

Как использовать  
мультимодальные  
трансформеры в  
исследованиях и не  
страдать?

## Обо мне



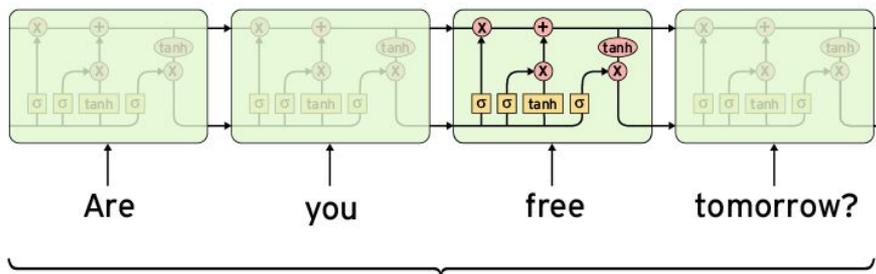
Исследую и применяю гигантские мультимодальные трансформеры, пишу статьи, имплементирую архитектуры



Коротко о трансформерах и почему их стоит полюбить  
вам

# Помните этого парня? (seq2seq спп-rnn)

ENCODER



Are

you

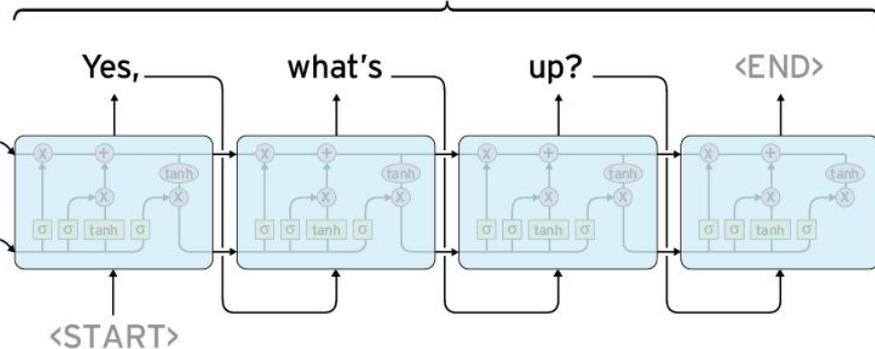
free

tomorrow?

Incoming Email

thought vector

Reply



<START>

Yes,

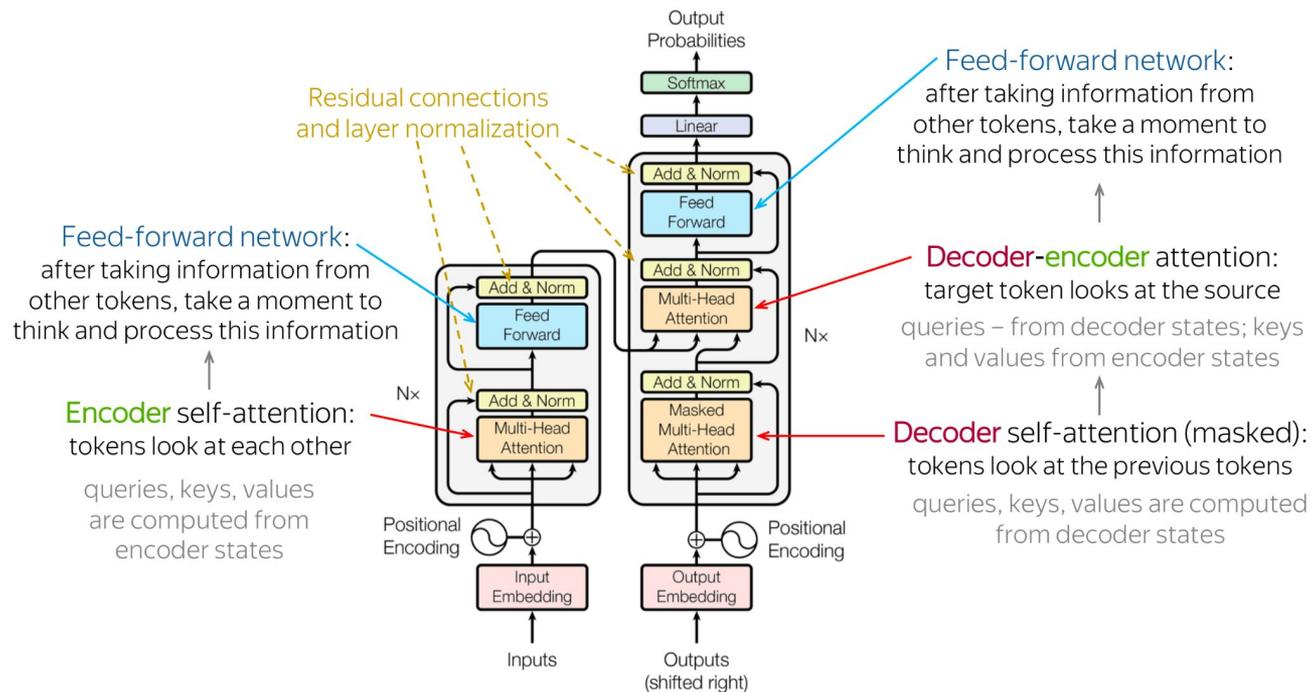
what's

up?

<END>

DECODER

# Теперь(с 2018) он выглядит так



[https://lena-voita.github.io/nlp\\_course/seq2seq\\_and\\_attention.html](https://lena-voita.github.io/nlp_course/seq2seq_and_attention.html)



