



**ІТМО**

**Современные архитектуры  
нейронных сетей**

**Мультимодальность**

# Что такое модальность

ІТМО



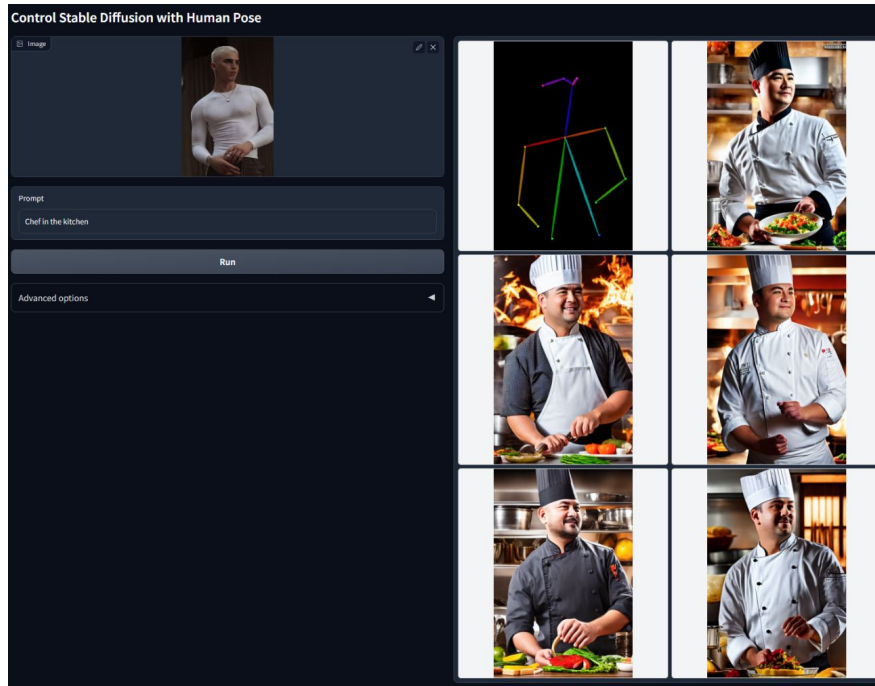
Однородные данные



Разнородные данные

# Что такое модальность

ИТМО

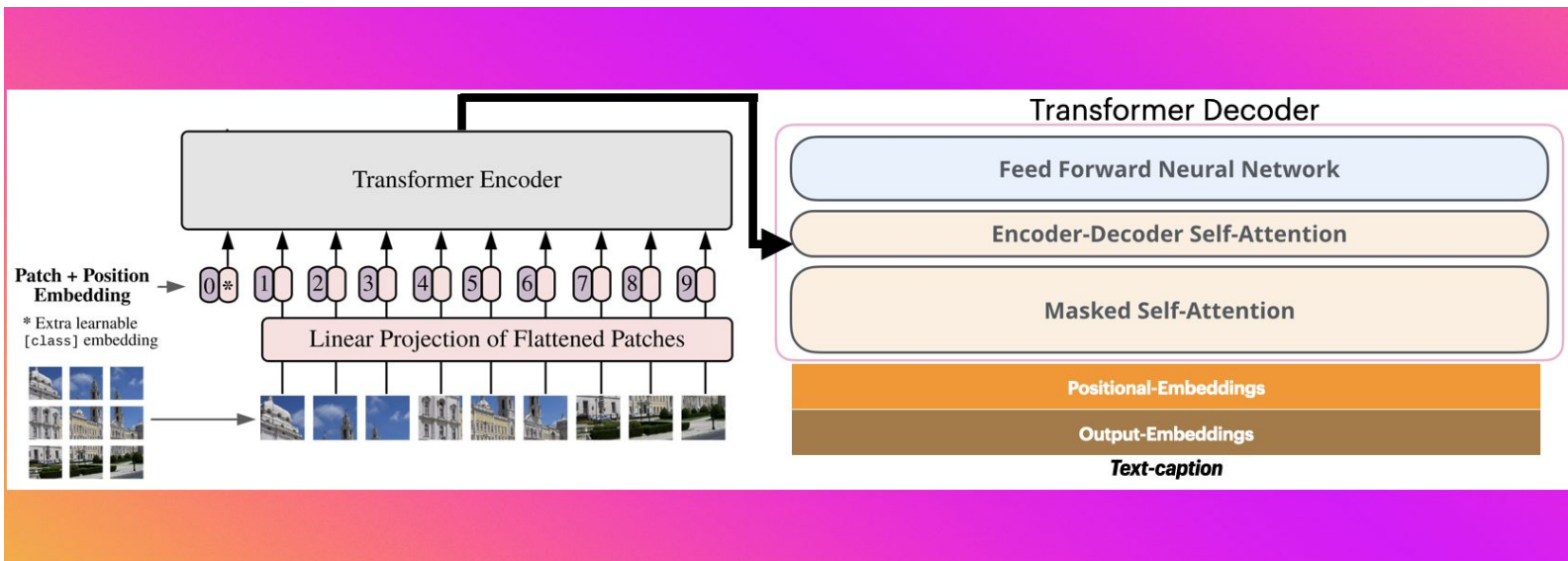


# Чем больше модальностей - тем лучше

ІІТМО



