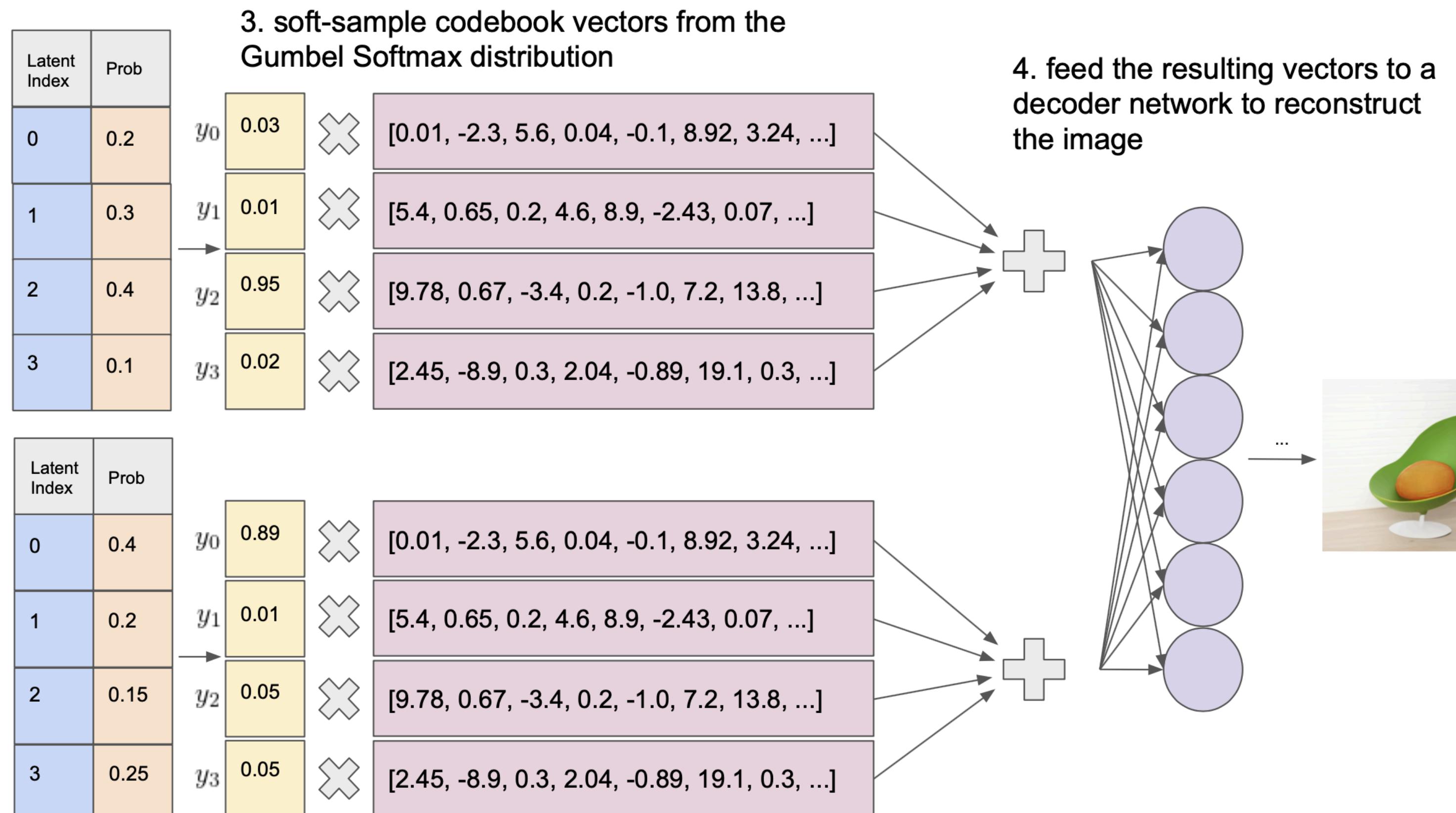
Отличия от VQ-VAE:

- чтобы текли градиенты - используем дифференцируемое близкое распределение (Gumbel Softmax Relaxation)



