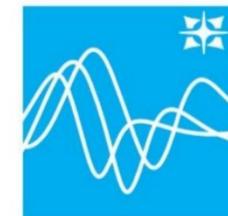


# Приложения моделей Transformer

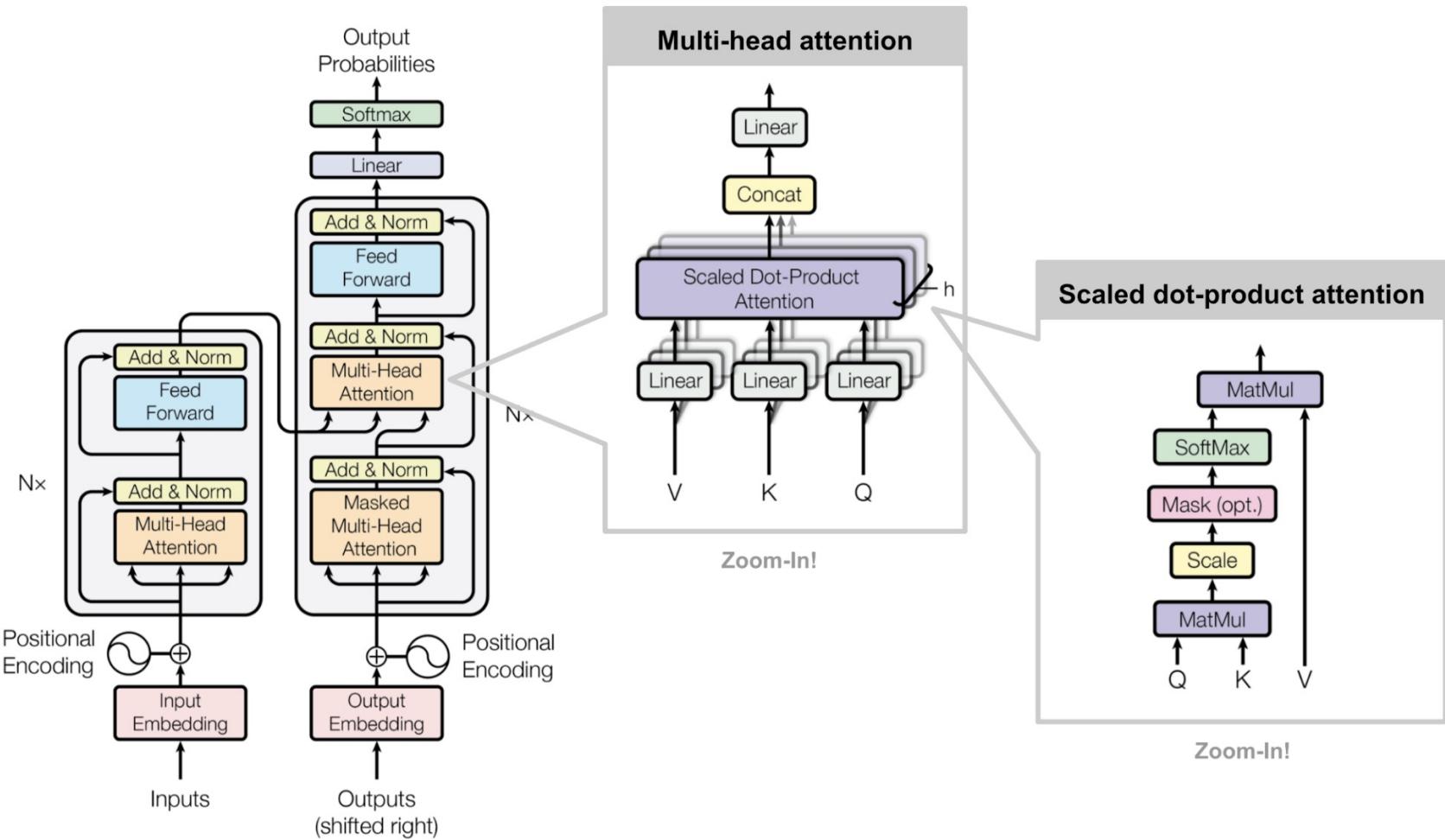


Кафедра  
технологий  
проектирования  
сложных  
технических  
систем

# План лекции

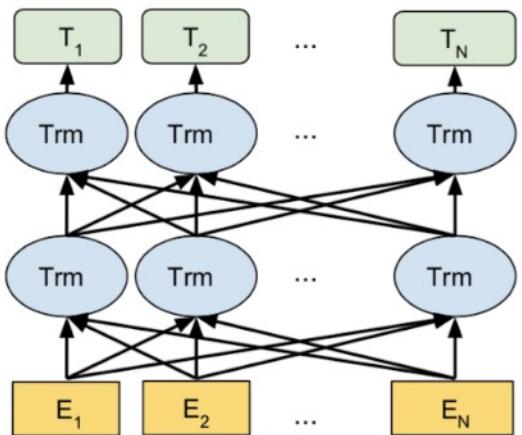
- Attention (повторение)
- CLIP framework
  - Transformer for images
- Question answering (QA)
  - Text generation
  - Close & Open domain

# Transformer

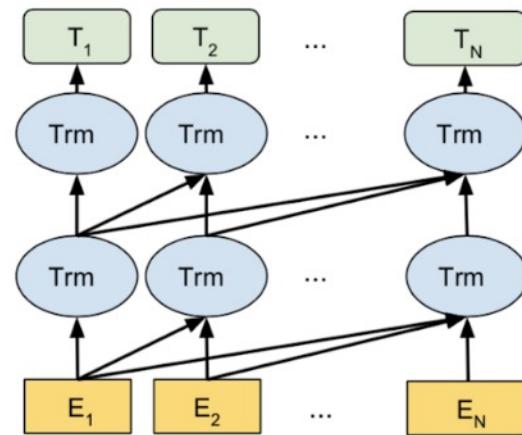


# Transformer basis

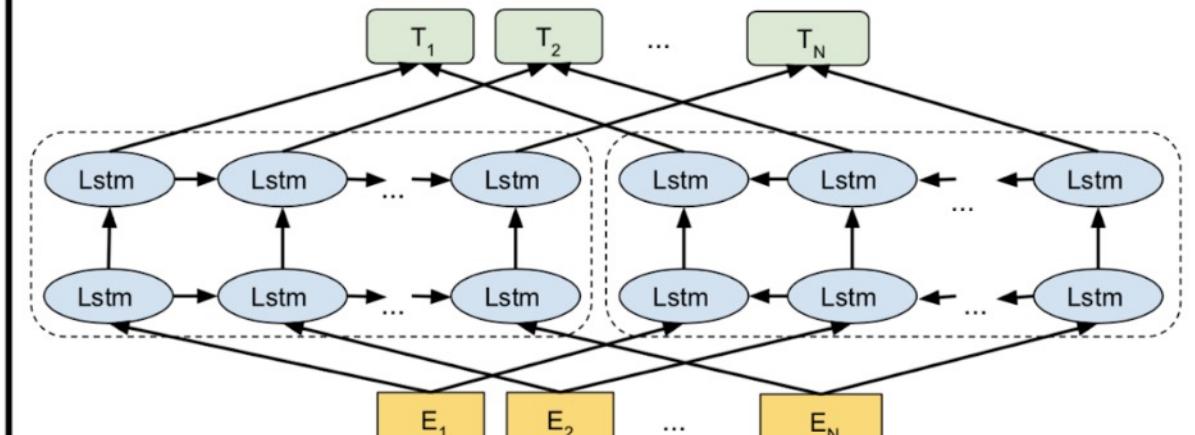
BERT (Ours)



OpenAI GPT



ELMo



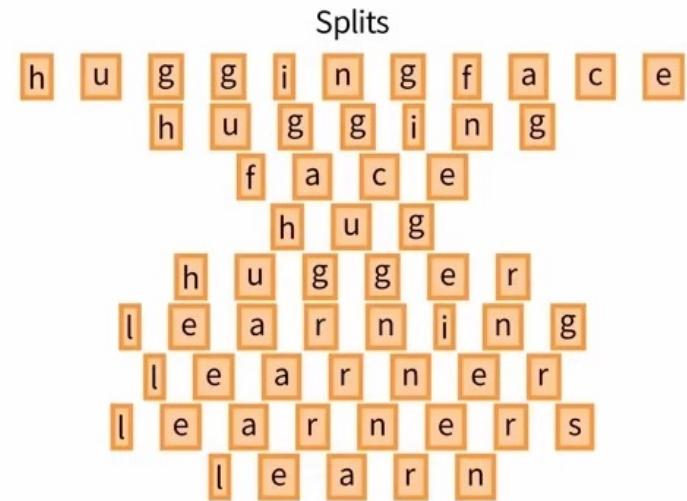
# Byte Pair Encoding

this is the hugging face course. this chapter is about tokenization. this section shows several tokenizer algorithms.

this is the hugging face course . this chapter is about tokenization . this section shows several tokenizer algorithms .