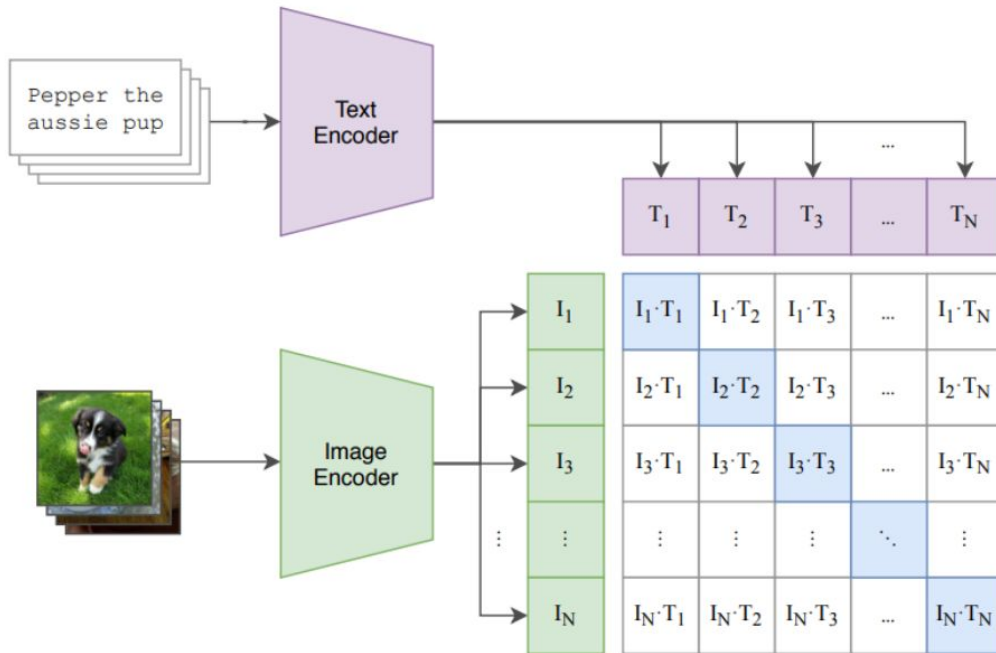
# Image Captioning



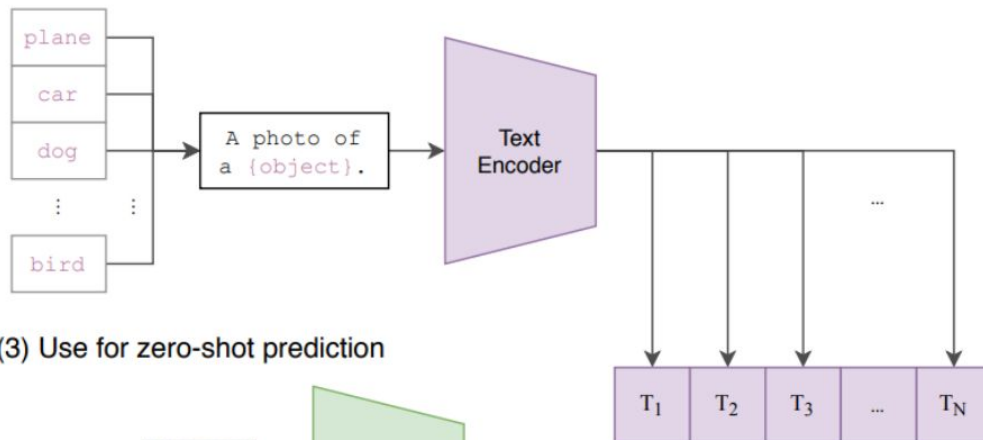
$\{I_k\}_{k=1}^{\text{batch\_size}}$  - векторы изображений

$\{\tau_k\}_{k=1}^{\text{batch\_size}}$  - векторы текста

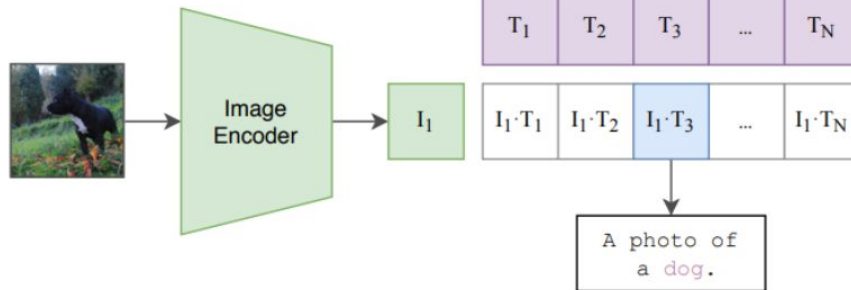
$$\text{cosine\_similarity}(I_k, \tau_j) = \frac{(I_k, \tau_j)}{\|I_k\| \cdot \|\tau_j\|}$$