build pending license Apache-2.0 website online release v2.11.0

State-of-the-art Natural Language Processing for PyTorch and TensorFlow 2.0

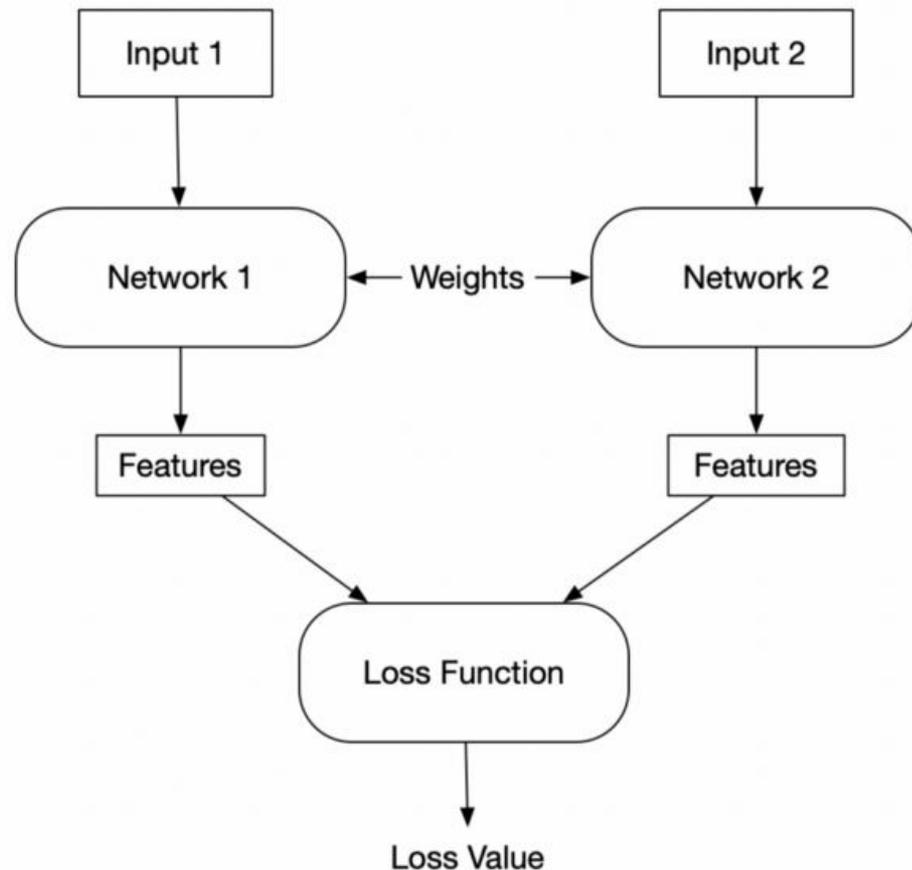
- Документация
- Курсы
- Запуск в облаке(бесплатно на небольших масштабах)
- Почти все модели работают “из коробки”

Мультиmodalность

# Сиамские сети

- Учим две сети с одним лоссом
- Нам не очень важно что у нас в разных энкодерах

Generic Siamese Model



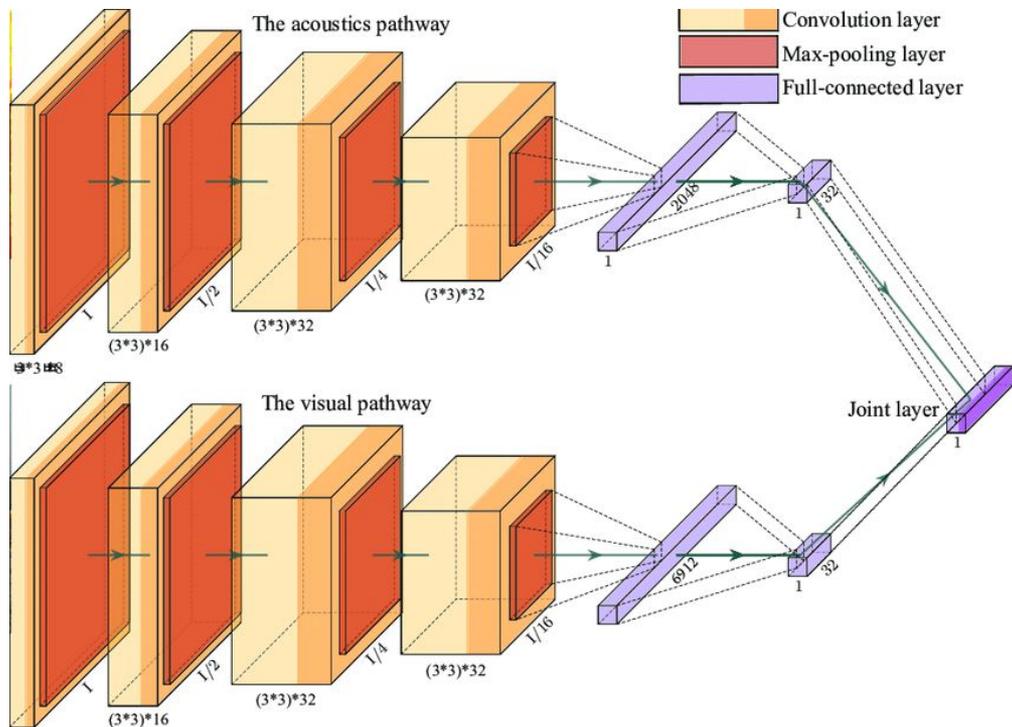
# Contrastive loss

- Максимизируем главную диагональ, нужен большой batch size для хороших моделей

$$\ell_{i,j} = -\log \frac{\exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_j)/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_k)/\tau)},$$

	T <sub>1</sub>	T <sub>2</sub>	T <sub>3</sub>	...	T <sub>N</sub>
I <sub>1</sub>	I <sub>1</sub> ·T <sub>1</sub>	I <sub>1</sub> ·T <sub>2</sub>	I <sub>1</sub> ·T <sub>3</sub>	...	I <sub>1</sub> ·T <sub>N</sub>
I <sub>2</sub>	I <sub>2</sub> ·T <sub>1</sub>	I <sub>2</sub> ·T <sub>2</sub>	I <sub>2</sub> ·T <sub>3</sub>	...	I <sub>2</sub> ·T <sub>N</sub>
I <sub>3</sub>	I <sub>3</sub> ·T <sub>1</sub>	I <sub>3</sub> ·T <sub>2</sub>	I <sub>3</sub> ·T <sub>3</sub>	...	I <sub>3</sub> ·T <sub>N</sub>
⋮	⋮	⋮	⋮	⋮	⋮
I <sub>N</sub>	I <sub>N</sub> ·T <sub>1</sub>	I <sub>N</sub> ·T <sub>2</sub>	I <sub>N</sub> ·T <sub>3</sub>	...	I <sub>N</sub> ·T <sub>N</sub>