dVAE

$$\text{Loss: } -\mathbb{E}_{z \sim q(z|x)} \log(p(x|z)) + KL(q(z|x) || p(z))$$



dVAE

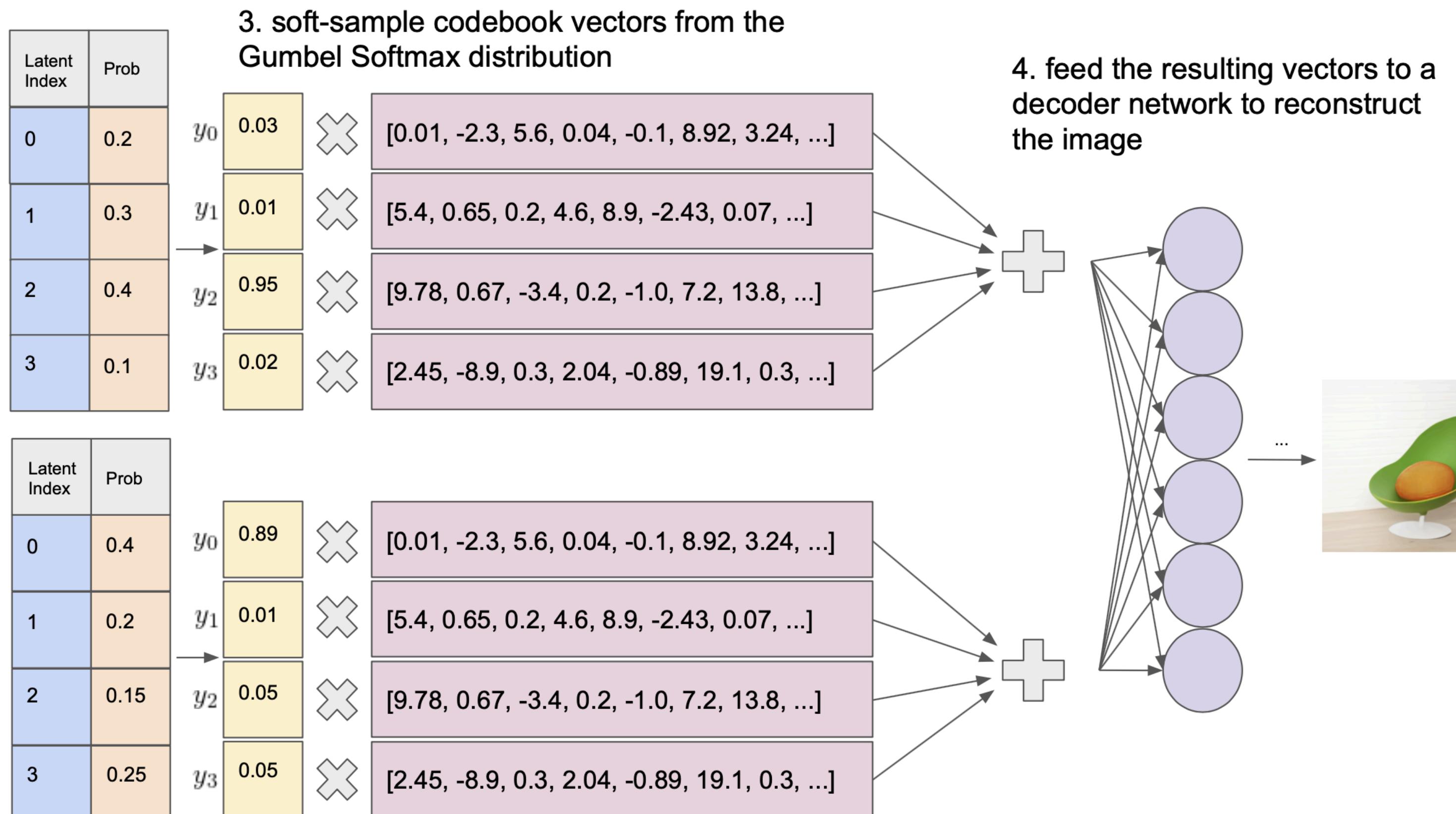
relaxation distribution
(from encoder)

uniform distribution

Loss: $-\mathbb{E}_{z \sim q(z|x)} \log(p(x|z)) + KL(q(z|x) || p(z))$

$KL - \text{“расстояние” между распределениями}$

$D_{KL}(P || Q) = \sum_{x \in \mathcal{X}} P(x) \log\left(\frac{P(x)}{Q(x)}\right).$