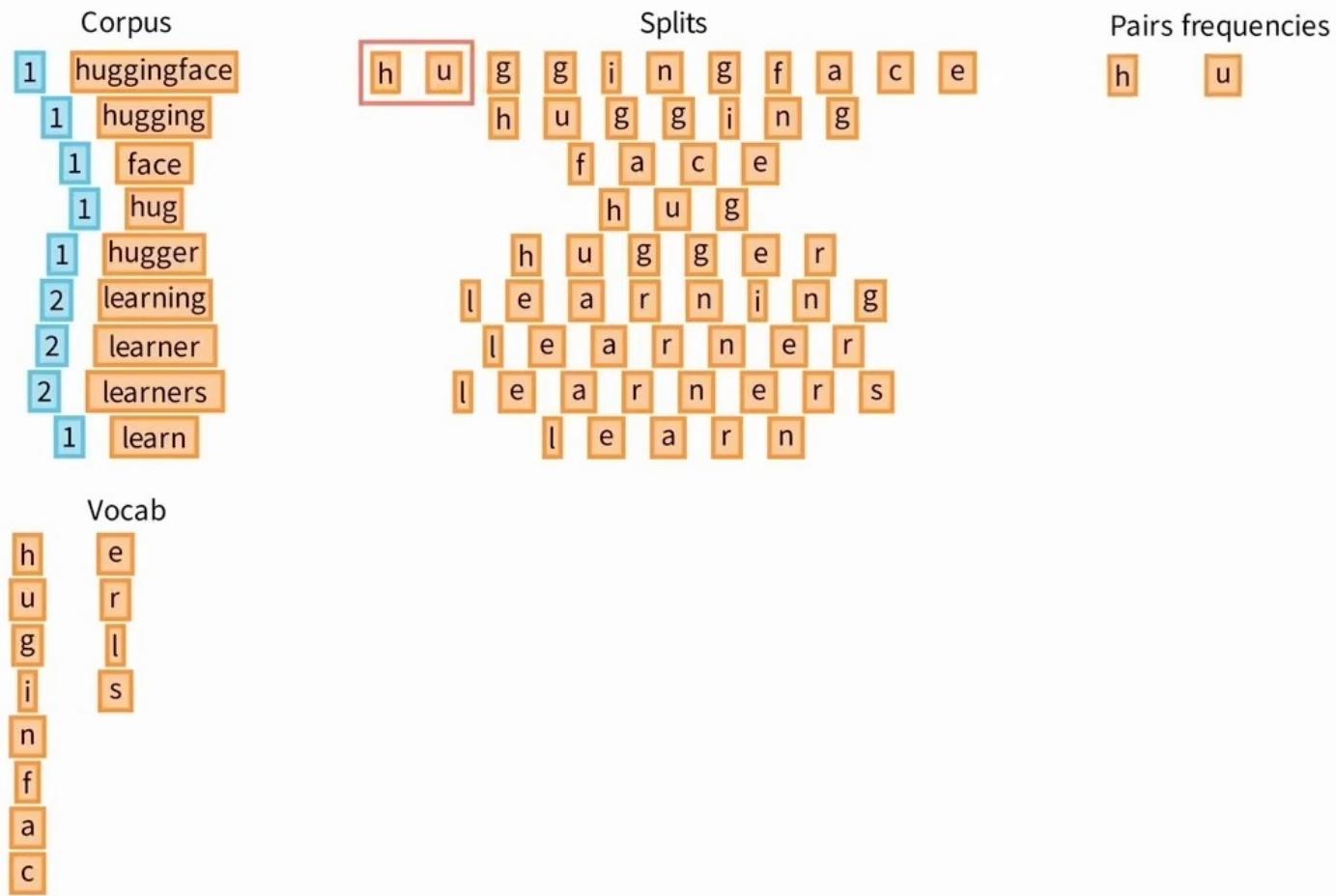
3x this 2x is 1x the  
1x hugging 1x face  
1x course 3x . 1x chapter  
1x about 1x tokenization  
1x section 1x shows  
1x several 1x tokenizer  
1x algorithms

Corpus	
1	huggingface
1	hugging
1	face
1	hug
1	hugger
2	learning
2	learner
2	learners
1	learn