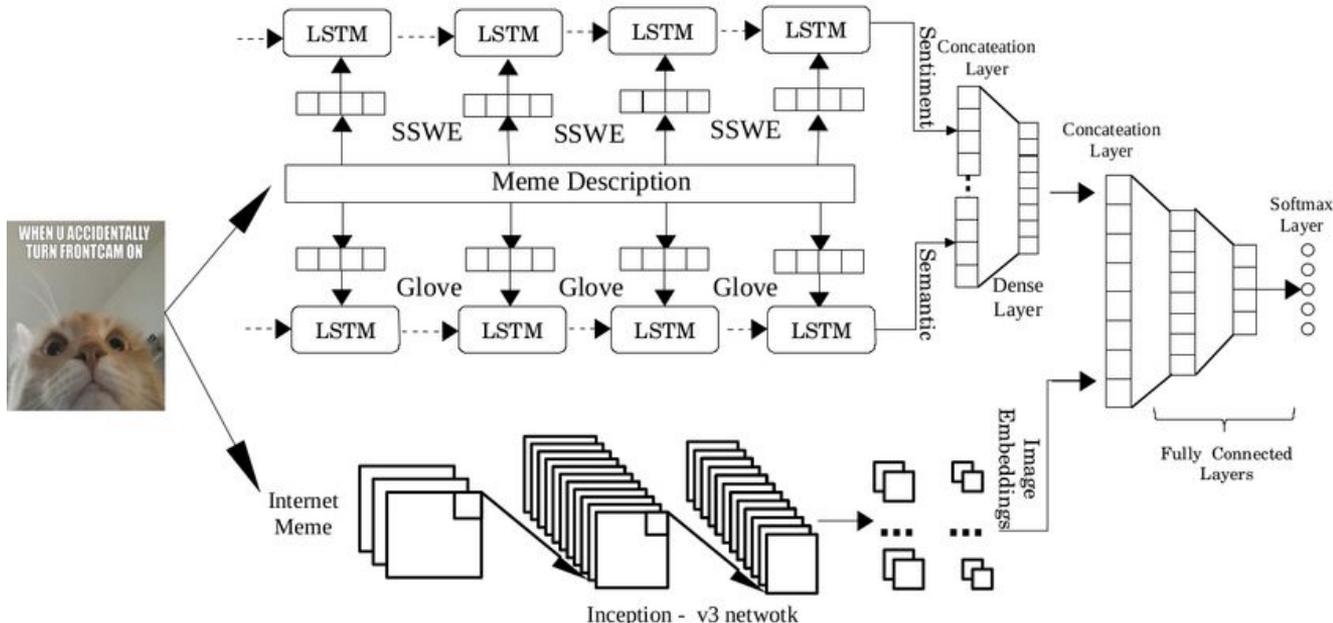
# Да кто такая эта ваша мультимодальность?



## Ключевые отличия

- Используем на каждую модальность свой энкодер
- Соединяем их (DSSM like) / кормим в декодер (Seq2Seq)

# Немного про решаемые задачи

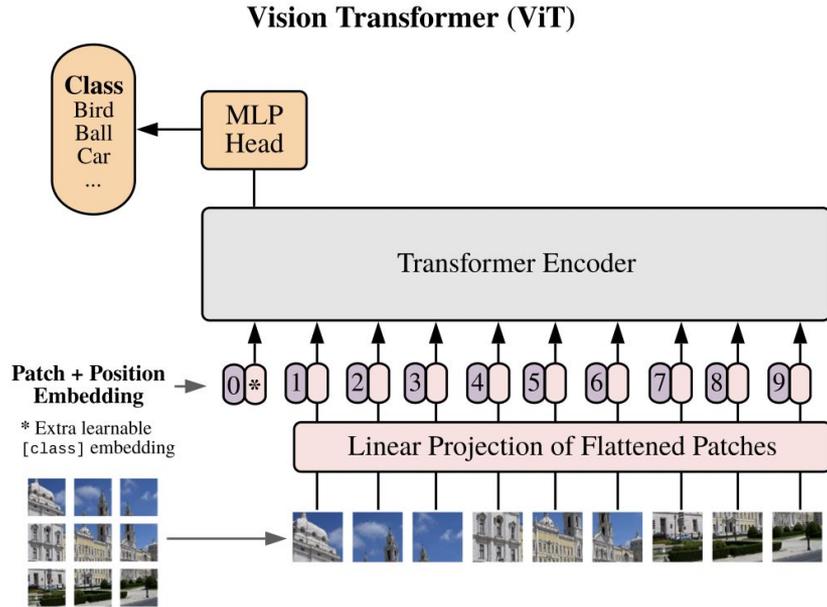


## Типы входных данных

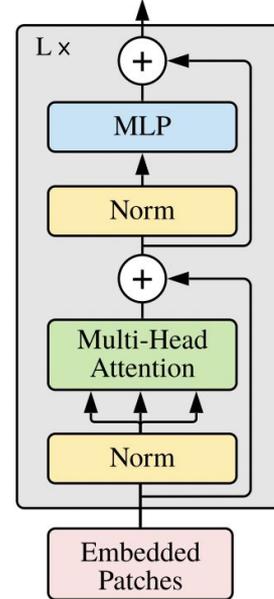
- На самом деле речь может идти и про картинки и про тексты и про видео(!)

[https://www.researchgate.net/figure/Architecture-of-Multimodal-Neural-Network-II\\_fig4\\_344603358](https://www.researchgate.net/figure/Architecture-of-Multimodal-Neural-Network-II_fig4_344603358)

# Visual transformer



## Transformer Encoder



## Идея

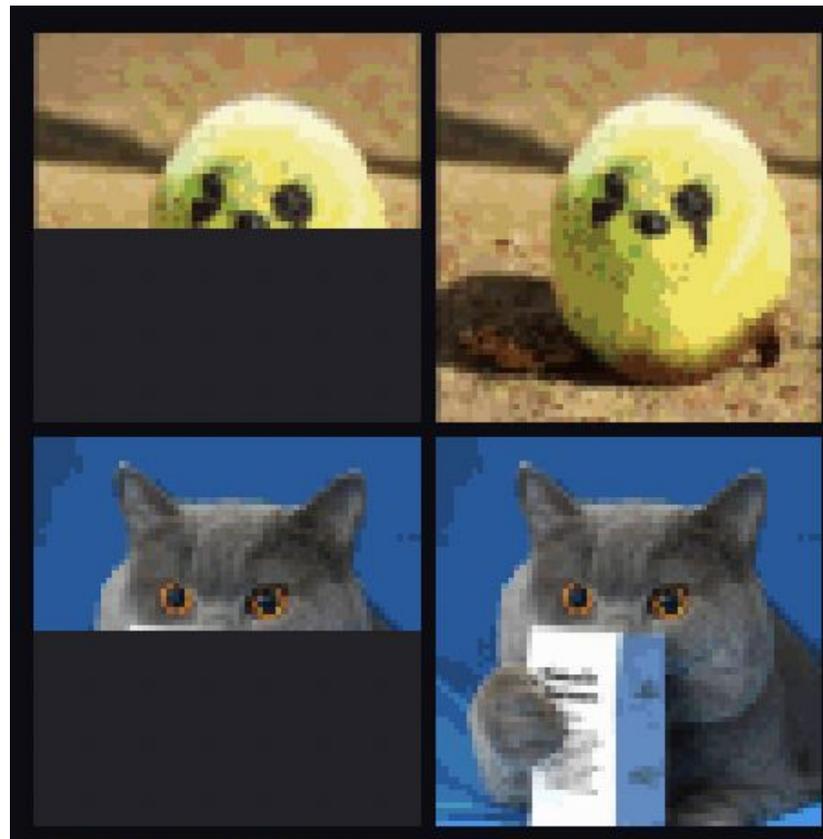
- Давайте разрежем изображение на визуальные токены ( $256*256 \rightarrow 1024$ )
- Обучим на токенах классификатор изображений