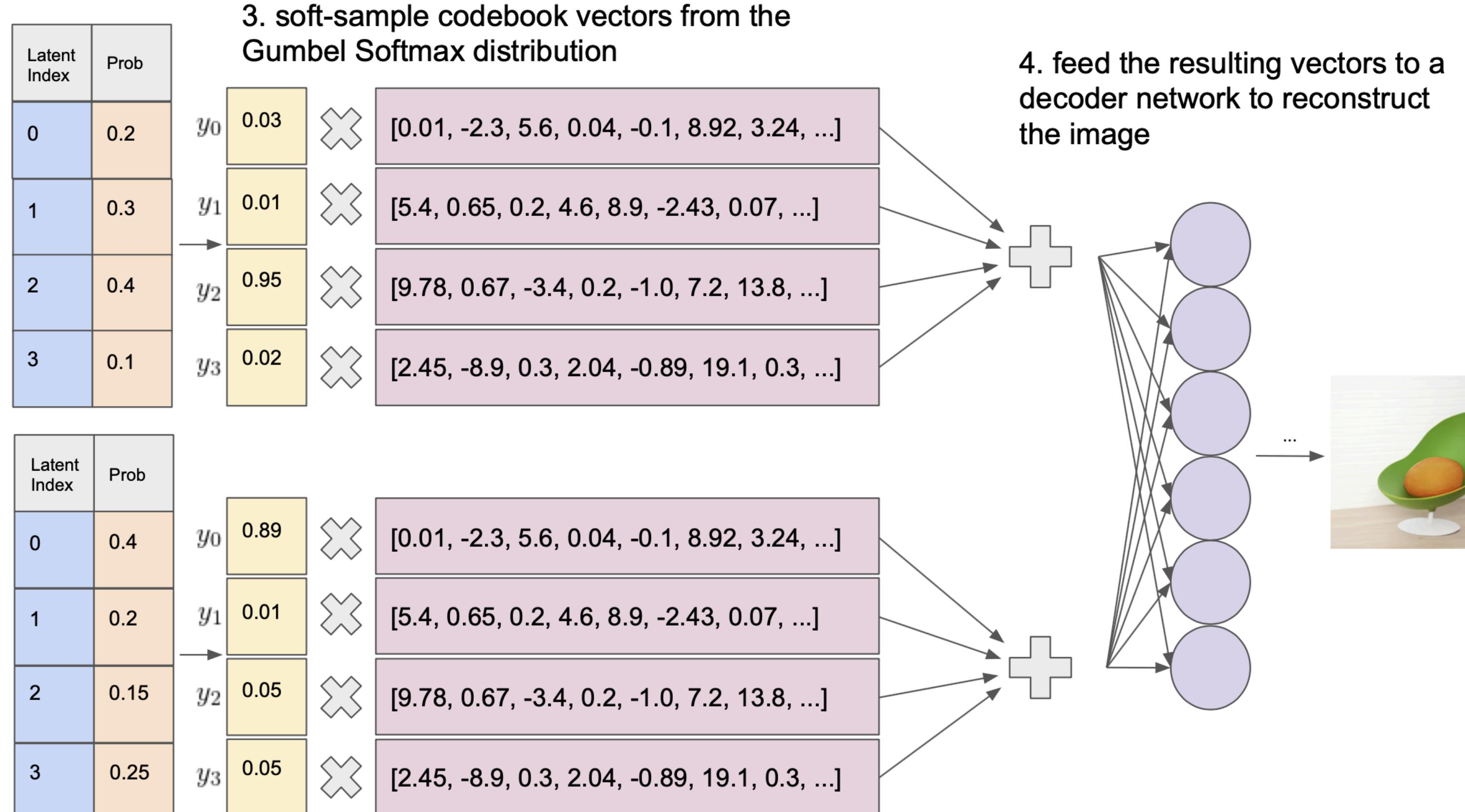


dVAE

Decoder output

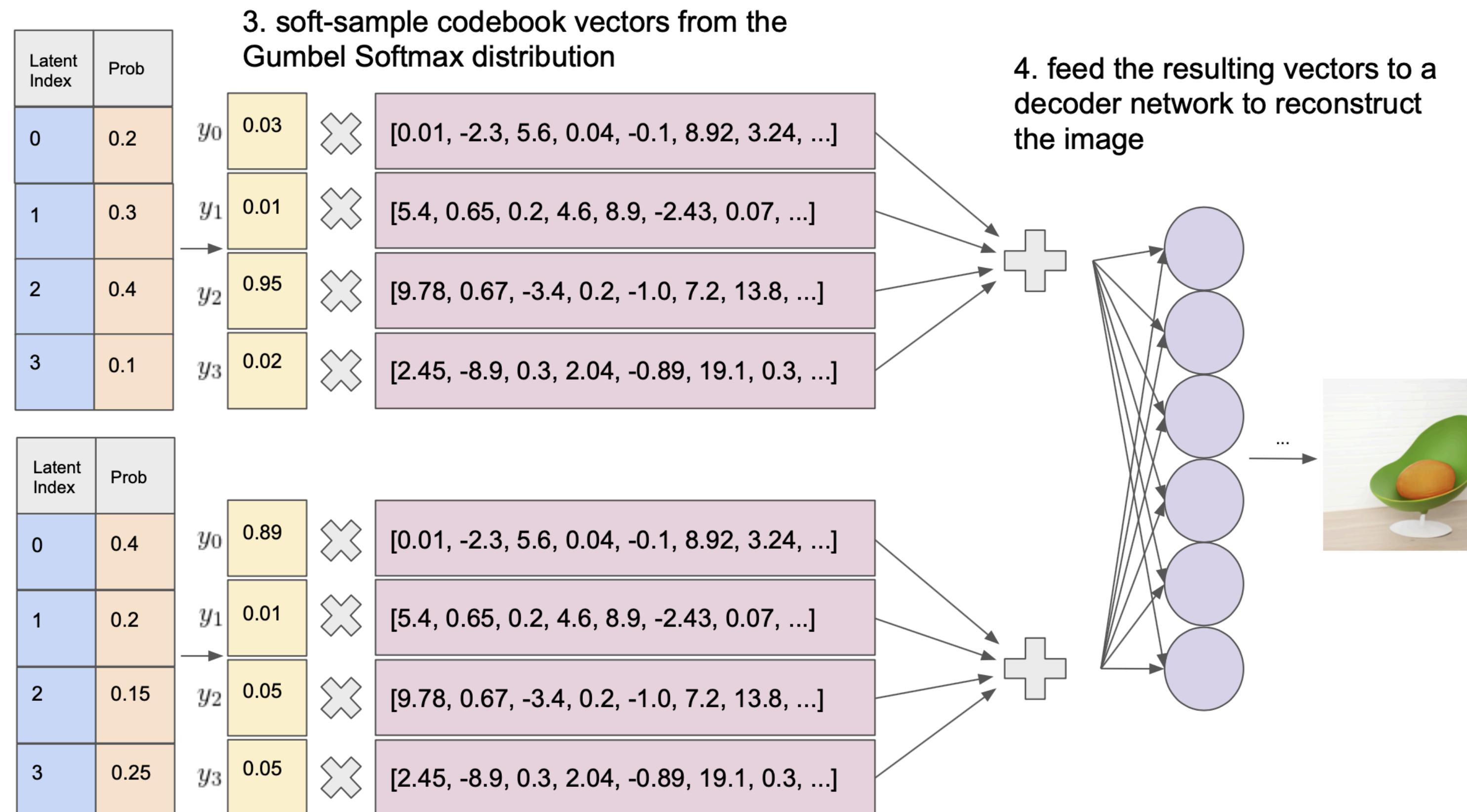
$$\text{Loss: } -\mathbb{E}_{z \sim q(z|x)} \log(p(x|z)) + KL(q(z|x) || p(z))$$