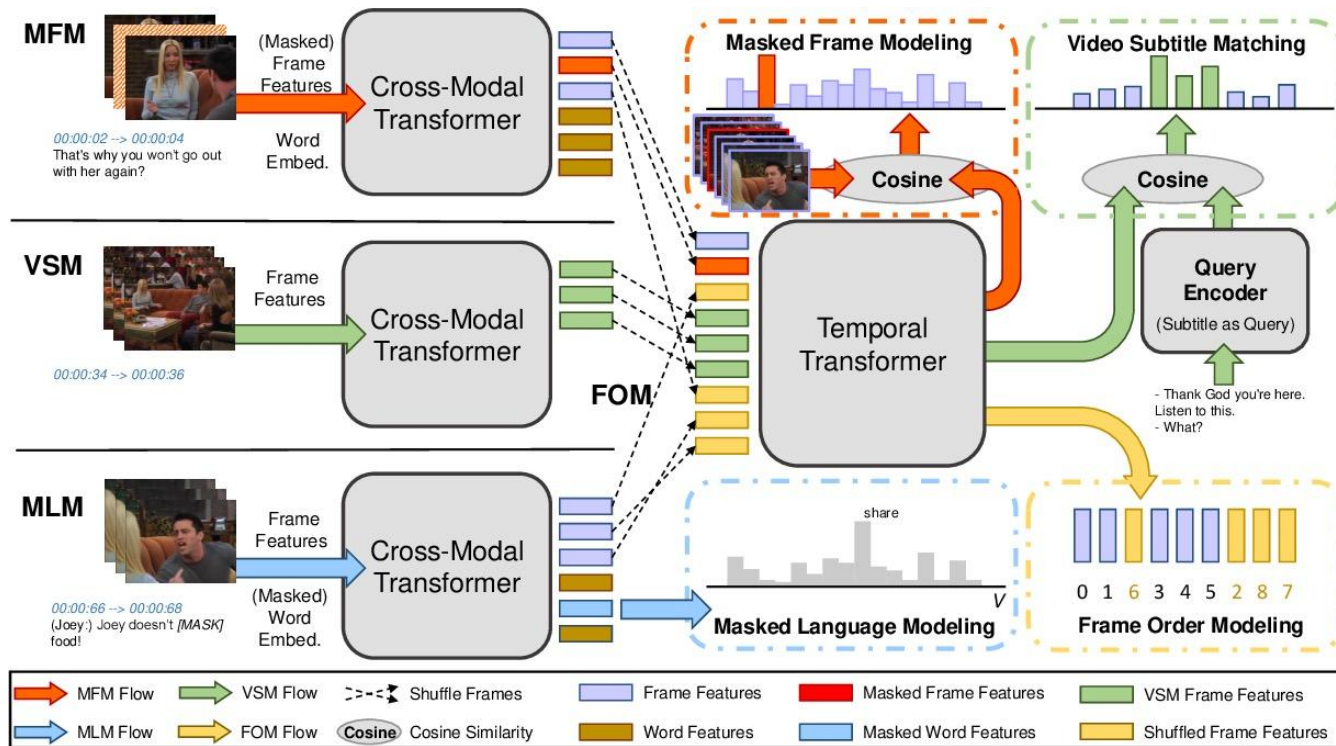
(2) Create dataset classifier from label text