# Token based

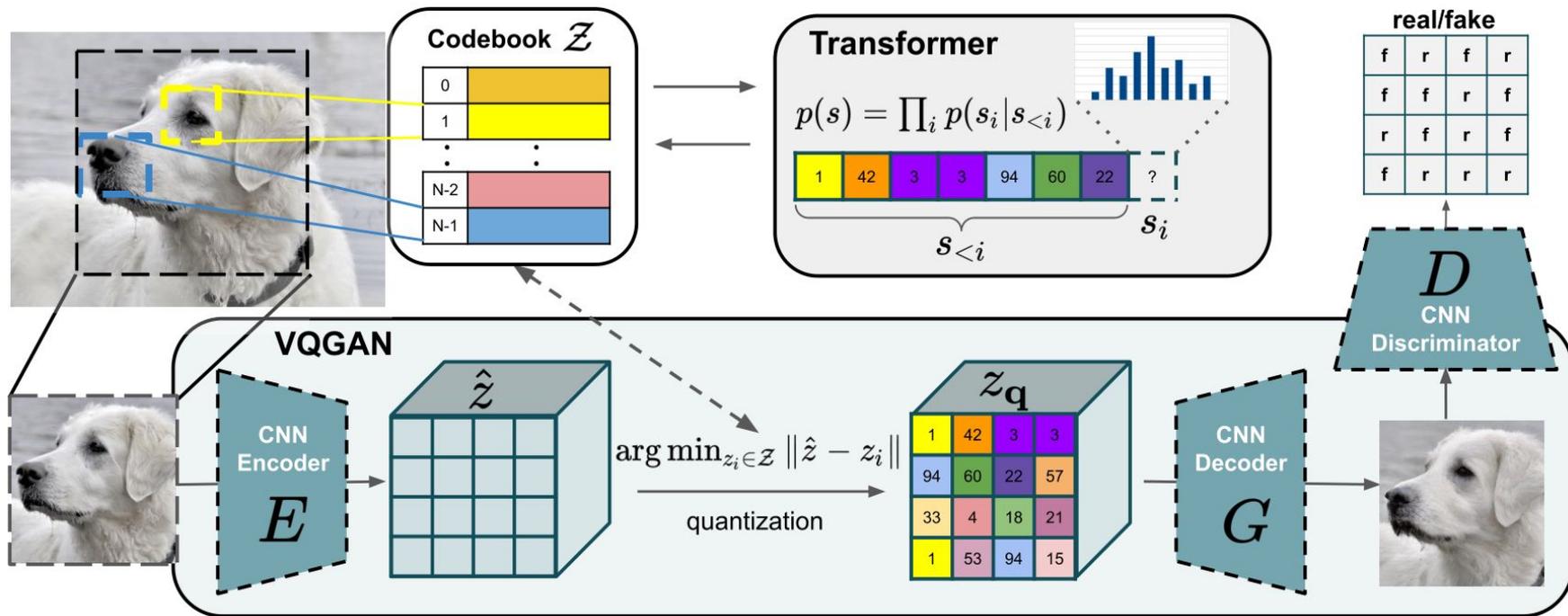
- Image GPT
- DALLE 1 like
- Rudolph
- BLIP
- Сила и слабость image tokens

# Image Gpt

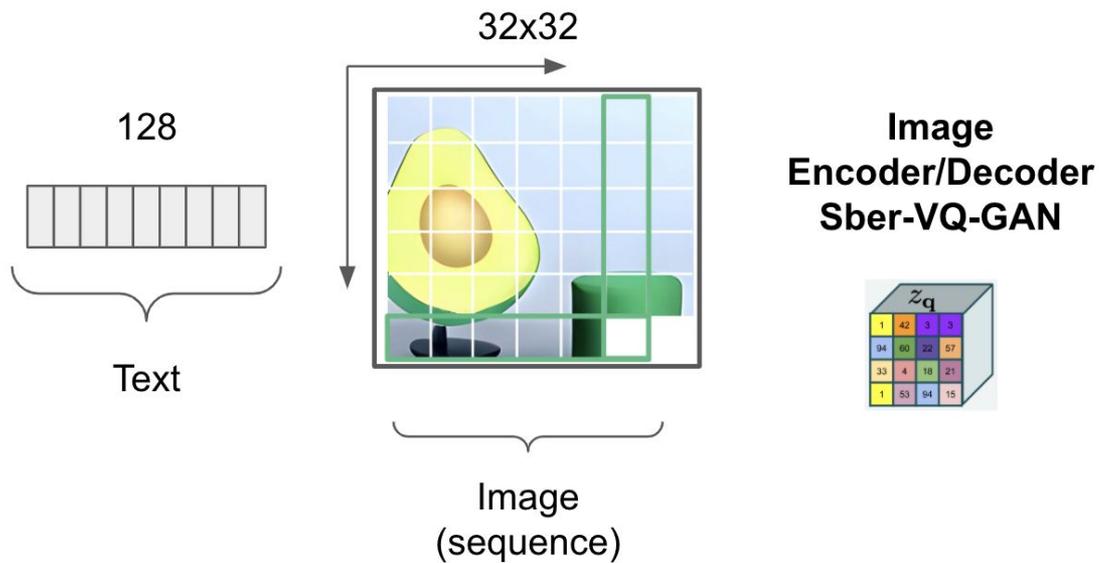
Предсказываем следующий пиксель картинки, каждый пиксель - токен.



