



Факультет компьютерных наук

Мультимодальные нейронные  
сети

Москва 2025

# Энкодеры



# Что такое энкодеры

Энкодер – модель или компонент модели, которая преобразует входные данные в векторное представление.



# Какие бывают энкодеры

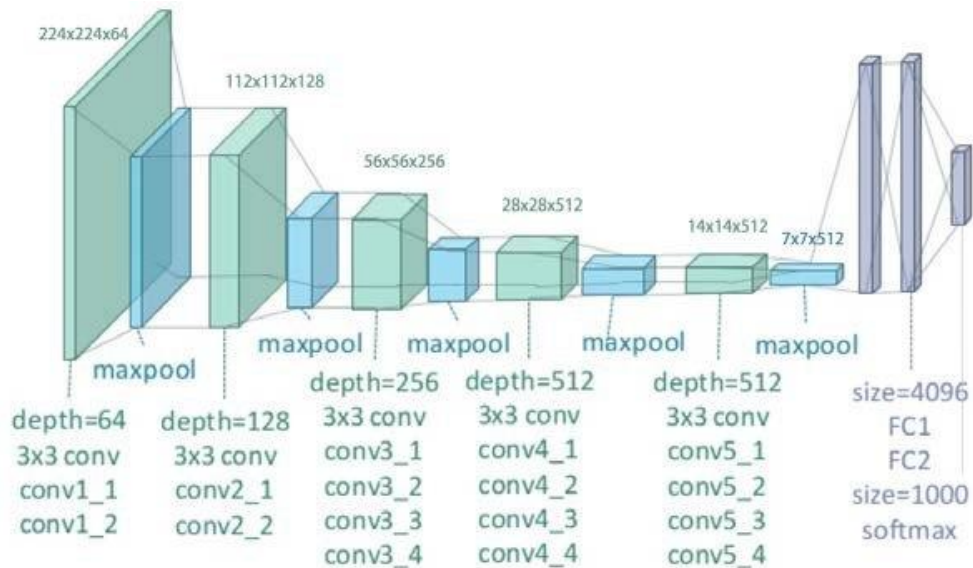
Скрытый слой обученной модели. Например, последний слой ResNet

Обучение специальной модели энкодера. Например, BERT

Обучение автокодировщика, например VAE

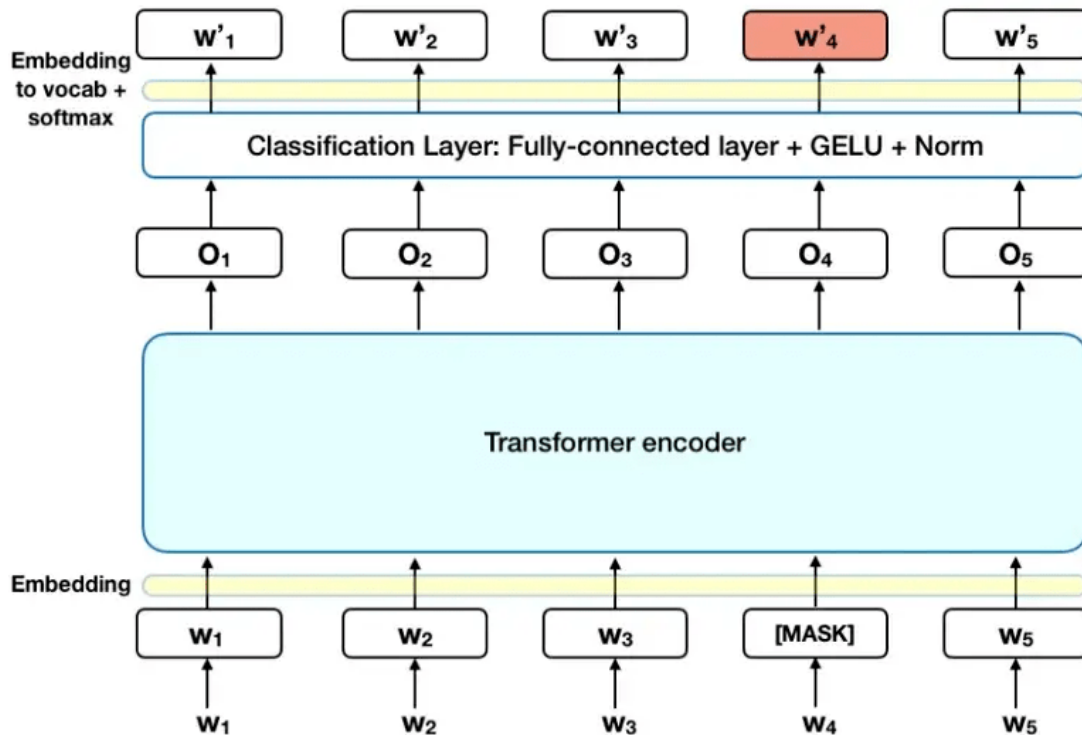
Обученный энкодер seq2seq модели. Например, Whisper или FridaT5