(3) Use for zero-shot prediction



# Кроссмодальность трансформеров ИТМО



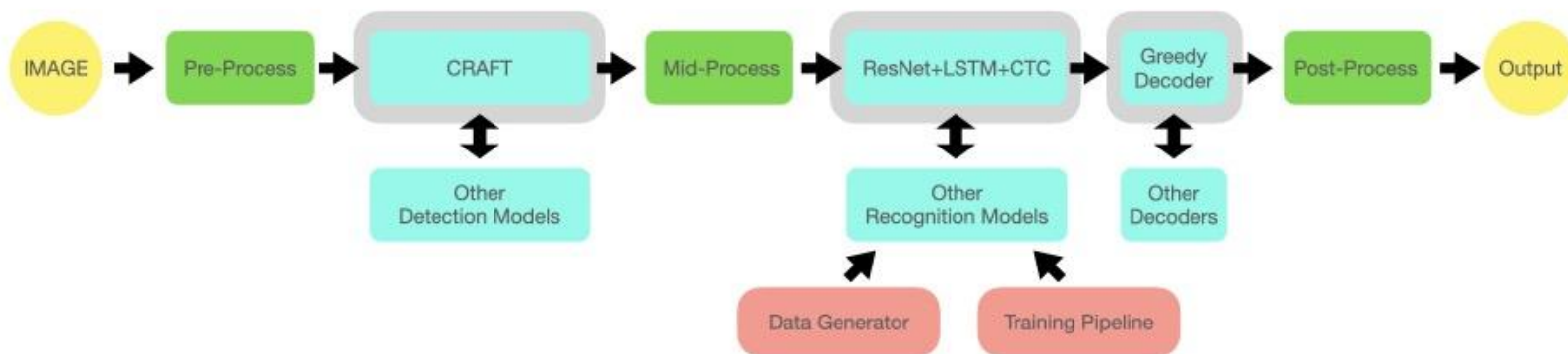
Картинка для  
устрашения



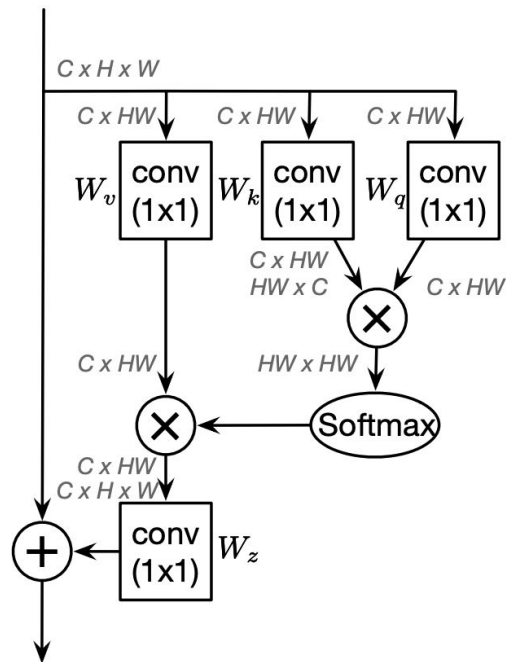
# Свертки vs трансформеры



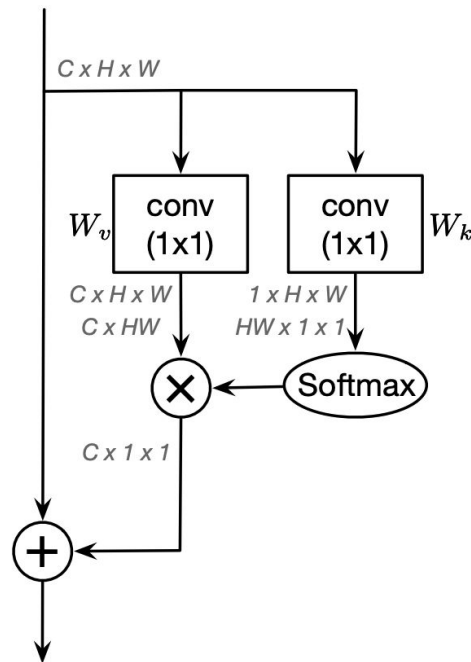
## EasyOCR Framework



# Свертки vs трансформеры

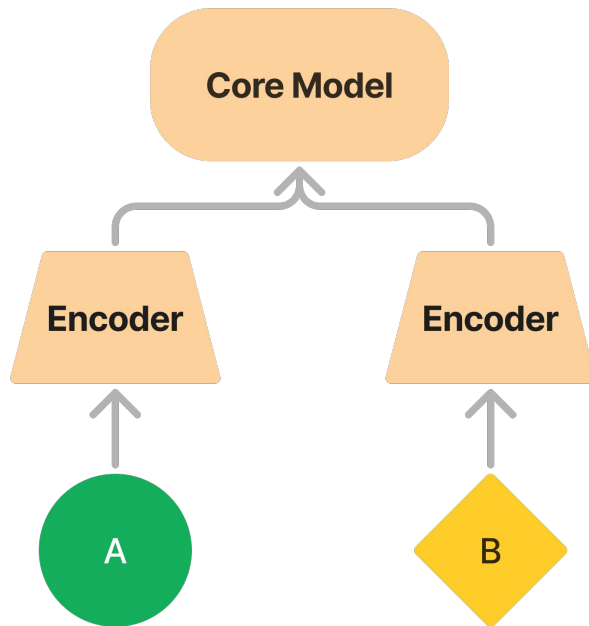


(a) NL block



(b) Simplified NL block (Eqn 2)

# Смешивание модальностей. Early Fusion



# Смешивание модальностей. Early Fusion



language model (Vicuna v1.5 13B)



vision-language connector (MLP)



tokenizer &  
embedding

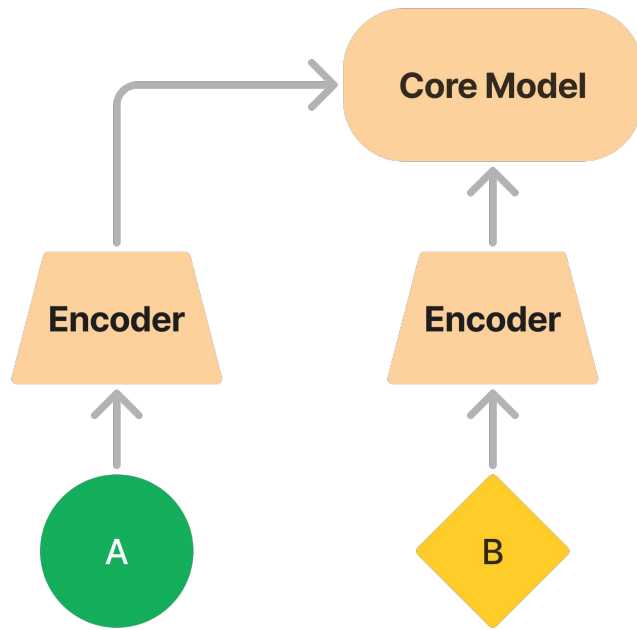
vision encoder (CLIP ViT-L/336px)