# Vqgan



# DALLE 1

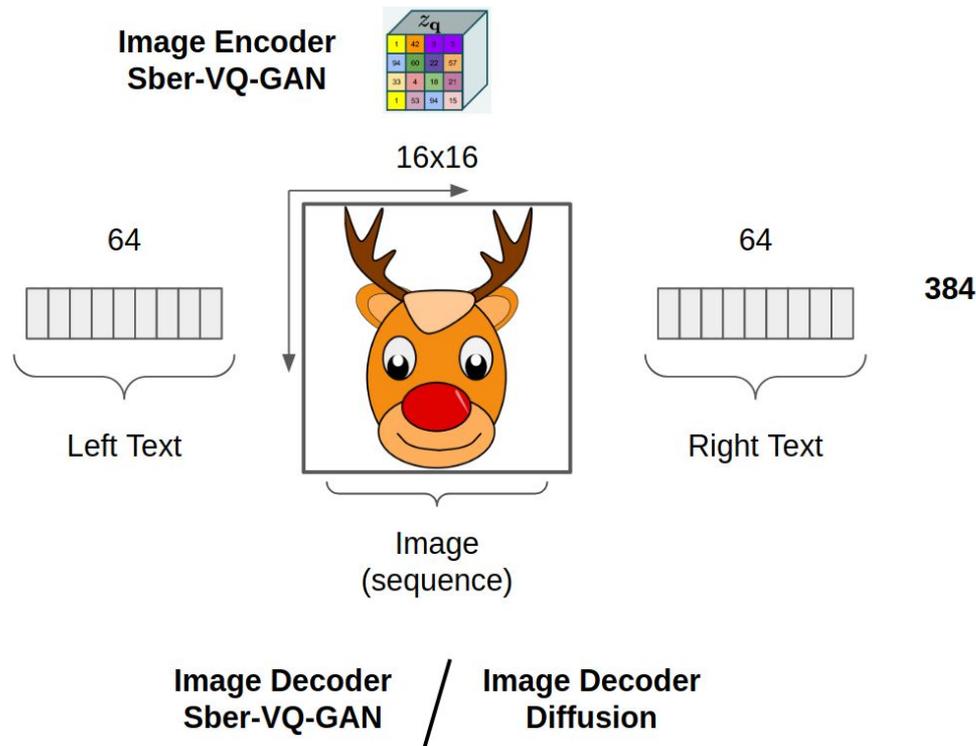


Учим GPT3 продолжать  
последовательность  
ТЕКСТОВЫХ ТОКЕНОВ,  
картиночными

<https://github.com/ai-forever/ru-dalle>

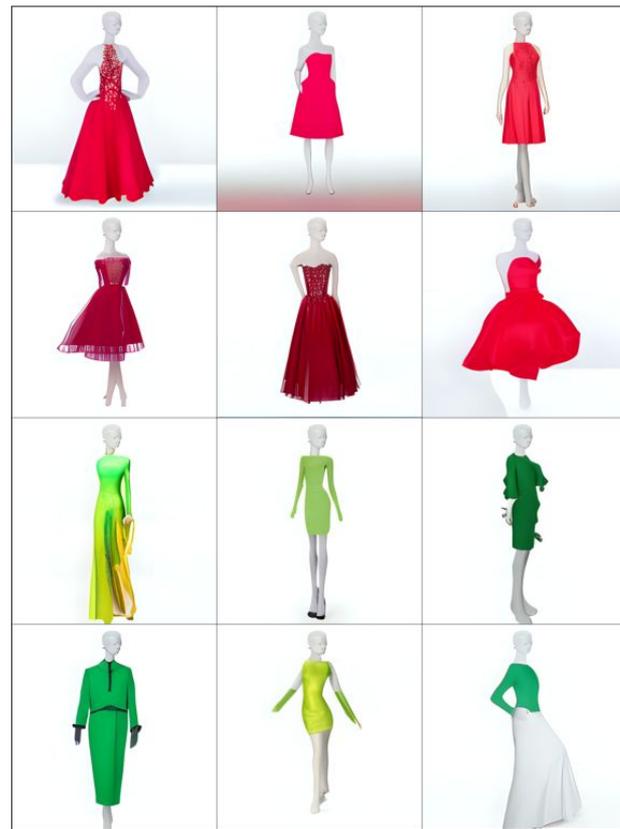
# Rudolph

Текстовые токены слева и справа, мы можем предсказывать ответы на вопросы по картинке, генерировать картинку



# Сила и слабость токенов в картинках

Так как мы оперируем картинками на уровне токенов, то мы можем дать текст, начало картинки и по текстовому запросу продолжить картинку до конца



[https://github.com/ai-forever/ru-dalle/blob/master/jupyter  
s/ruDALLE-image-prompts-dress-mannequins-V100.ipyn  
b](https://github.com/ai-forever/ru-dalle/blob/master/jupyter%20s/ruDALLE-image-prompts-dress-mannequins-V100.ipynb)

# Multi Encoder

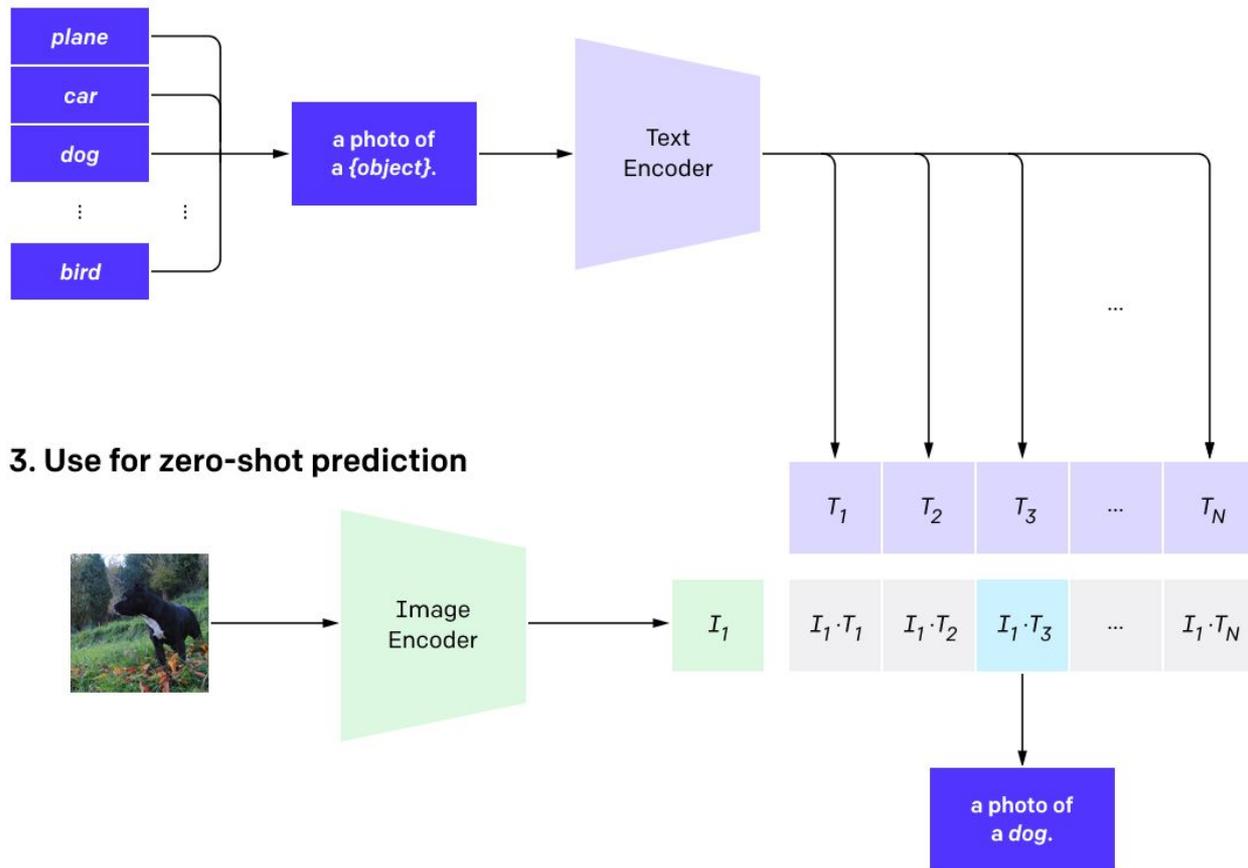
- CLIP
- KM BART
- FUSION BRAIN
- GATO
- DALLE2, ImageGen, ...