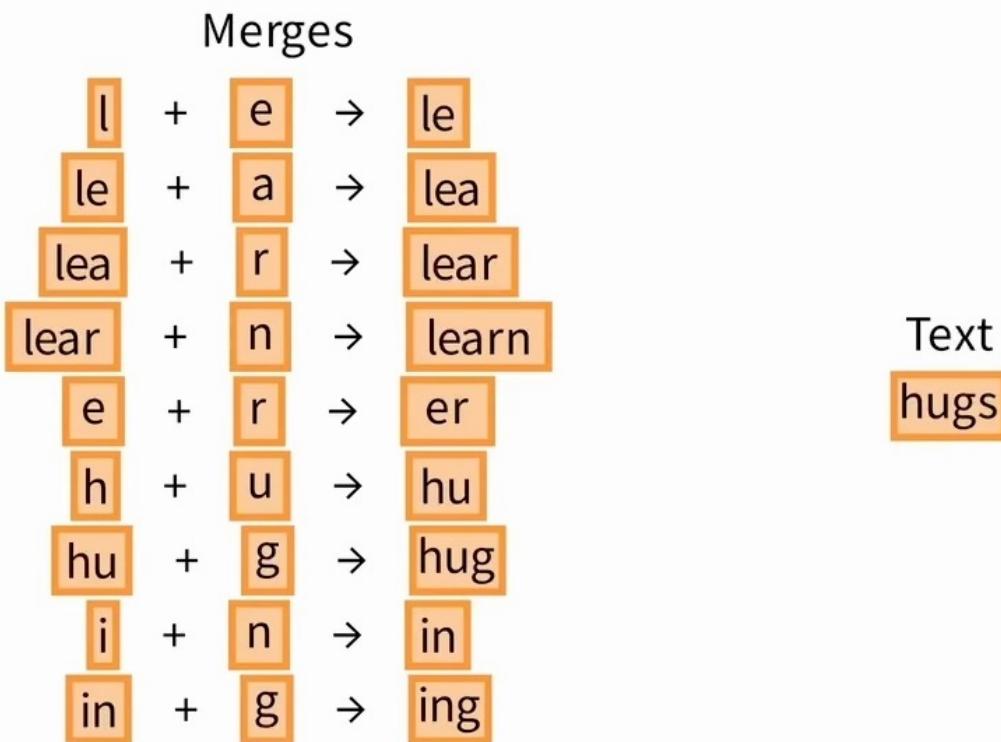


Vocab	
h	e
u	r
g	i
i	s
n	f
f	a
a	c

# Byte Pair Encoding

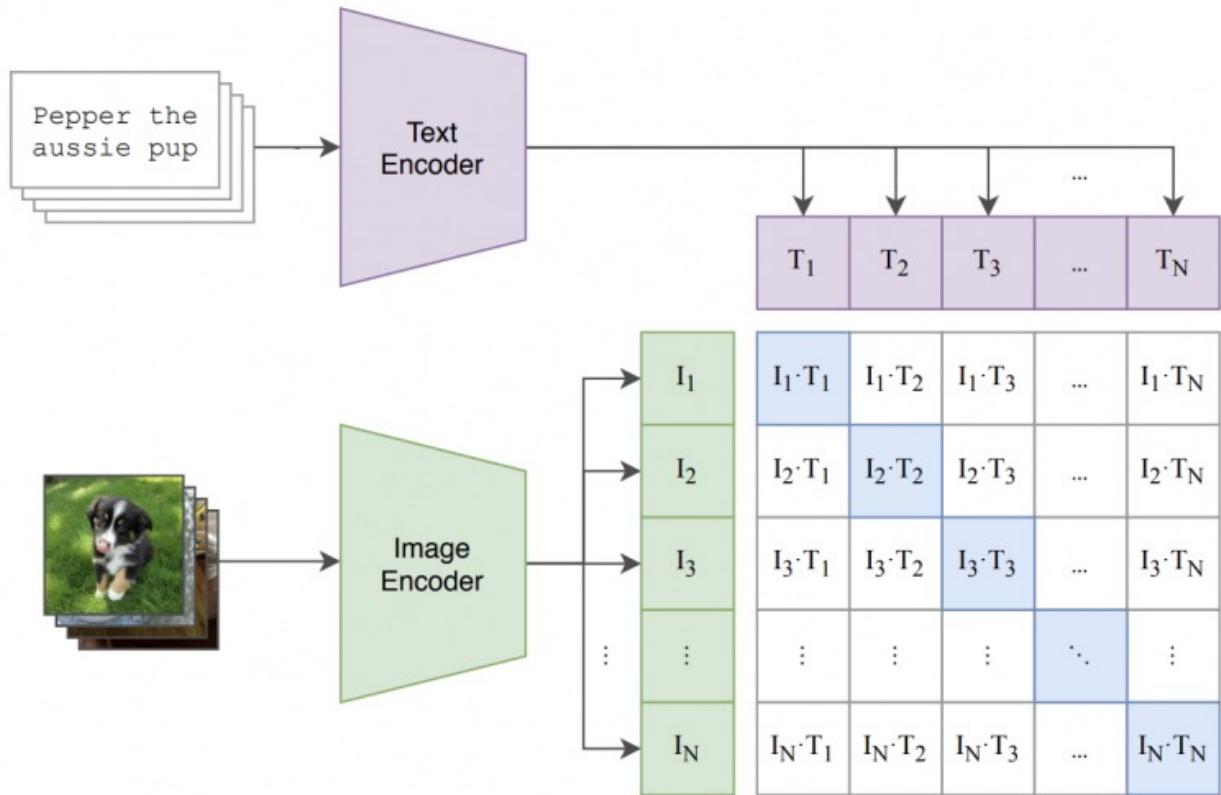


# Byte Pair Encoding

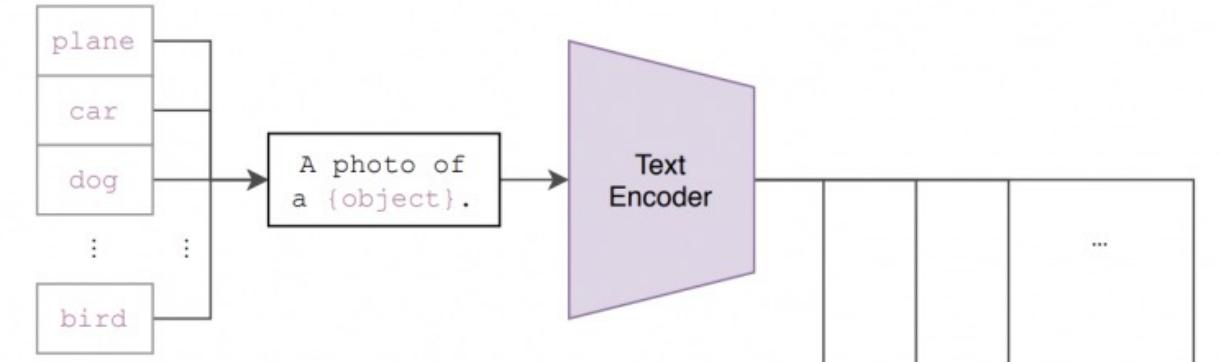


# CLIP

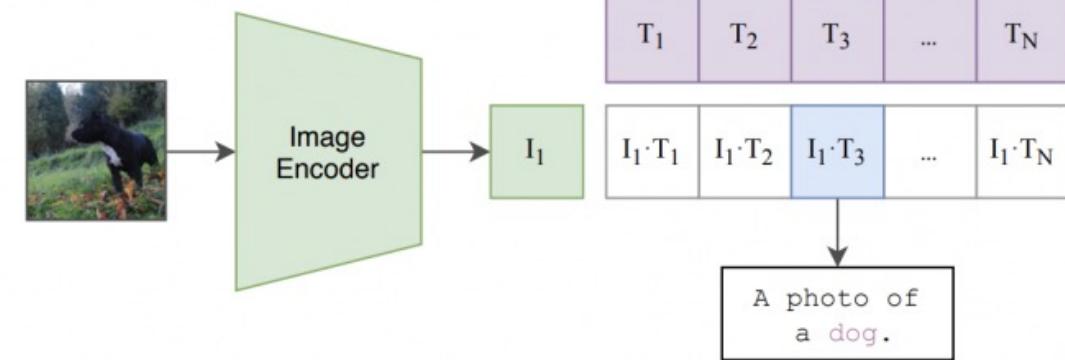
(1) Contrastive pre-training



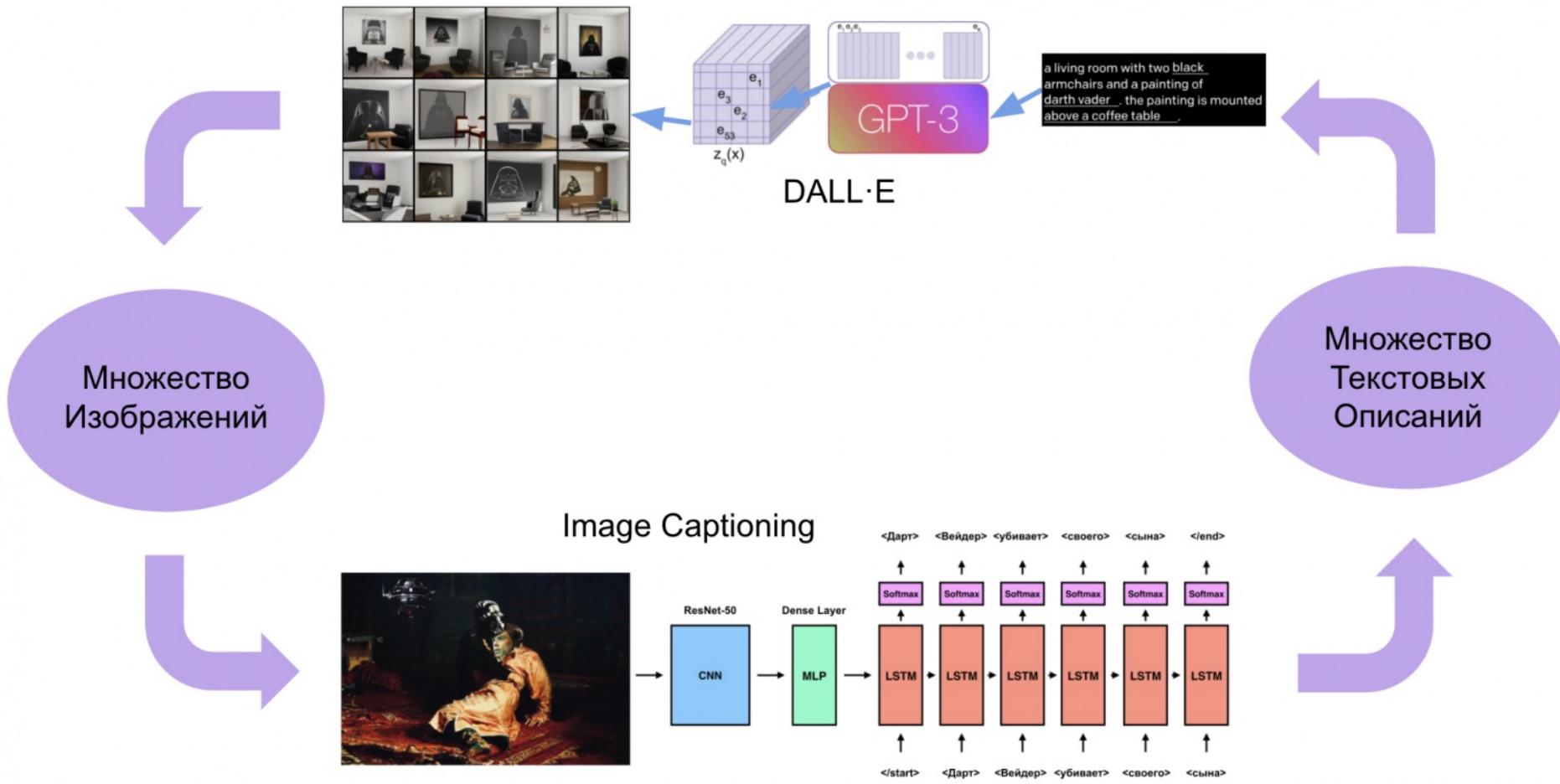
(2) Create dataset classifier from label text