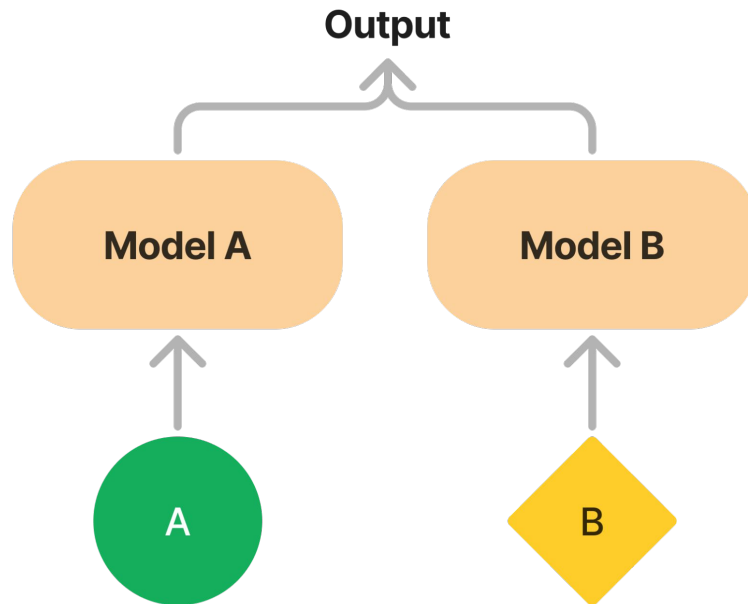
User: what is  
unusual about  
this image?

# Смешивание модальностей. Deep Fusion

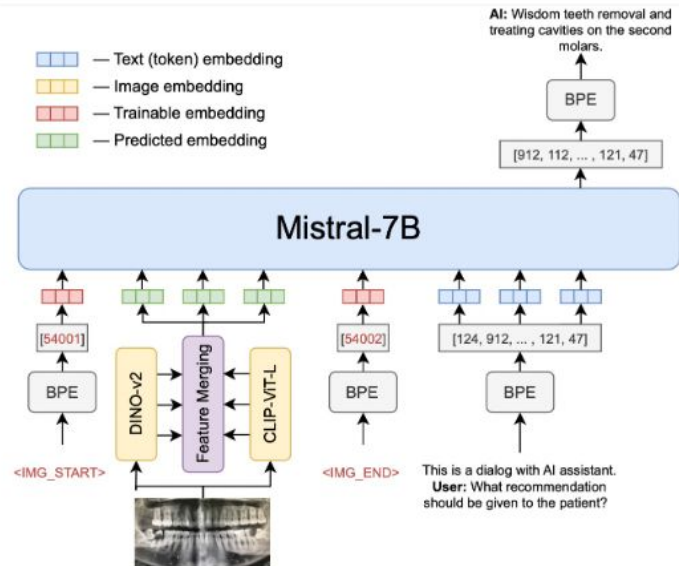
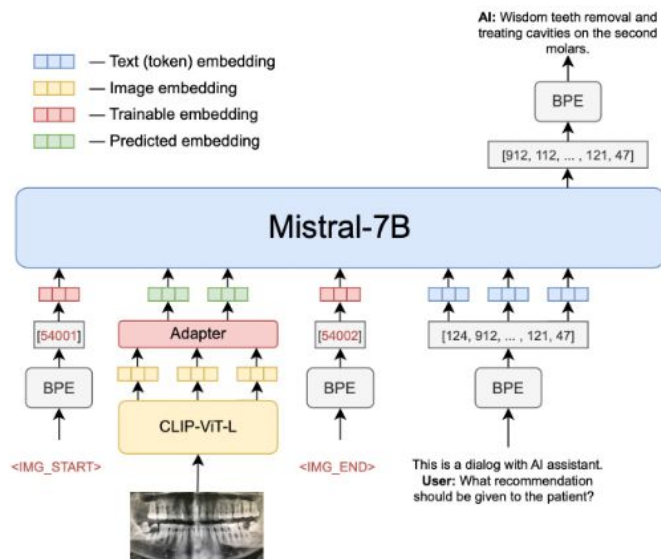
ИТМО



# Смешивание модальностей. Late Fusion



# Смешивание модальностей. Merge of Encoders



- Pixel-shuffle
- Аддитивное смешивание

$$z = \sum_{i=1}^k w_i x_i,$$

- Мультипликативное смешивание

$$z = w(x_1 \times x_2)$$

- Gated Fusion

$$z = \sum_{i=1}^k g_i(x_1, x_2, \dots, x_k) x_i$$





# Captioning

ViTMO



language model (Vicuna v1.5 13B)



vision-language connector (MLP)