# CLIP

## Основные особенности

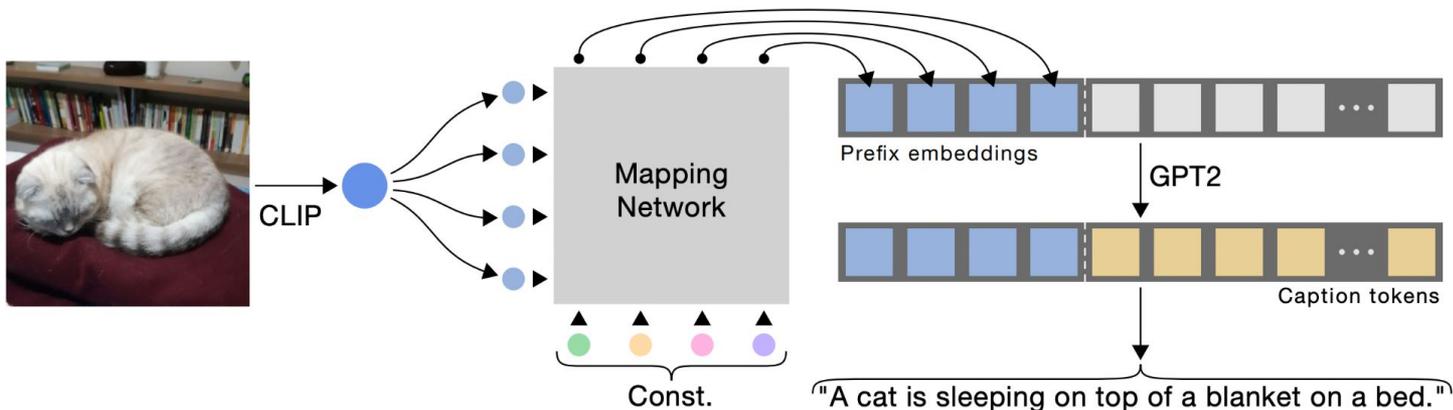
- Используем BERT для текстового энкодера
- Используем ResNet50 / ViT16 в качестве картиночного энкодера

### 2. Create dataset classifier from label text



<https://openai.com/blog/clip/>

# CLIP позволяет эффективно кодировать картинку в пространство текста



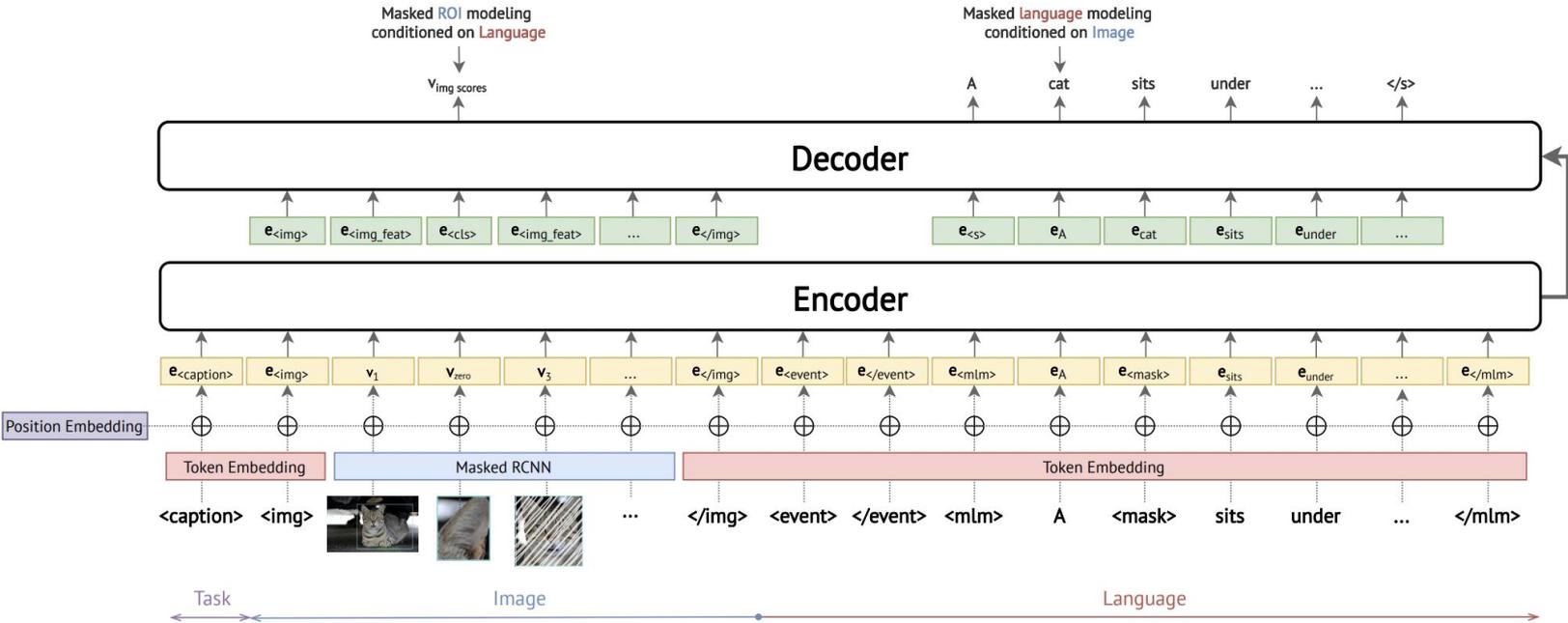
## Задачи

- Если добавить декодер(GPT), то мы получим архитектуру для решения image2text, VQA,

[https://github.com/rmokady/CLIP\\_prefix\\_caption](https://github.com/rmokady/CLIP_prefix_caption)