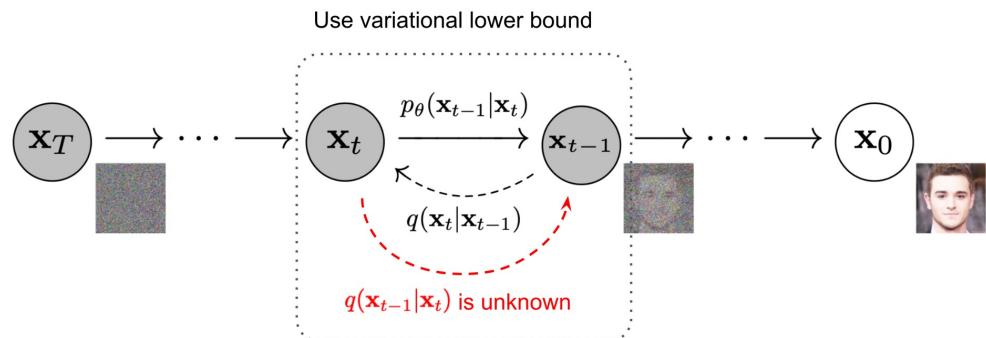
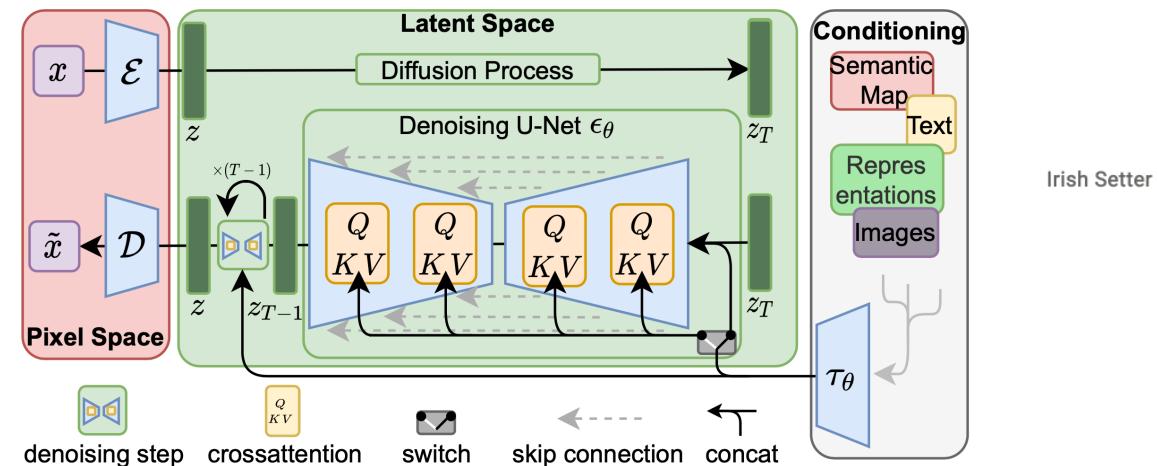
(3) Use for zero-shot prediction



# CLIP + Image generation



# Diffusion models (recap)



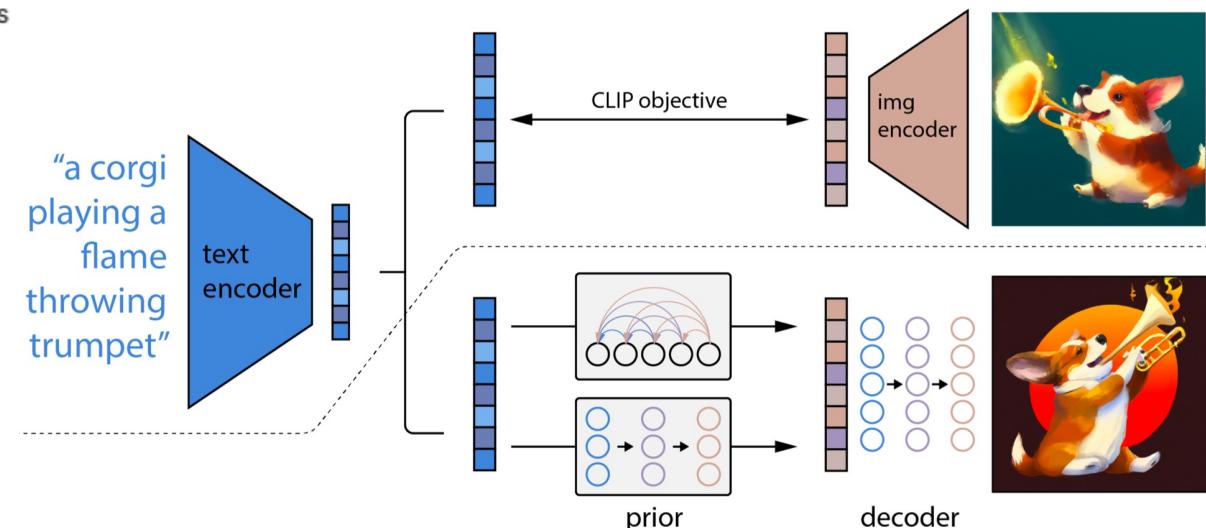
# Promt Image Generation

The two-stage diffusion model **unCLIP** (Ramesh et al. 2022) heavily utilizes the CLIP text encoder to produce text-guided images at high quality. Given a pretrained CLIP model  $\mathbf{c}$  and paired training data for the diffusion model,  $(\mathbf{x}, y)$ , where  $\mathbf{x}$  is an image and  $y$  is the corresponding caption, we can compute the CLIP text and image embedding,  $\mathbf{c}^t(y)$  and  $\mathbf{c}^i(\mathbf{x})$ , respectively. The unCLIP learns two models in parallel:

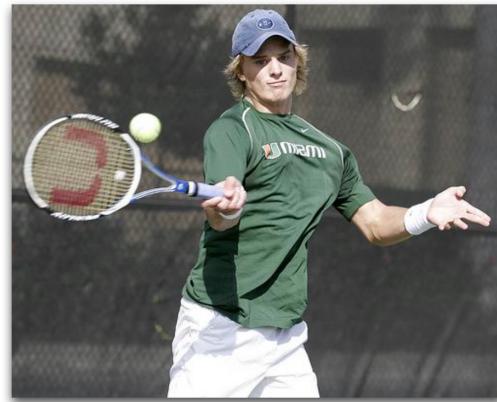
- A prior model  $P(\mathbf{c}^i|y)$ : outputs CLIP image embedding  $\mathbf{c}^i$  given the text  $y$ .
- A decoder  $P(\mathbf{x}|\mathbf{c}^i, [y])$ : generates the image  $\mathbf{x}$  given CLIP image embedding  $\mathbf{c}^i$  and optionally the original text  $y$ .