Полученная картинка должна быть похожа на исходную



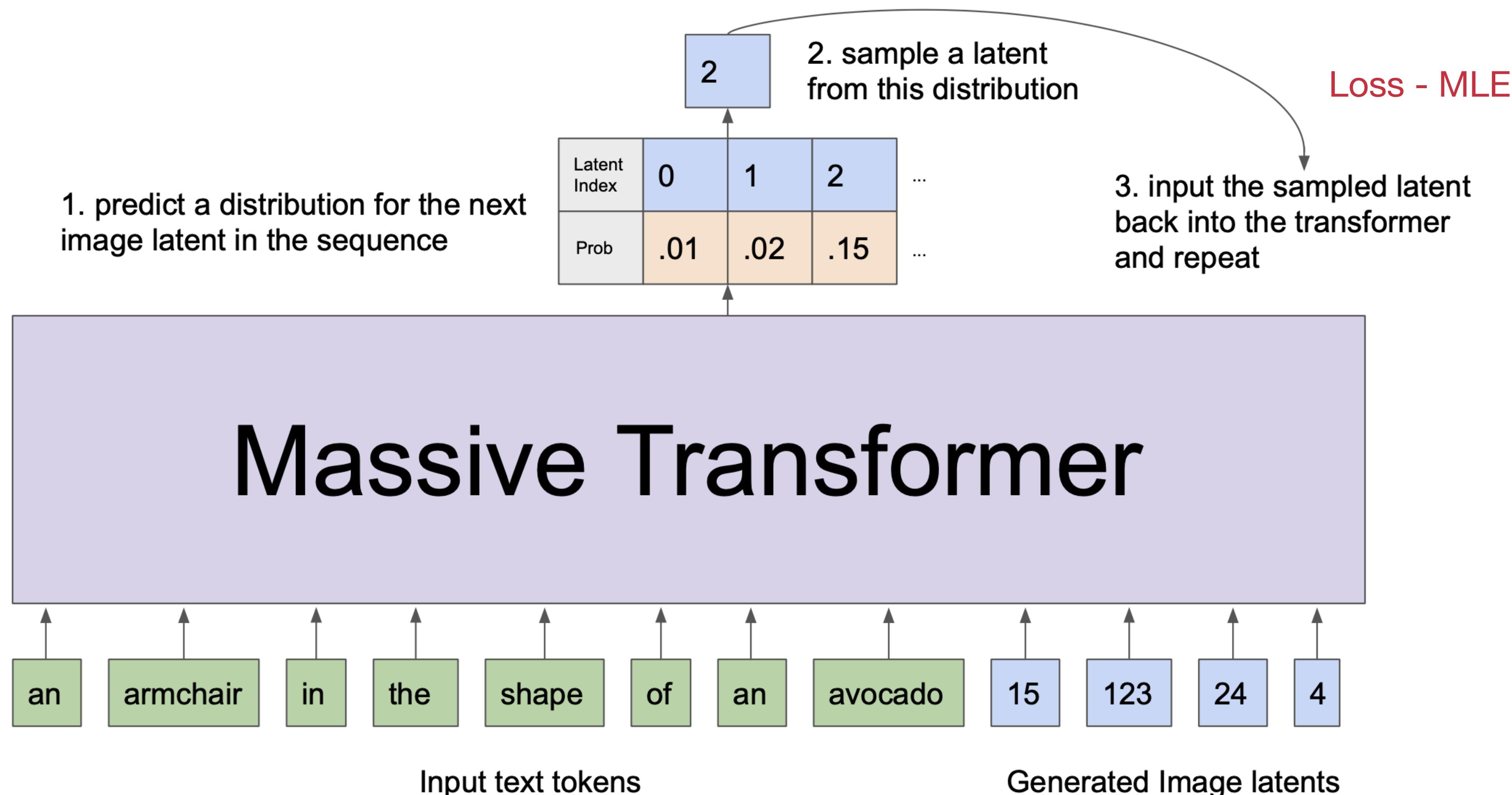
dVAE

Если есть распределение над индексами codebook -
можно генерировать картинку (decoder dVAE)