<https://github.com/AlexWortega/ruImageCaptioning>

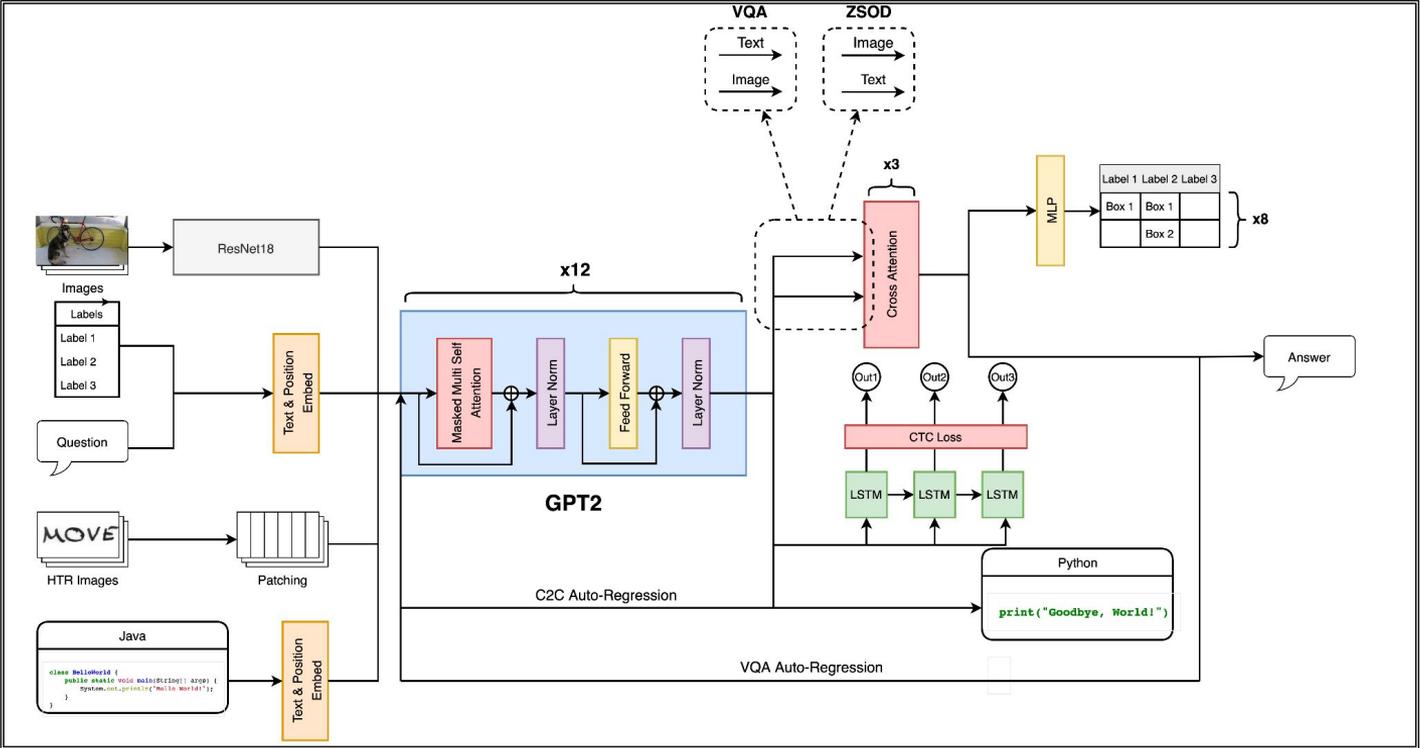
# KM BART



<https://arxiv.org/pdf/2101.00419.pdf>

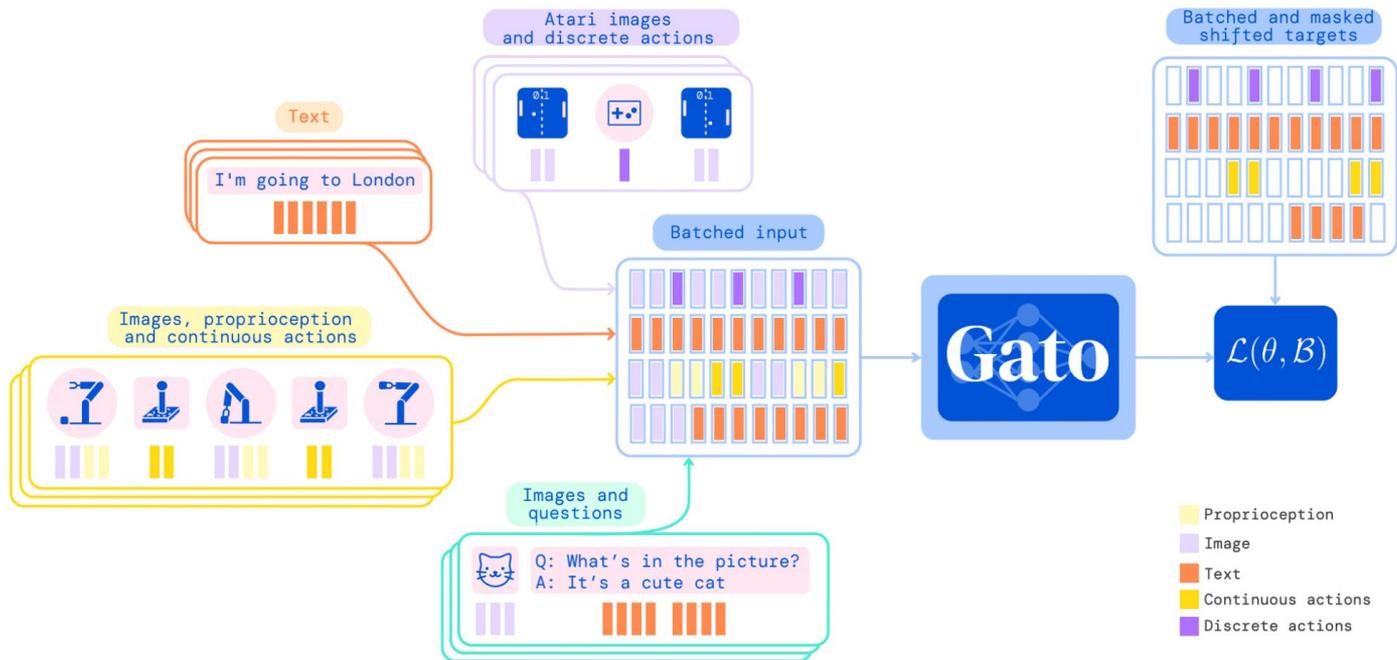
<https://github.com/FomalhautB/KM-BART>

# Fusion Brain



[https://github.com/ai-forever/fusion\\_brain\\_aig2021](https://github.com/ai-forever/fusion_brain_aig2021)