unCLIP follows a two-stage image generation process:

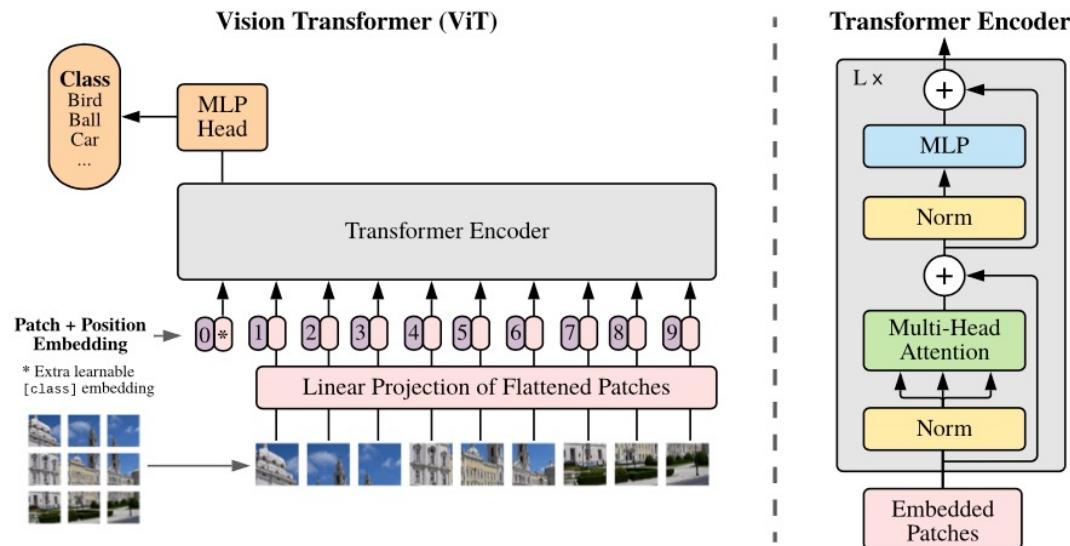
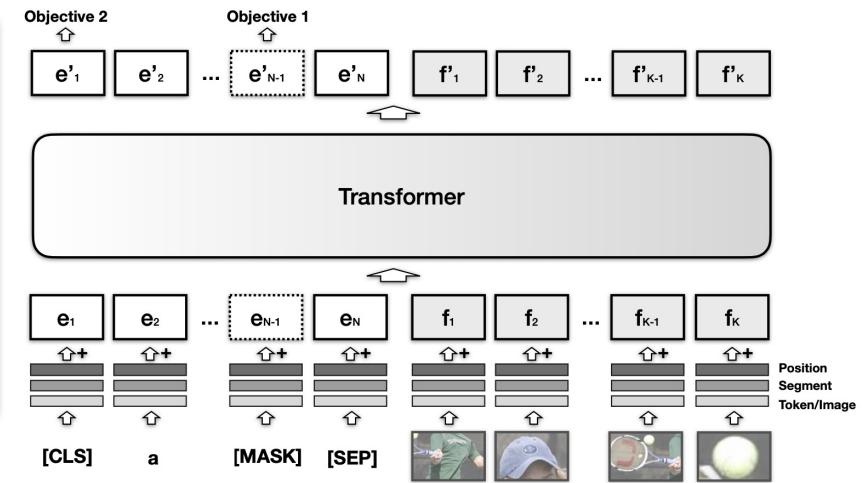
1. Given a text  $y$ , a CLIP model is first used to generate a text embedding  $\mathbf{c}^t(y)$ . Using CLIP latent space enables zero-shot image manipulation via text.
2. A diffusion or autoregressive prior  $P(\mathbf{c}^i|y)$  processes this CLIP text embedding to construct an image prior and then a diffusion decoder  $P(\mathbf{x}|\mathbf{c}^i, [y])$  generates an image, conditioned on the prior. This decoder can also generate image variations conditioned on an image input, preserving its style and semantics.