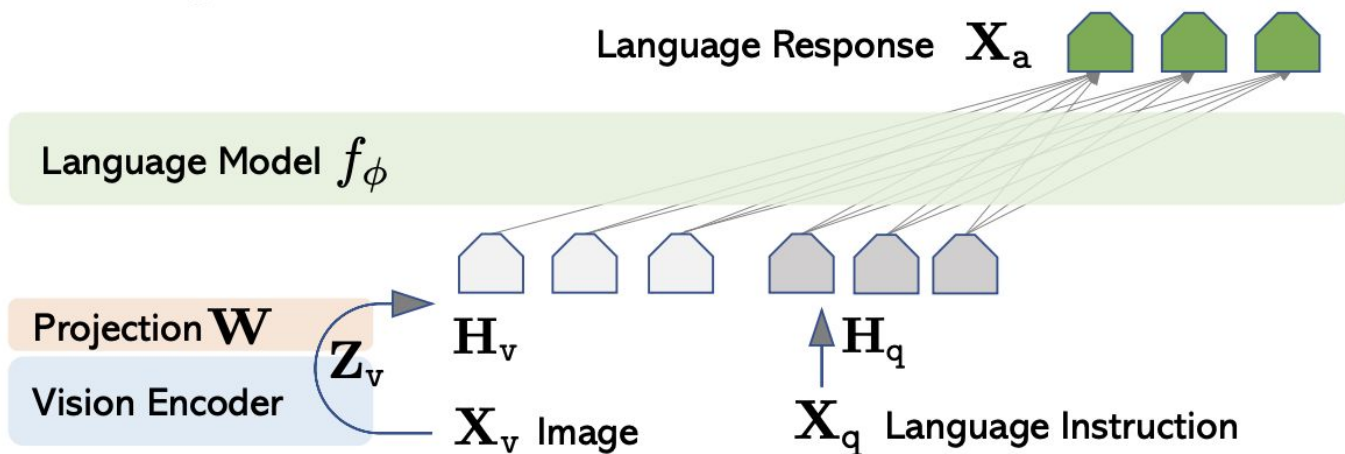
tokenizer &  
embedding

vision encoder (CLIP ViT-L/336px)



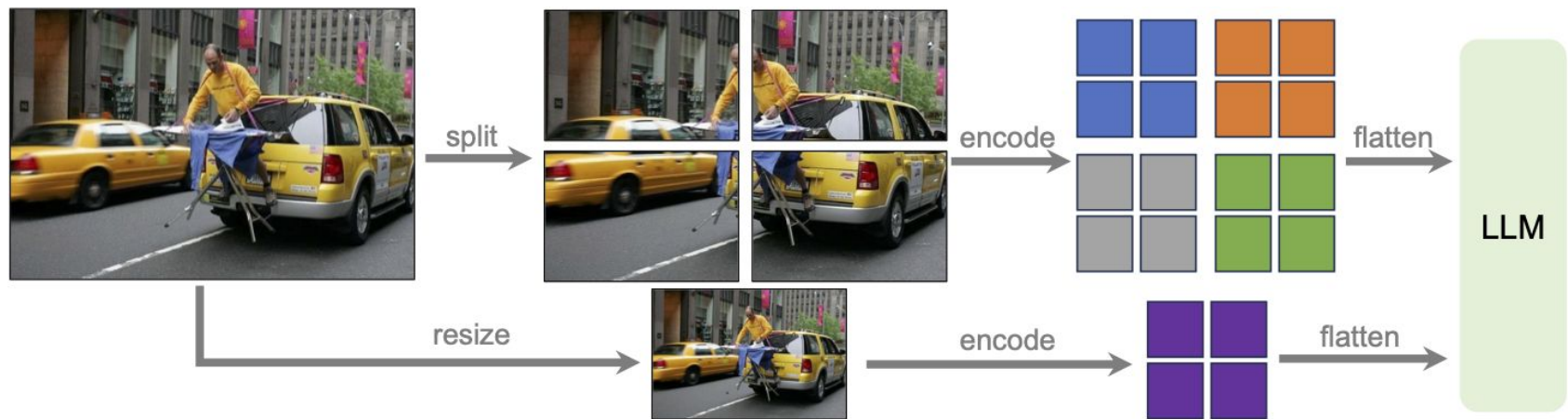
User: what is  
unusual about  
this image?