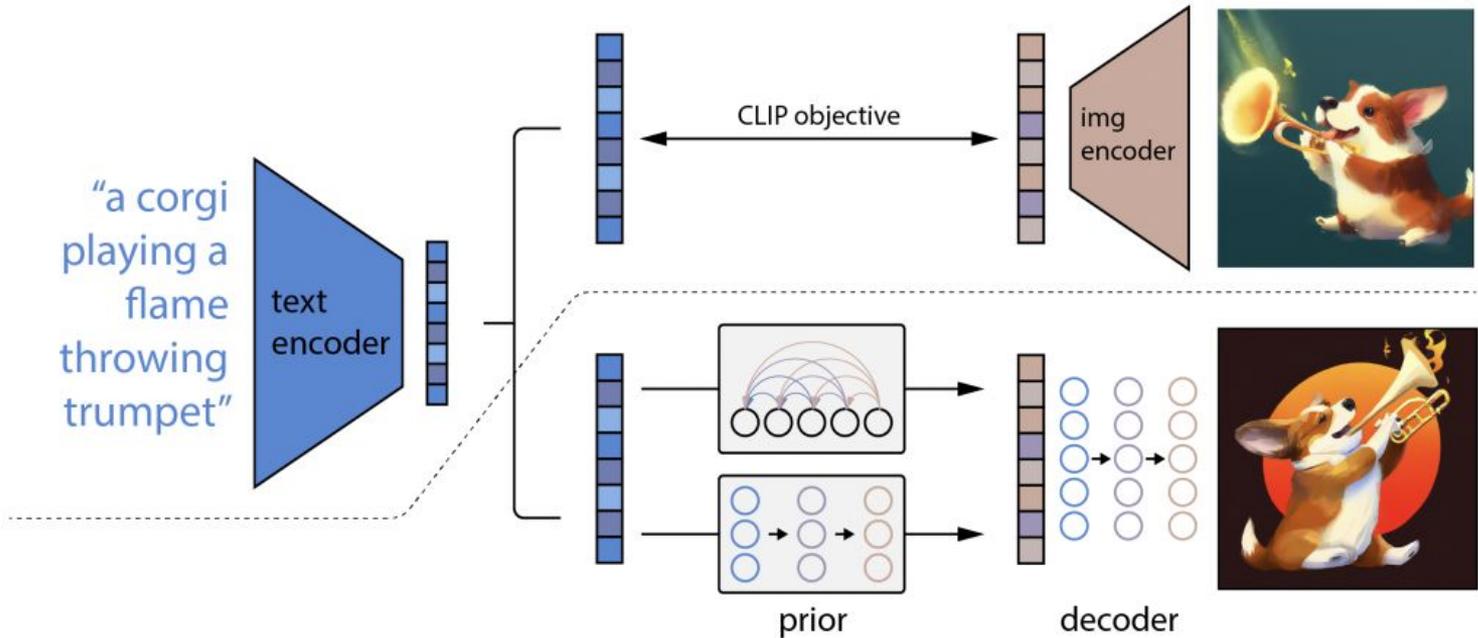
# GATO



## идея

- берем много задач и учим GPT like

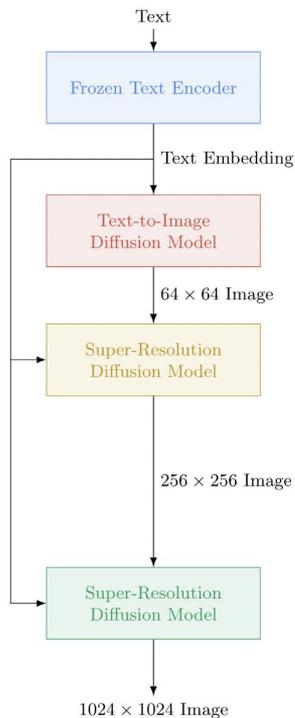
# DALLE 2



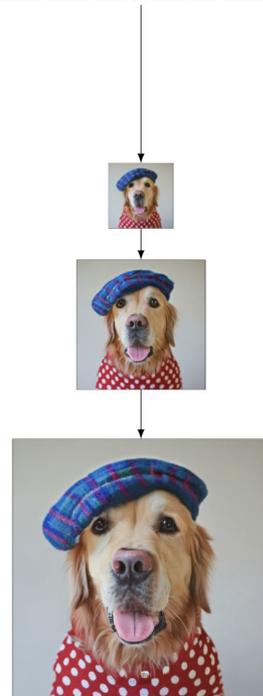
<https://cdn.openai.com/papers/dall-e-2.pdf>

# ImageGen

Используем энкодер T5  
вместо CLIP text encoder



“A Golden Retriever dog wearing a blue checkered beret and red dotted turtleneck.”



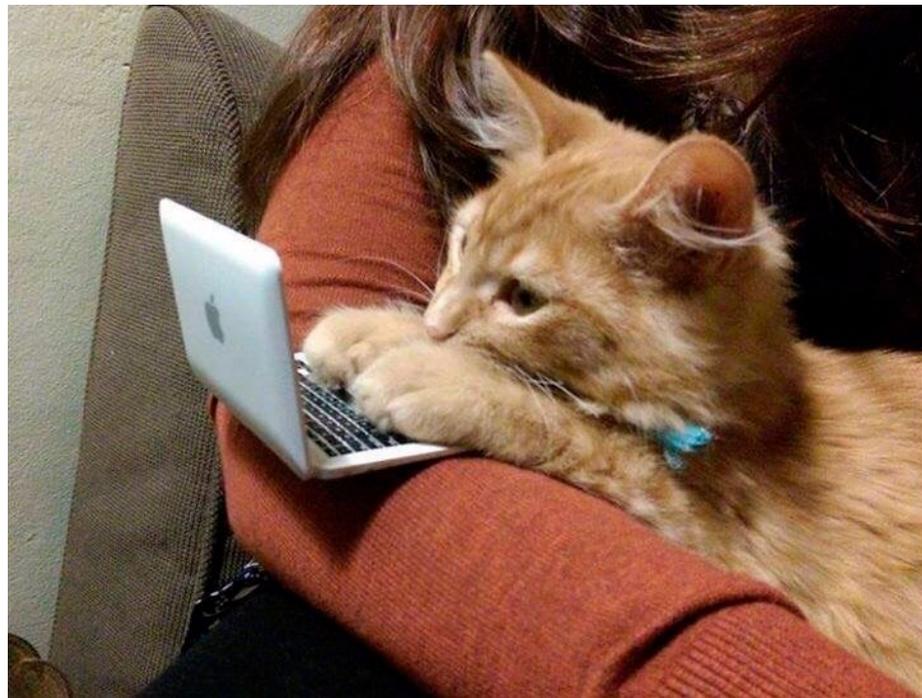
Спасибо за внимание!

# Мои контакты



<https://t.me/Alexwortegae>

Давайте что нибудь потыкаем



[https://colab.research.google.com/drive/12YRRzhl5cHER\\_U2F-buQxif8GIhMPWq3?usp=sharing#scrollTo=OJTpZacBn0OI](https://colab.research.google.com/drive/12YRRzhl5cHER_U2F-buQxif8GIhMPWq3?usp=sharing#scrollTo=OJTpZacBn0OI)