# скрытый слой

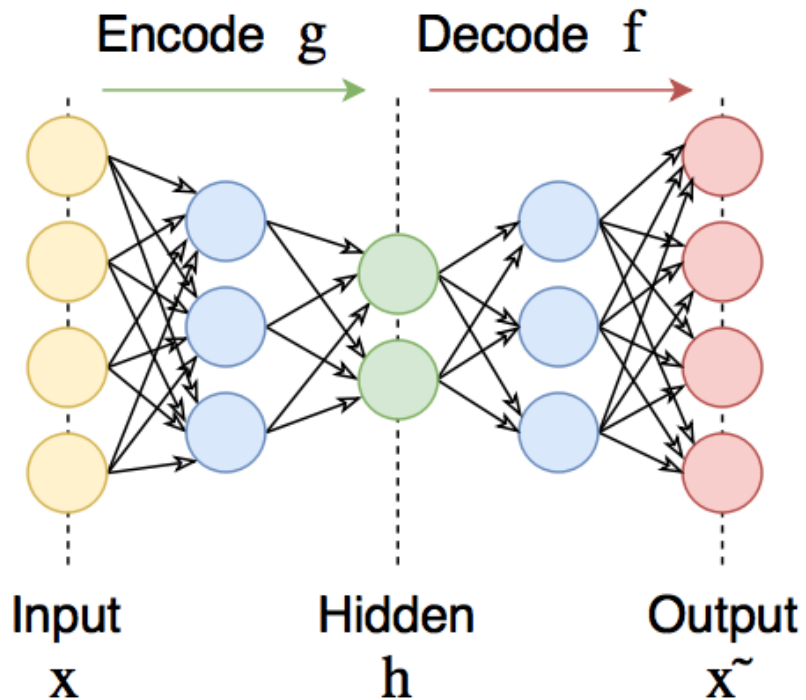




# BERT

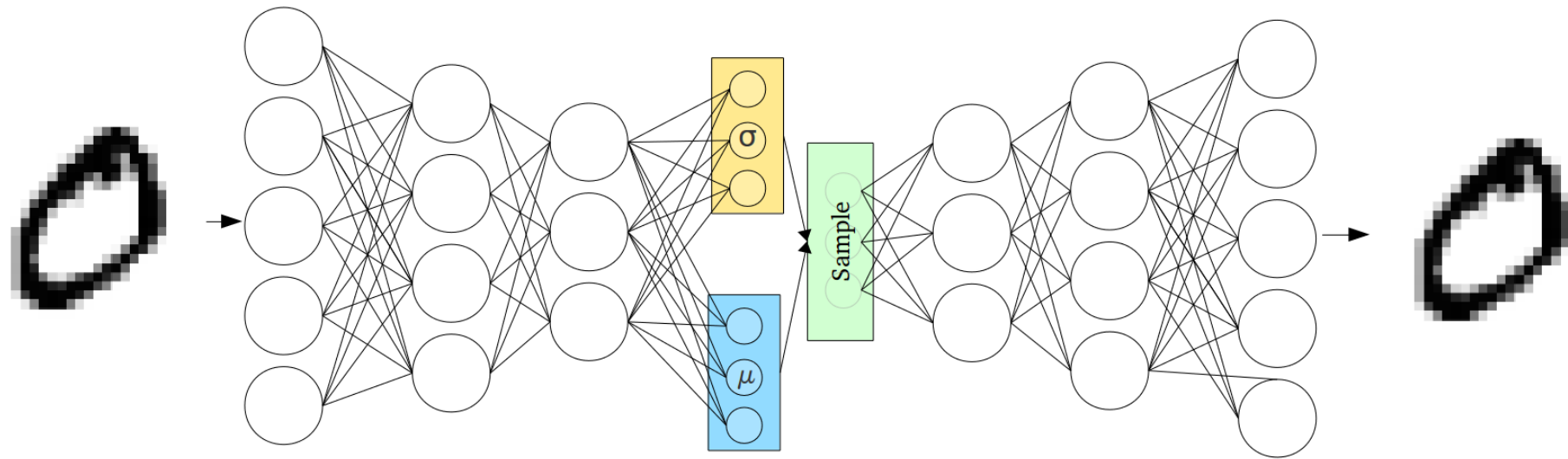


# Автокодировщик





# Вариационный автокодировщик



# Энкодер seq2seq модели

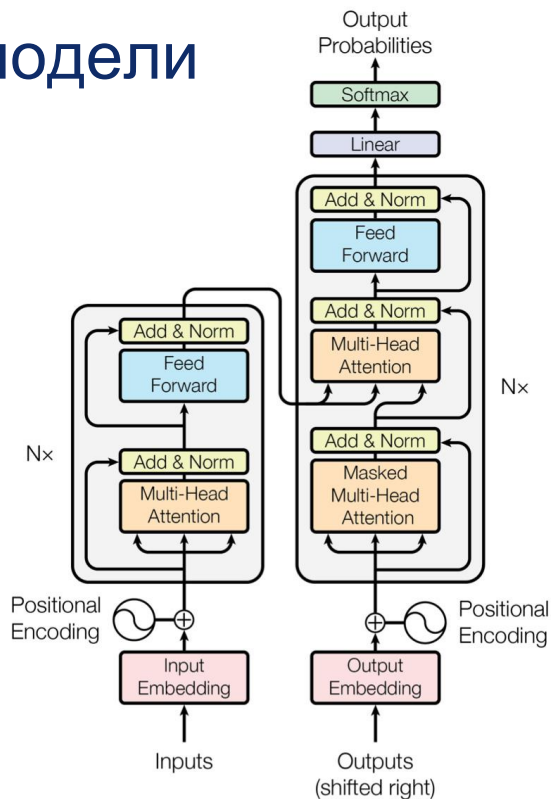


Figure 1: The Transformer - model architecture.





# Энкодеры изображений

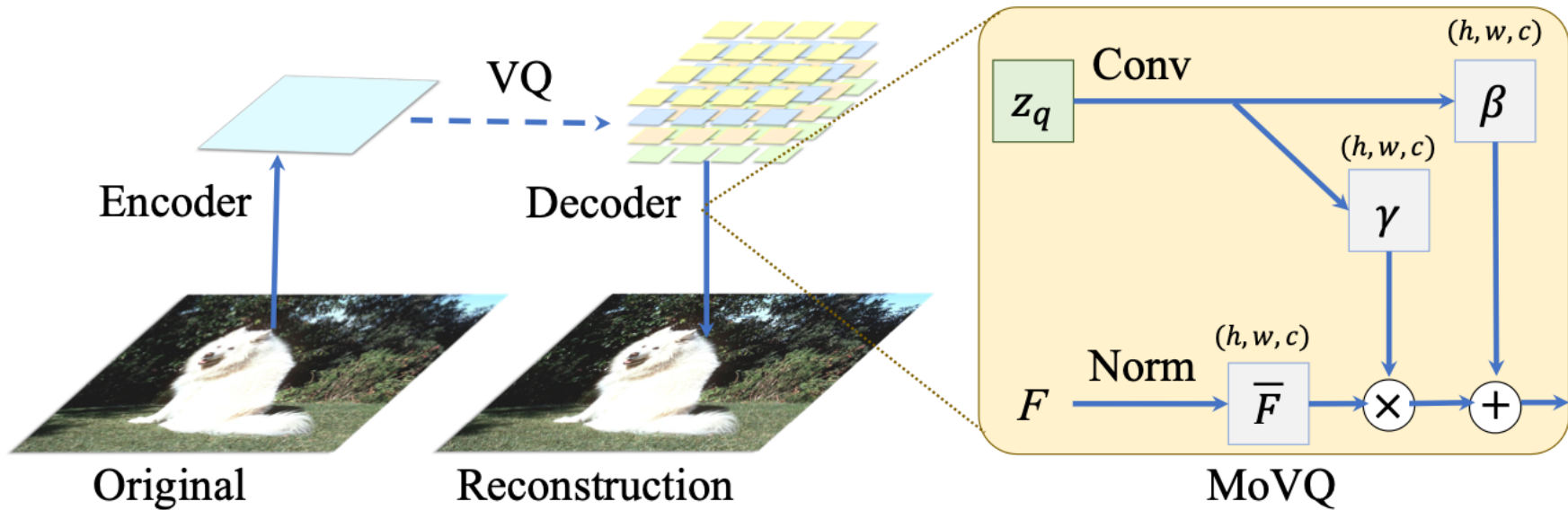
image to vector: ResNet, AE, VaE

image to latent: VQGAN

Image to sequence: ViT

# VQGAN

## Stage 1: Modulated VQ Encoder and Decoder Training





# VQGAN

отображение входного изображения в скрытое представление  
квантизация скрытого представления  
реконструкция изображения из многомерного массива токенов