# Early Fusion VLM (LLaVA)

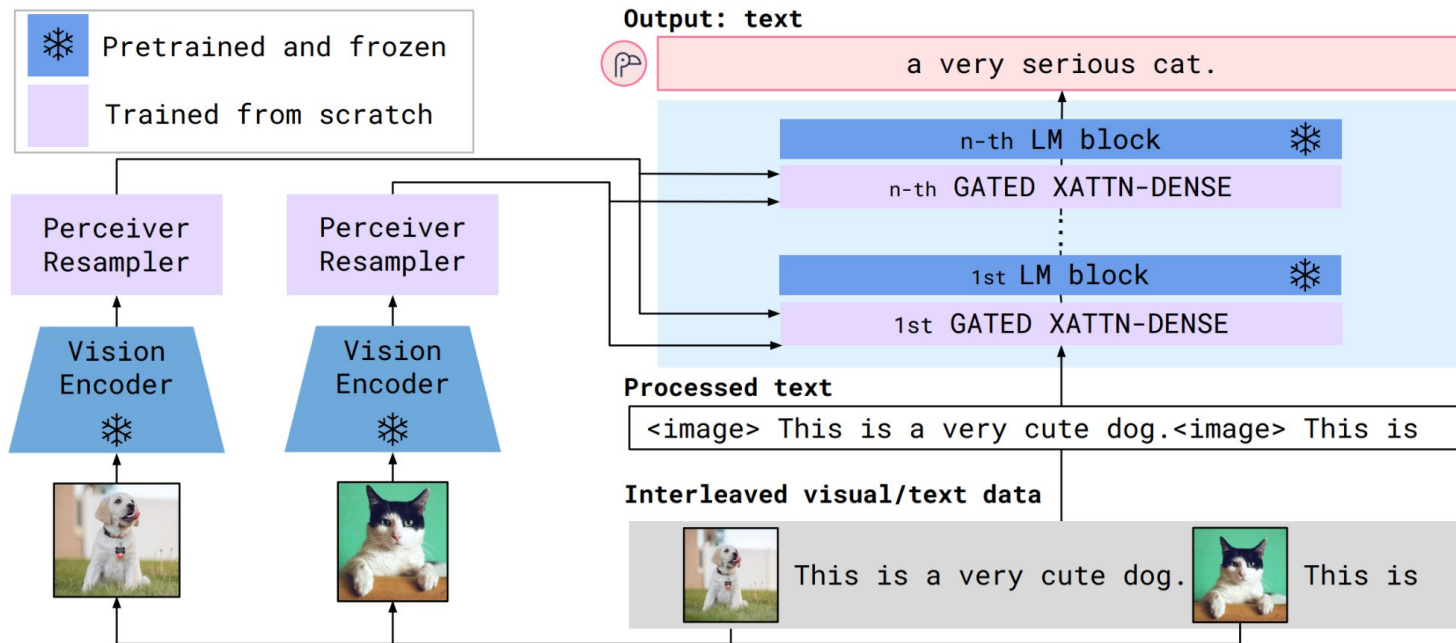


# LLaVA Next

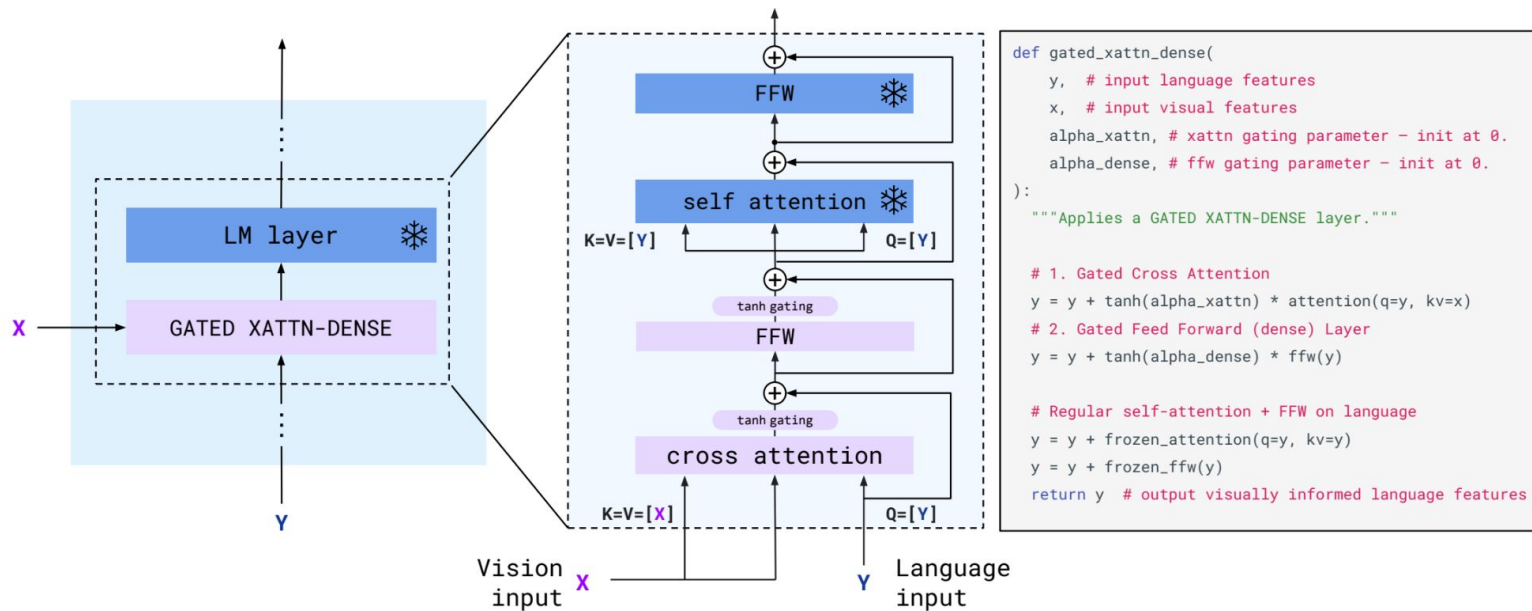
MITMO



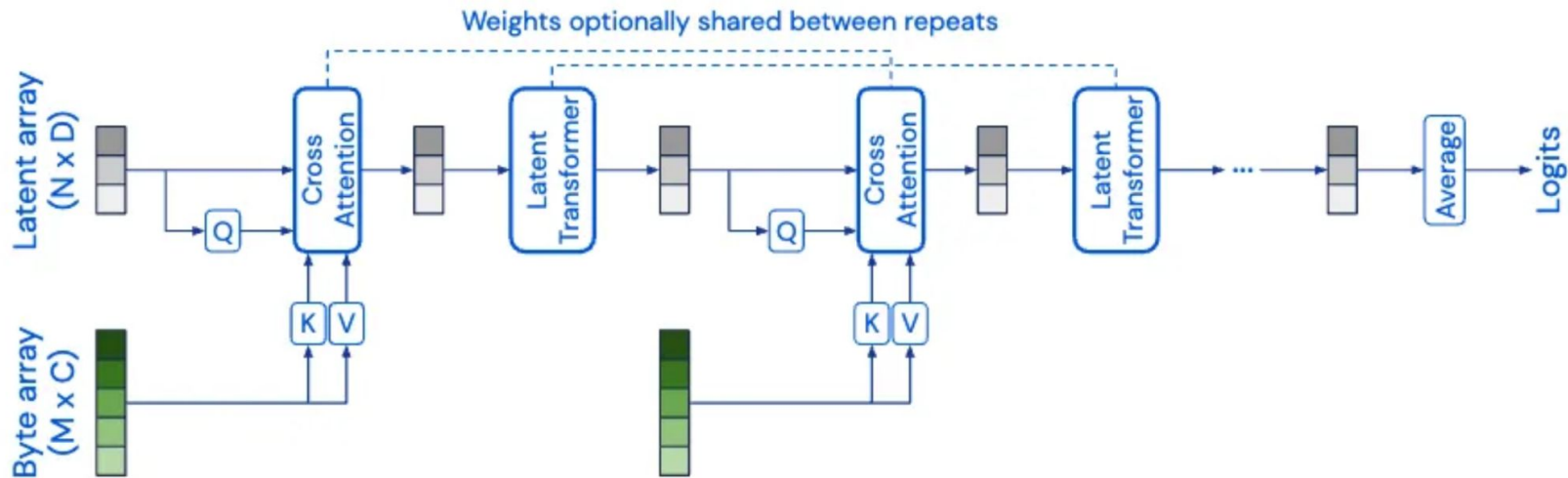
# Deep Fusion VLM (Flamingo)



# Deep Fusion VLM (Flamingo XAttn)



# Perceiver Resampler



# Visual Question Answering



C: A dog with goggles is in a motorcycle side car.

Q: Is motorcycle moving or still?

A: It's parked

Q: What kind of dog is it?

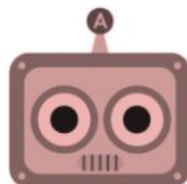