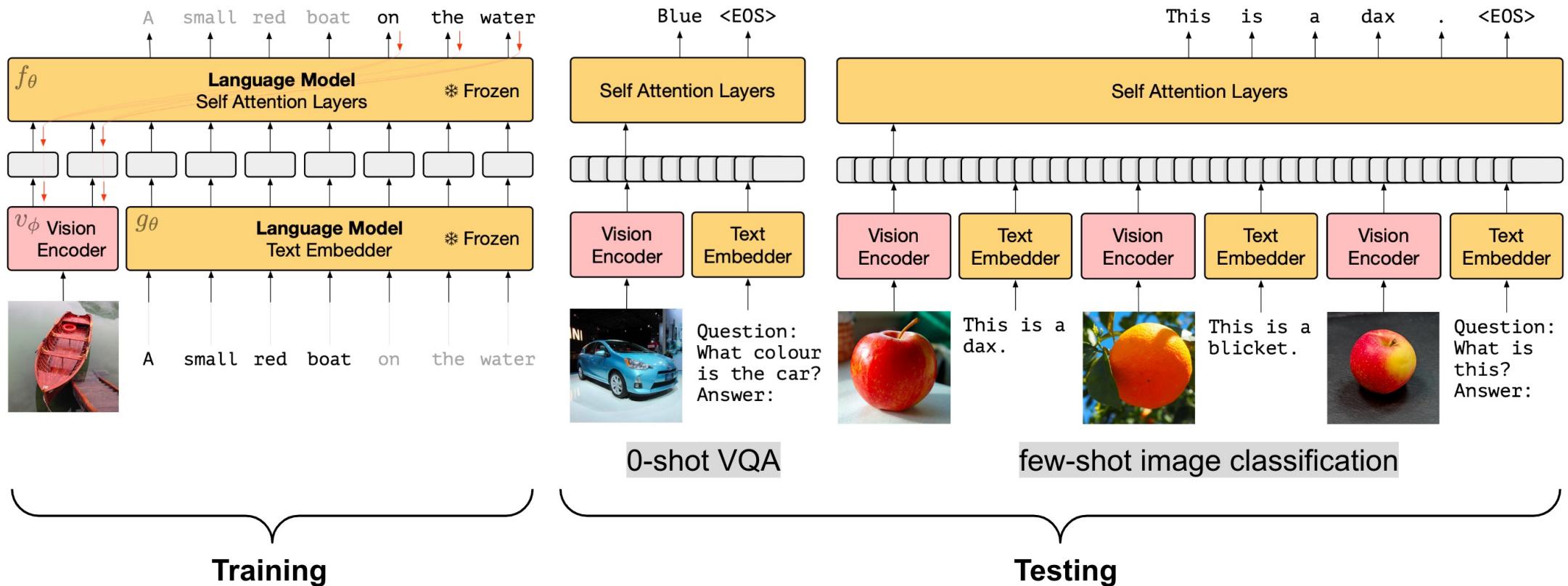
# Visual Transformers



A person hits a ball with a tennis racket



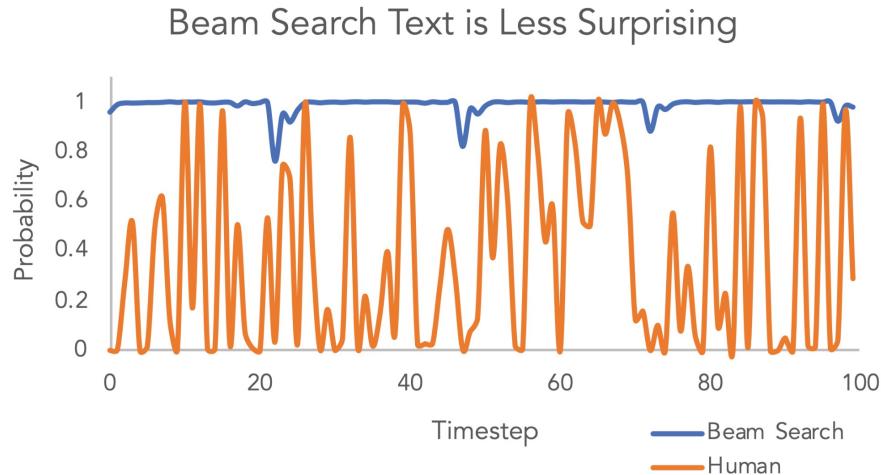
# Learned Image Embedding



# Question & Answering

## Text generation

$$p_i \propto \frac{\exp(o_i/T)}{\sum_j \exp(o_j/T)}$$



- Greedy search
- Beam search
- Top-k sampling

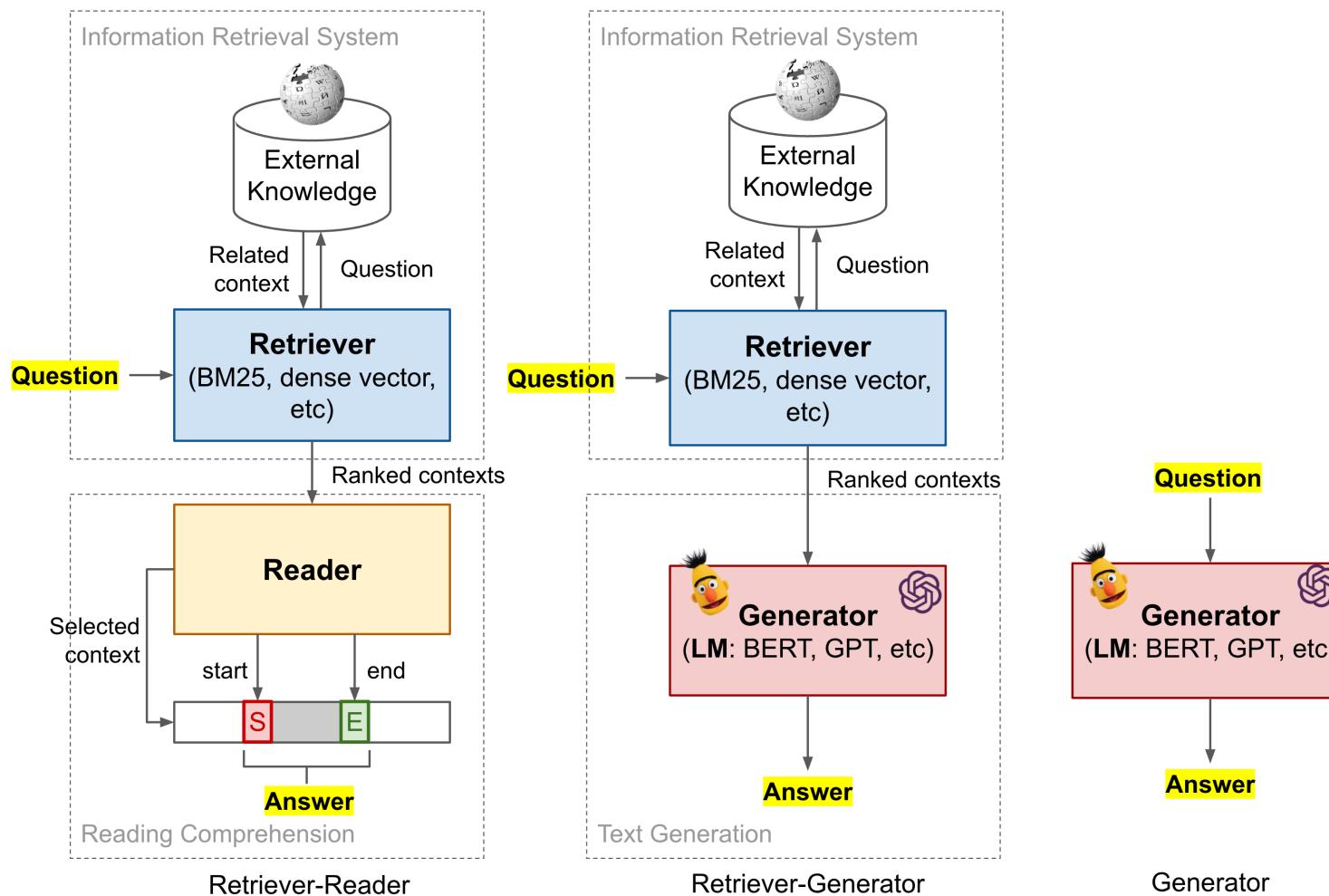
## Guided Decoding

$$\mathbf{y}^* = \arg \max_{\mathbf{y} \in \mathcal{Y}} \left( \underbrace{\log p_\theta(\mathbf{y}|\mathbf{x})}_{\text{MAP}} - \underbrace{\lambda \mathcal{R}(\mathbf{y})}_{\text{regularizer}} \right)$$

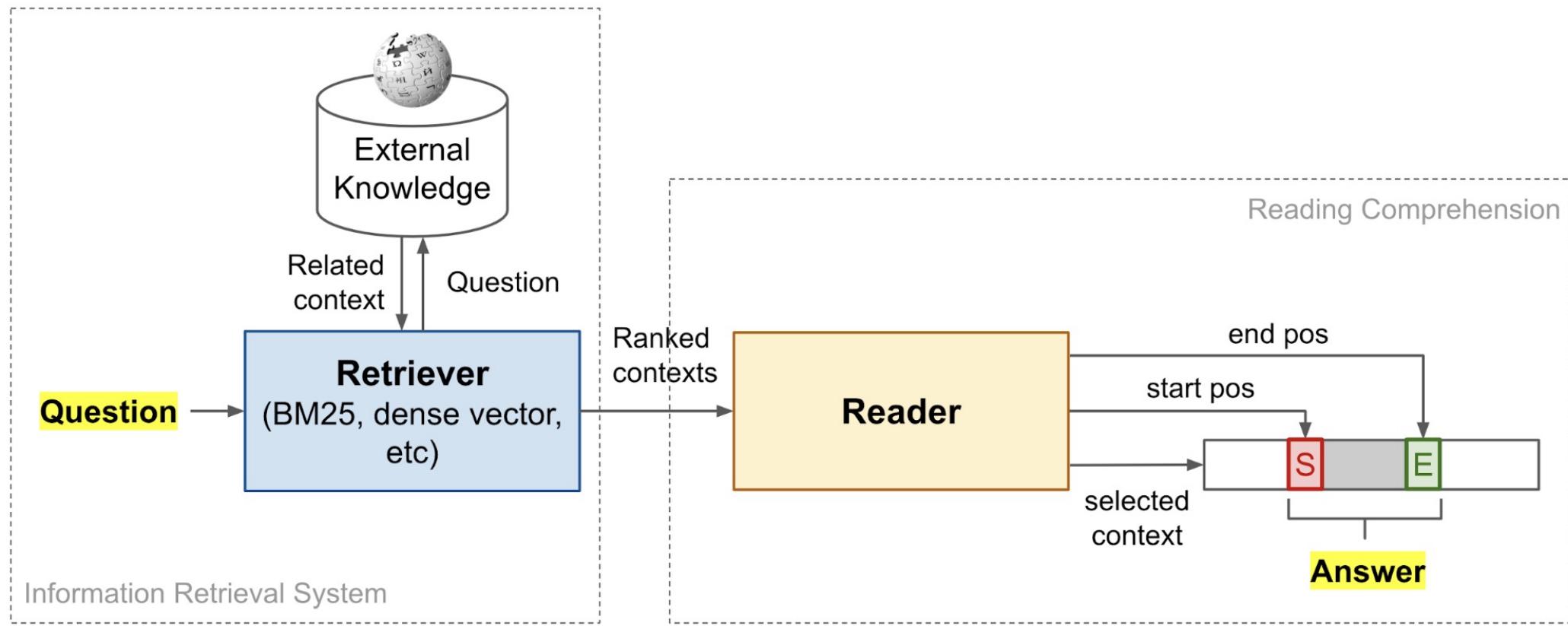
1. *Greedy*:  $\mathcal{R}_{\text{greedy}}(\mathbf{y}) = \sum_{t=1}^{|\mathbf{y}|} (u_t(y_t) - \min_{y' \in \mathcal{V}} u_t(y'))^2$ ; if set  $\lambda \rightarrow \infty$ , we have greedy search. Note that being greedy at each individual step does not guarantee global optimality.
2. *Variance regularizer*:  $\mathcal{R}_{\text{var}}(\mathbf{y}) = \frac{1}{|\mathbf{y}|} \sum_{t=1}^{|\mathbf{y}|} (u_t(y_t) - \bar{u})^2$ , where  $\bar{u}$  is the average surprisal over all timesteps. It directly encodes the UID hypothesis.
3. *Local consistency*:  $\mathcal{R}_{\text{local}}(\mathbf{y}) = \frac{1}{|\mathbf{y}|} \sum_{t=1}^{|\mathbf{y}|} (u_t(y_t) - u_{t-1}(y_{t-1}))^2$ ; this decoding regularizer encourages adjacent tokens to have similar surprisal.
4. *Max regularizer*:  $\mathcal{R}_{\text{max}}(\mathbf{y}) = \max_t u_t(y_t)$  penalizes the maximum compensation of surprisal.
5. *Squared regularizer*:  $\mathcal{R}_{\text{square}}(\mathbf{y}) = \sum_{t=1}^{|\mathbf{y}|} u_t(y_t)^2$  encourages all the tokens to have surprisal close to 0.