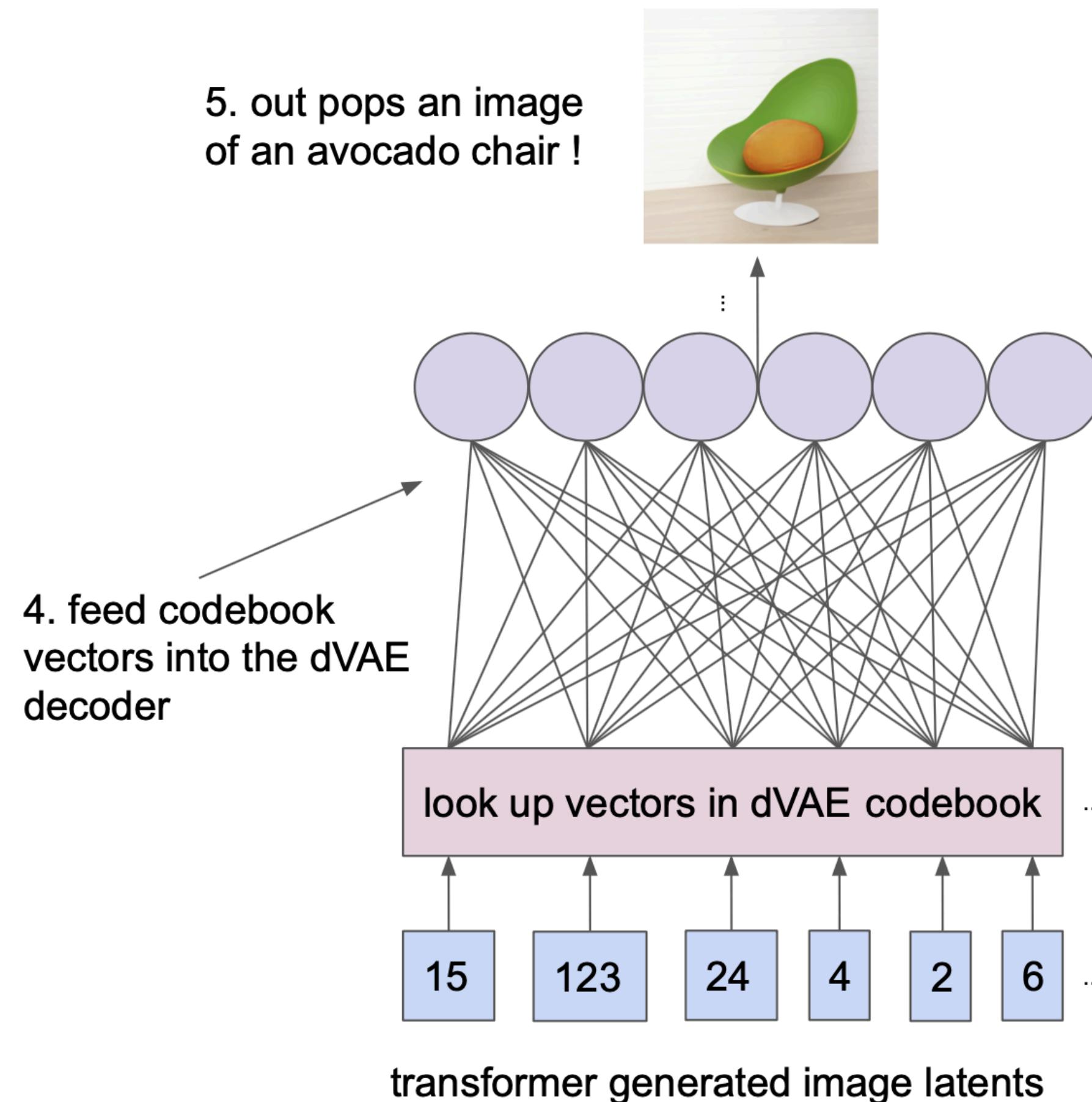
DALL-E

1. Transformer по входному тексту авторегрессионно генерирует индексы



DALL-E

2. dVAE по сгенерированным индексам генерирует картинку



DALL-E

TEXT PROMPT

an illustration of a baby daikon radish in a tutu walking a dog

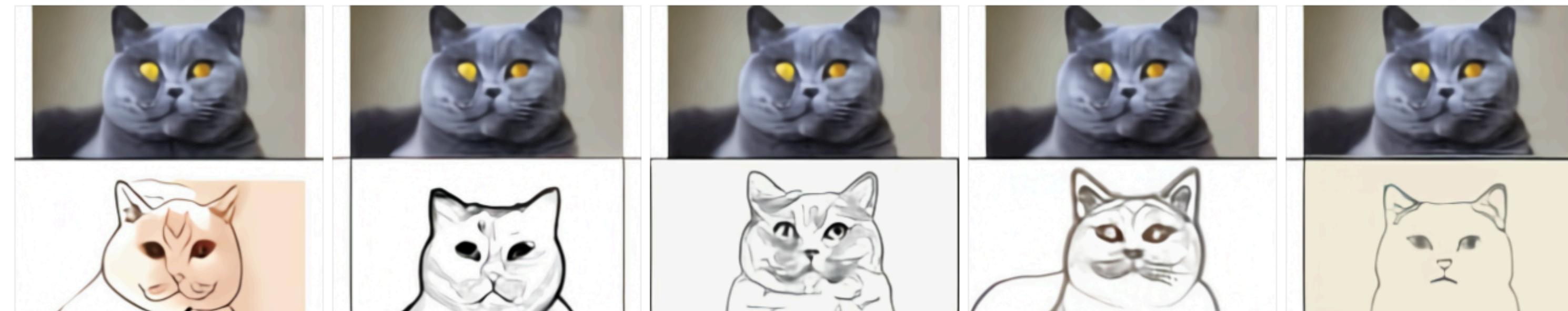
AI-GENERATED
IMAGES



TEXT & IMAGE
PROMPT

the exact same cat on the top as a sketch on the bottom

AI-GENERATED
IMAGES



DALL-E 2

Модель: CLIP + GLIDE (Denoising Diffusion Probabilistic Model)



a dolphin in an astronaut suit on saturn, artstation



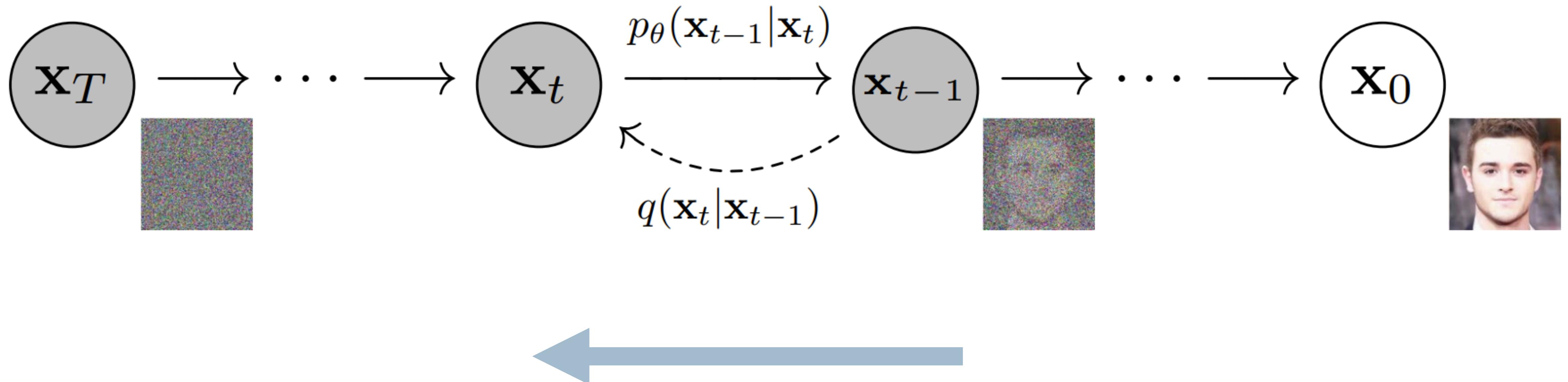
a propaganda poster depicting a cat dressed as french emperor napoleon holding a piece of cheese