A: Looks like beautiful pit bull mix

Q: What color is it?

Image

Dialog history

Question



Visual Dialog  
model

Answer

A: Light tan with white patch that  
runs up to bottom of his chin



# DALL-E

DALL·E 2 can take an image and create different variations of it inspired by the original.



ORIGINAL IMAGE



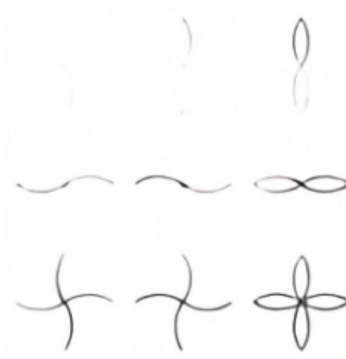
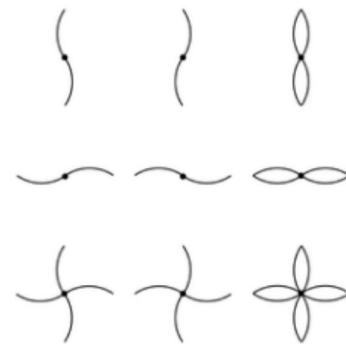
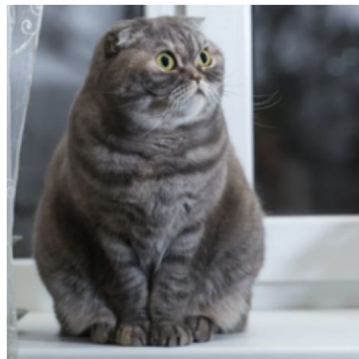
DALL·E 2 VARIATIONS





# Кодирование VAE

ИТМО



# VQ-GAN