"The uniform information density hypothesis (UID; Levy and Jaeger, 2007) states that—subject to the constraints of the grammar—humans prefer sentences that distribute information (in the sense of information theory) equally across the linguistic signal, e.g., a sentence."

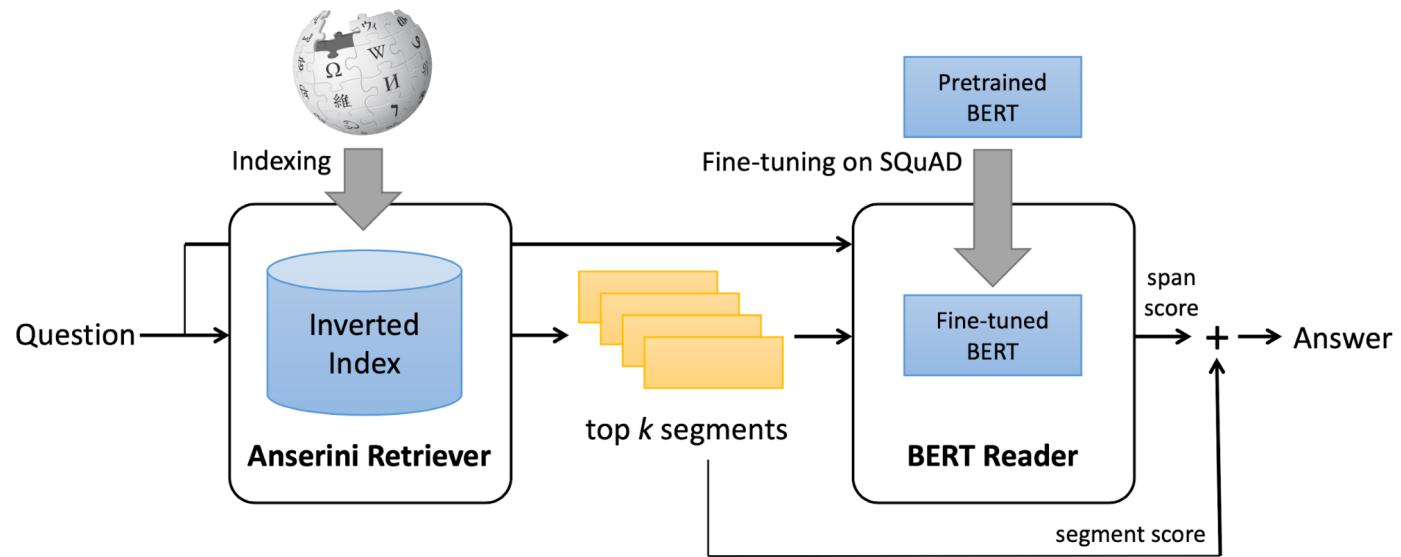
# Question & Answering



# Open-book QA: Retriever-Reader



# Open-book QA: Retriever-Reader



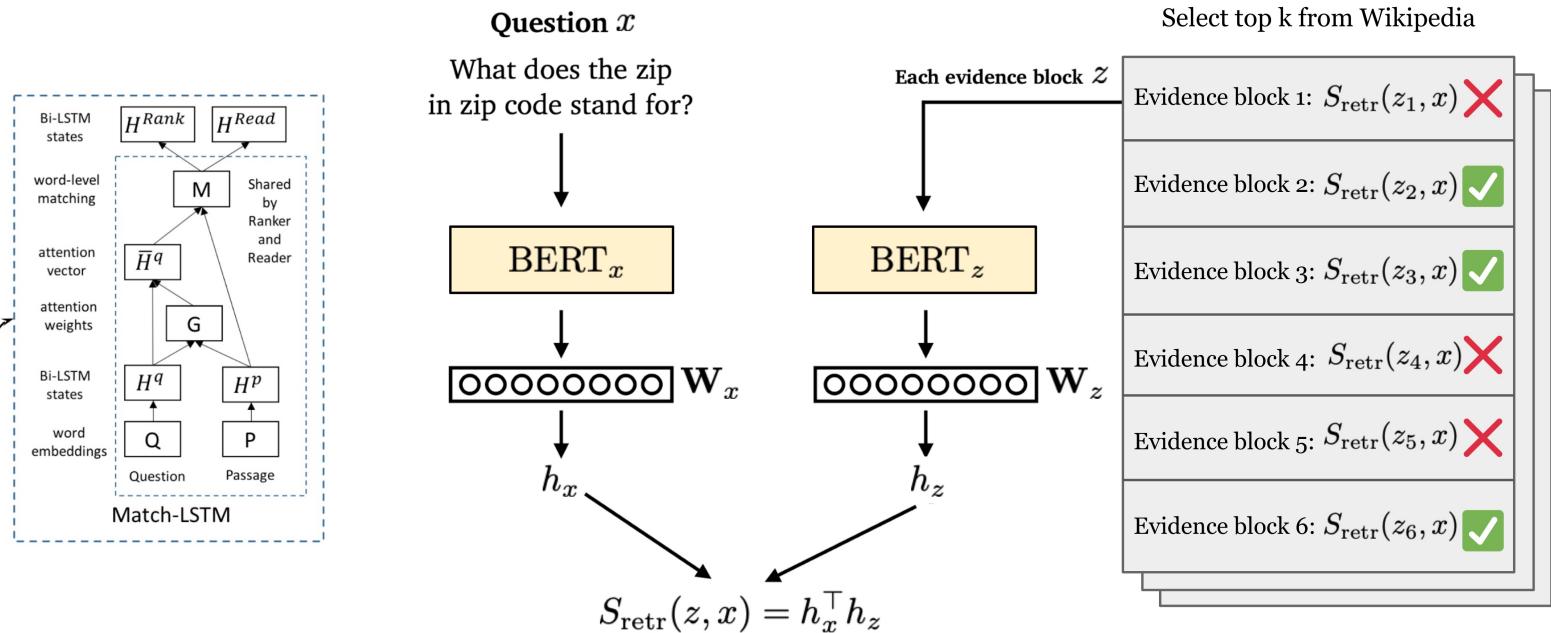
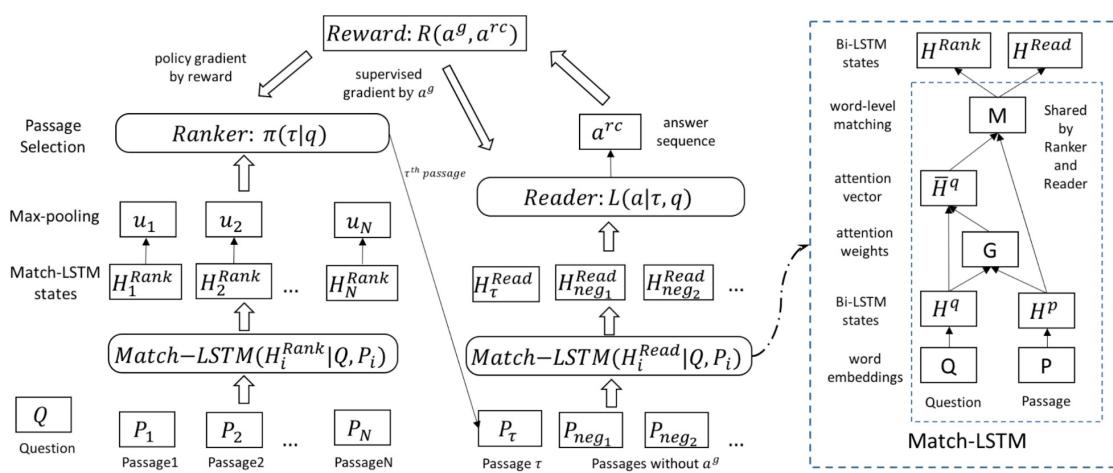
$$\text{tf-idf}(t, d, \mathcal{D}) = \text{tf}(t, d) \times \text{idf}(t, \mathcal{D})$$