a teddy bear on a skateboard in times square

DDPM: overview

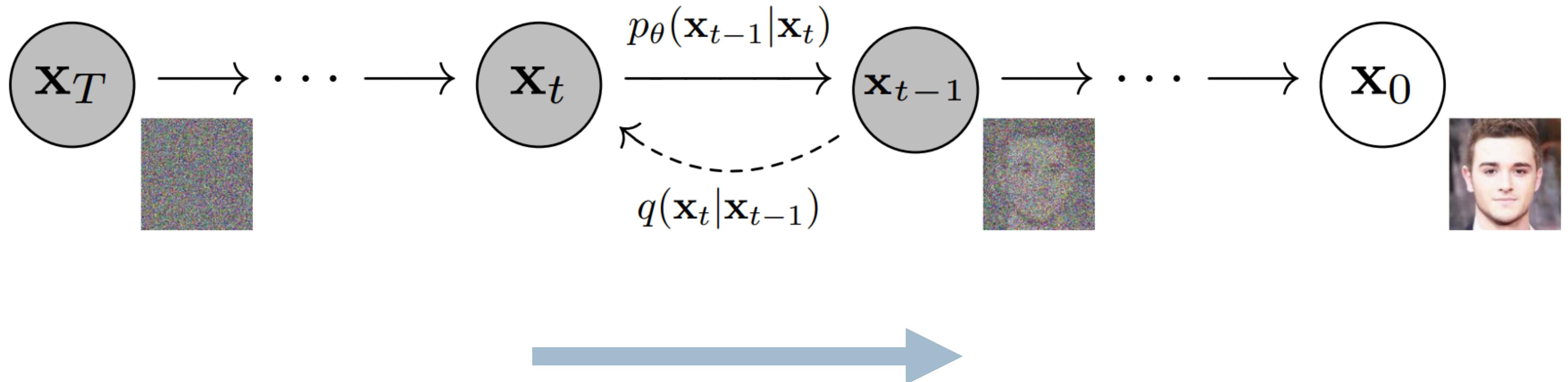
DDPM:



Добавляем случайный шум на каждом шаге

DDPM: overview

DDPM:

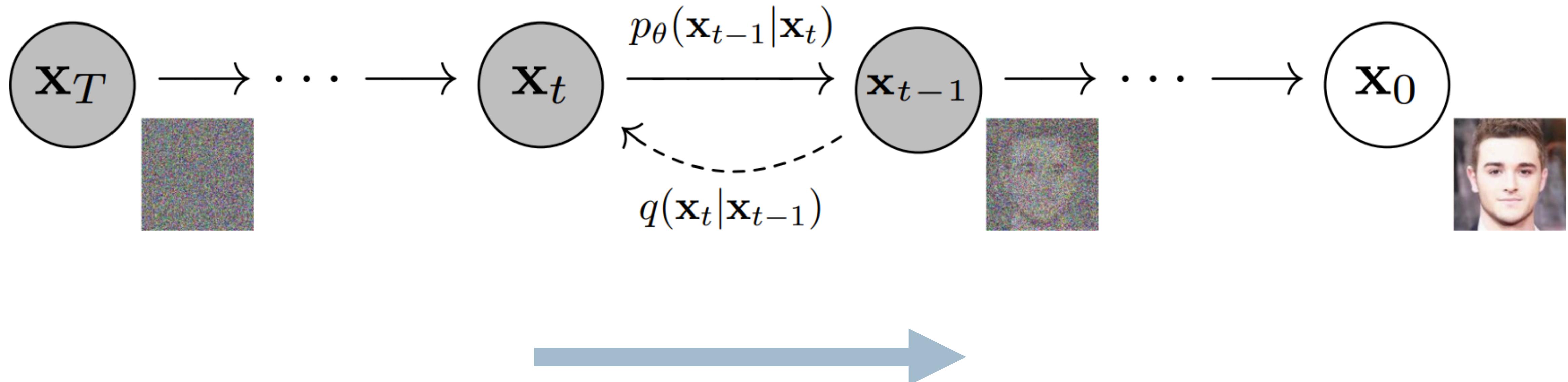


Восстанавливаем (убираем добавленный шум)

DDPM: overview

DDPM:

Markov chain



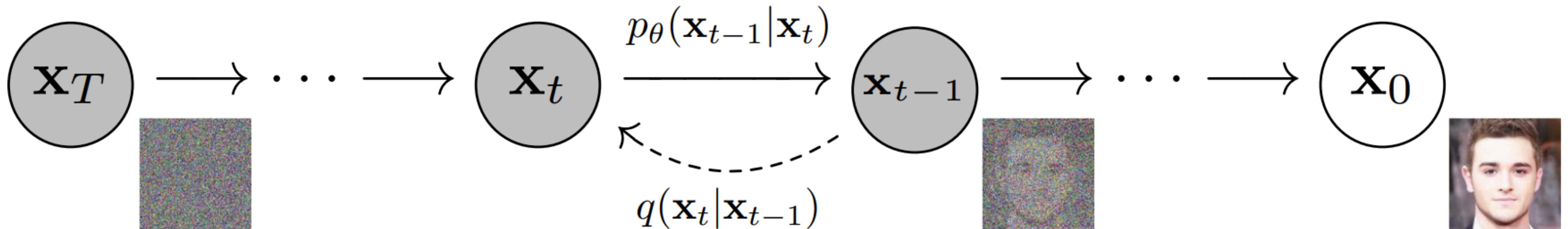
Восстанавливаем (убираем добавленный шум)

Каждый раз применяем модель типа UNet

DDPM: overview

DDPM:

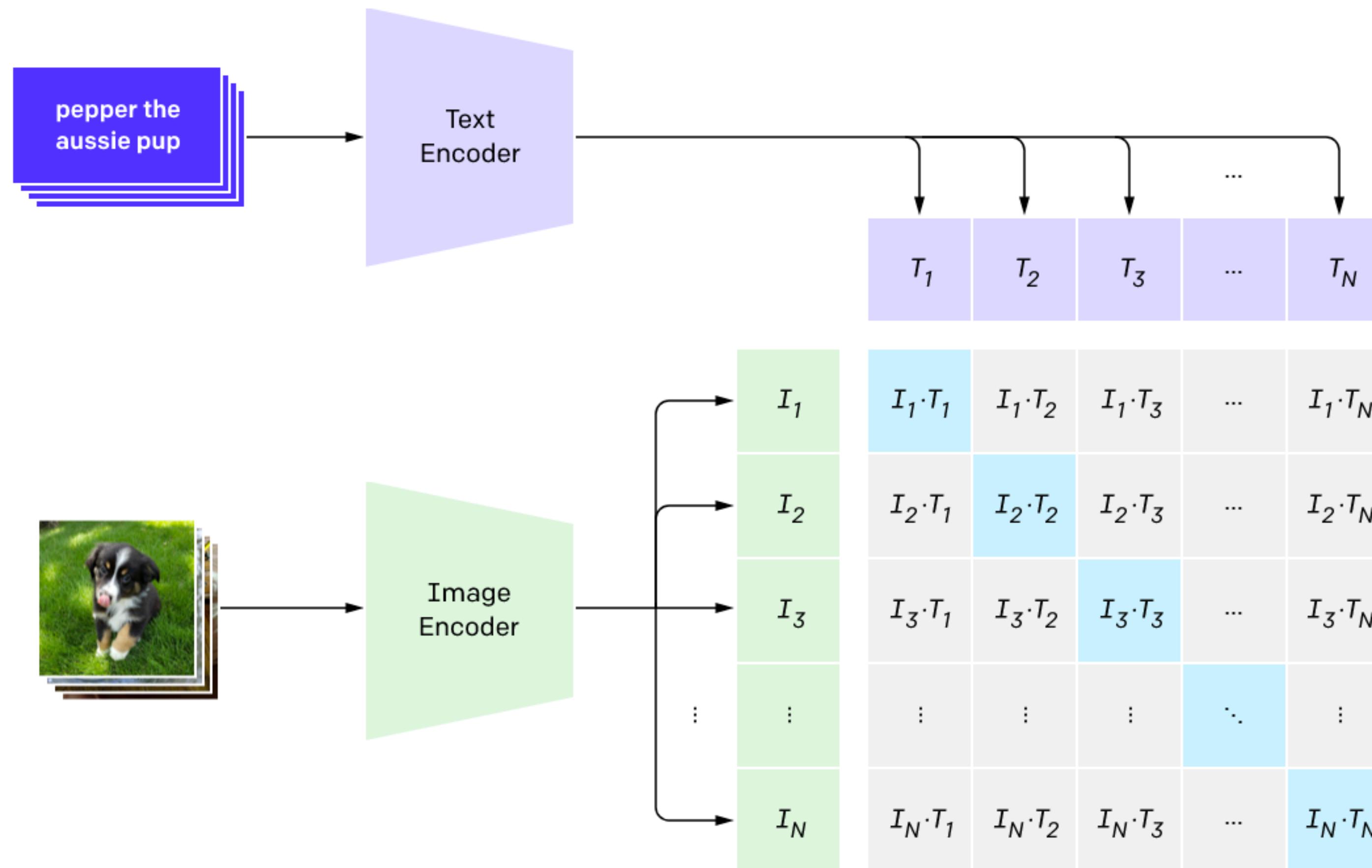
Markov chain



$$L_{\text{simple}} := E_{t \sim [1, T], x_0 \sim q(x_0), \epsilon \sim \mathcal{N}(0, \mathbf{I})} [||\epsilon - \epsilon_\theta(x_t, t)||^2]$$