$$\text{tf}(t, d) = \log(1 + \text{freq}(t, d))$$

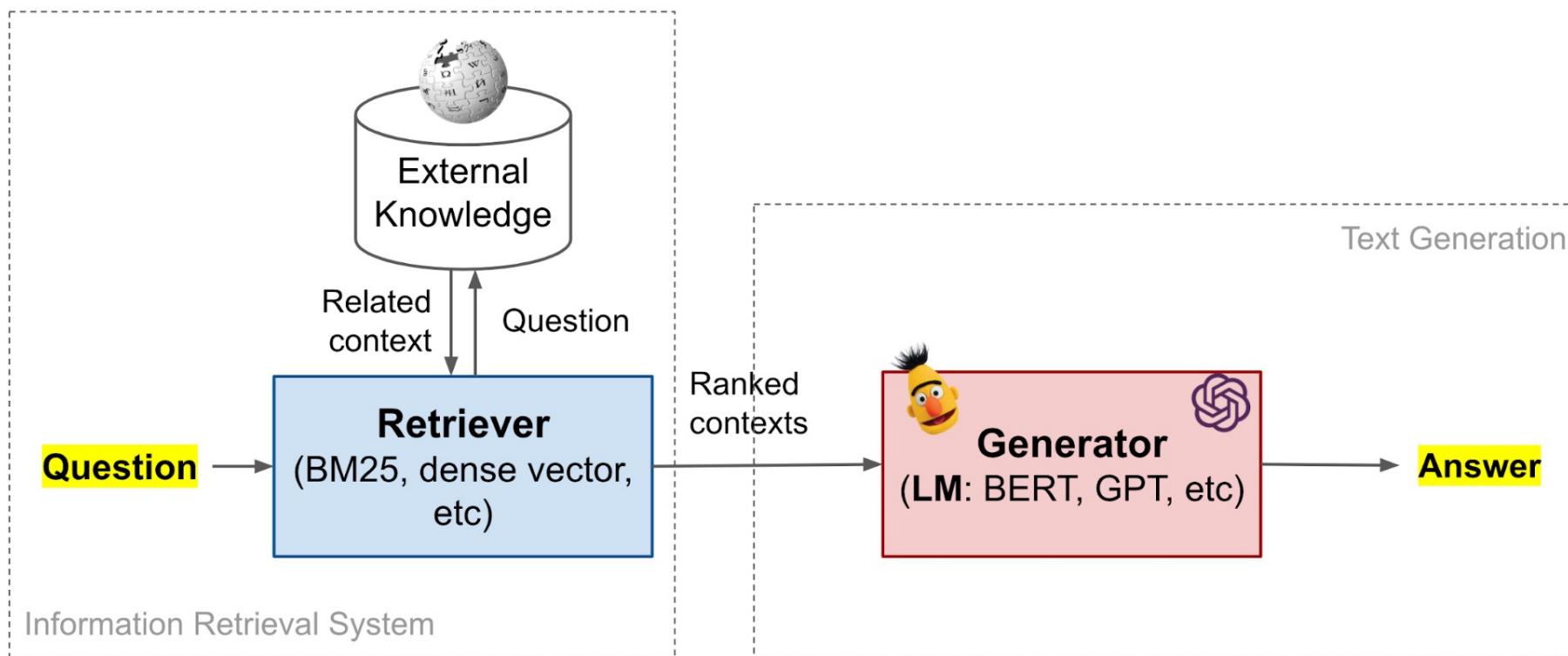
$$\text{idf}(t, \mathcal{D}) = \log\left(\frac{|\mathcal{D}|}{|d \in \mathcal{D} : t \in d|}\right)$$

# Open-book QA: Retriever-Reader

## End to end training

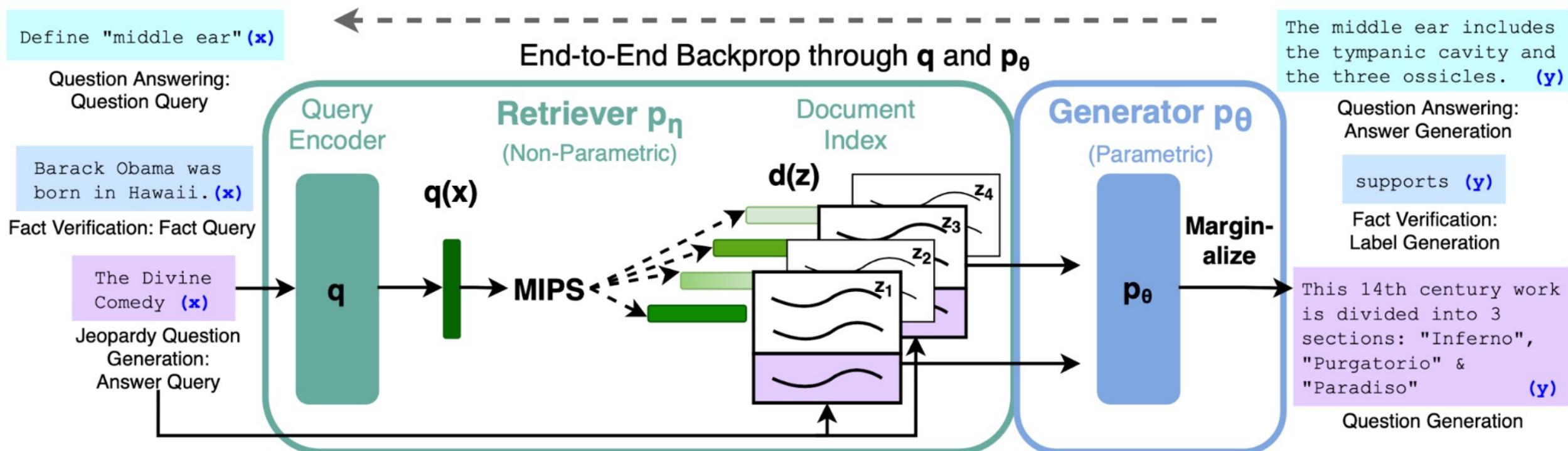


# Open-book QA: Retriever-Generator



# Open-book QA: Retriever-Generator

## End to end training



# Заключение

- Attention (повторение)
- CLIP framework
  - Transformer for images
- Question answering (QA)
  - Text generation
  - Close & Open domain