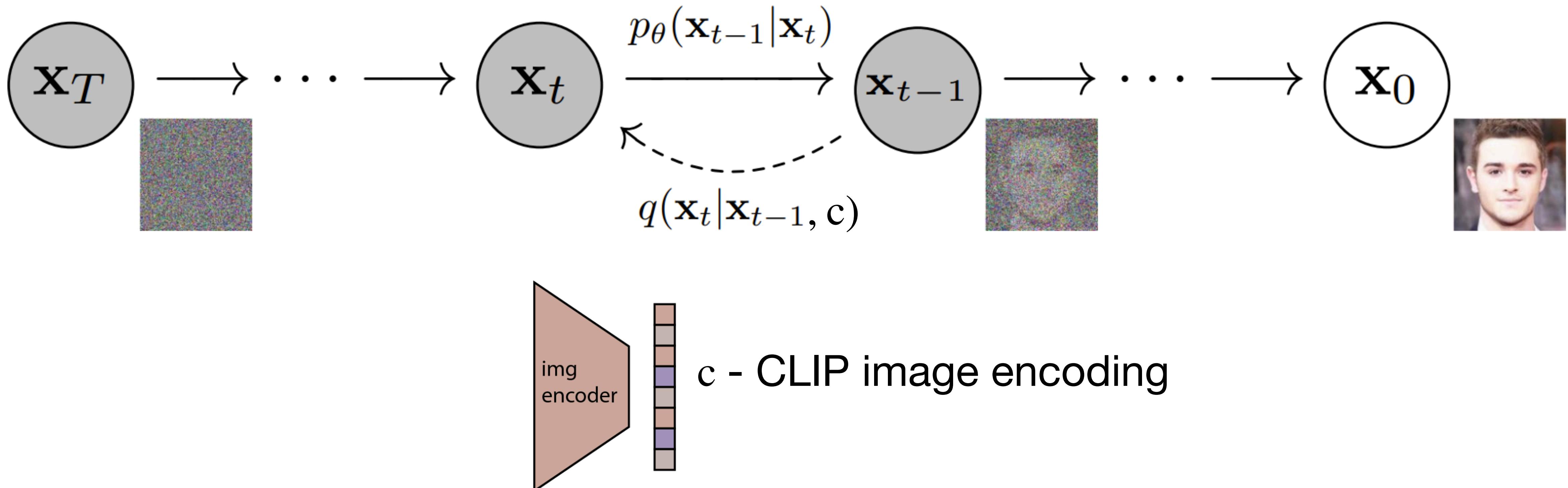
DALL-E 2

1. Обучаем CLIP



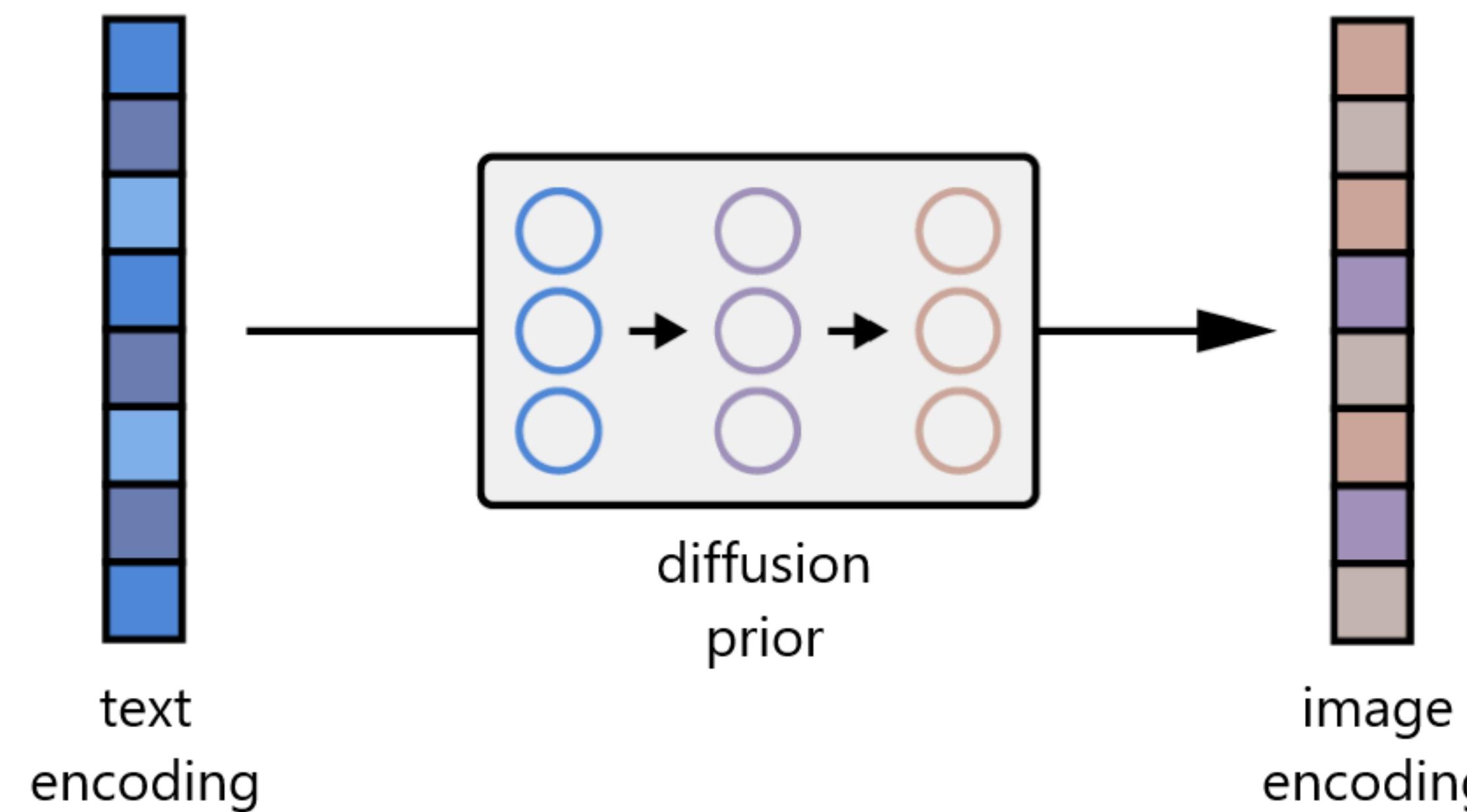
DALL-E 2

1. Обучаем CLIP
2. Обучаем DDPM обусловленную на CLIP image encoding (GLIDE)



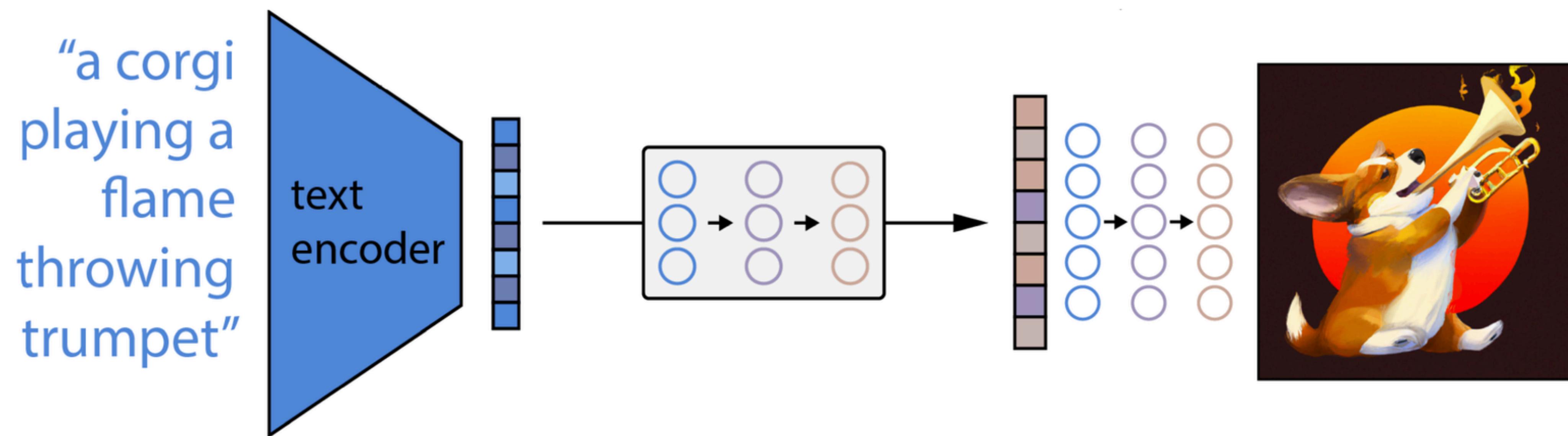
DALL-E 2

1. Обучаем CLIP
2. Обучаем DDPM обусловленную на CLIP image encoding (GLIDE)
3. Обучаем другую DDPM предсказывать CLIP image encoding (по text encoding)



DALL-E 2

Генерация:



DALL-E 2



an espresso machine that makes coffee from human souls, artstation



panda mad scientist mixing sparkling chemicals, artstation



a corgi's head depicted as an explosion of a nebula