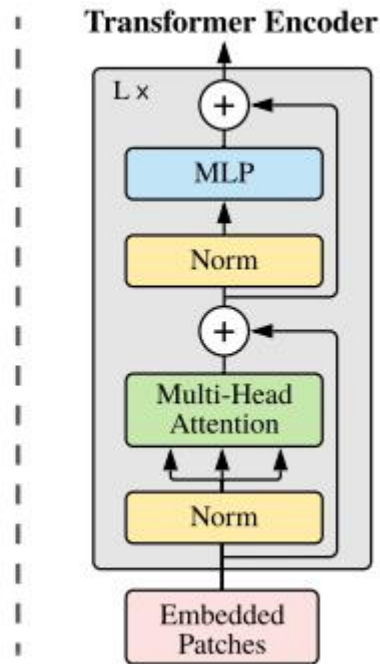
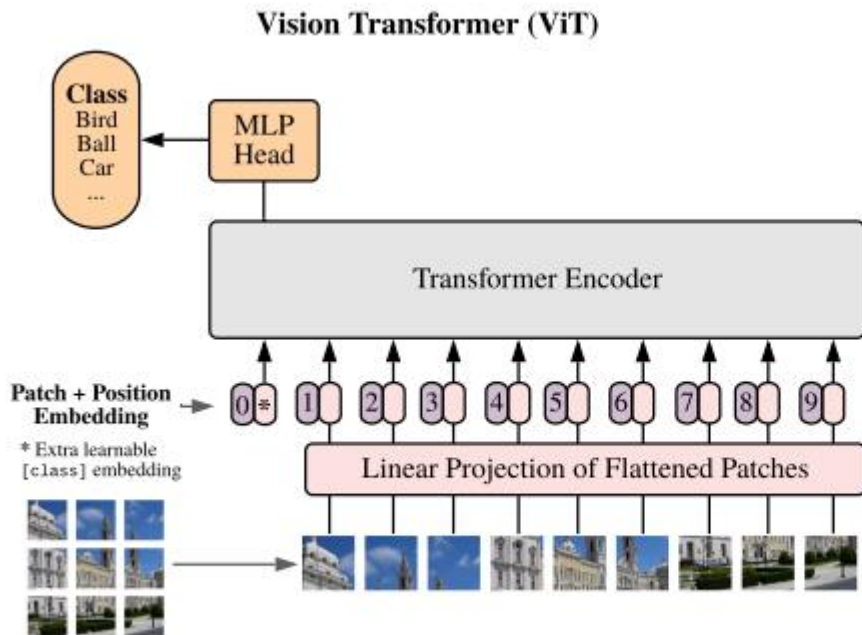


# VQGAN. задачи

генерация изображений

объединение модальностей для генерации изображений

# Vision transformer





# Vision transformer

	Ours-JFT (ViT-H/14)	Ours-JFT (ViT-L/16)	Ours-I21k (ViT-L/16)	BiT-L (ResNet152x4)	Noisy Student (EfficientNet-L2)
ImageNet	<b>88.55</b> $\pm 0.04$	87.76 $\pm 0.03$	85.30 $\pm 0.02$	87.54 $\pm 0.02$	88.4/88.5*
ImageNet ReaL	<b>90.72</b> $\pm 0.05$	90.54 $\pm 0.03$	88.62 $\pm 0.05$	90.54	90.55
CIFAR-10	<b>99.50</b> $\pm 0.06$	99.42 $\pm 0.03$	99.15 $\pm 0.03$	99.37 $\pm 0.06$	—
CIFAR-100	<b>94.55</b> $\pm 0.04$	93.90 $\pm 0.05$	93.25 $\pm 0.05$	93.51 $\pm 0.08$	—
Oxford-IIIT Pets	<b>97.56</b> $\pm 0.03$	97.32 $\pm 0.11$	94.67 $\pm 0.15$	96.62 $\pm 0.23$	—
Oxford Flowers-102	99.68 $\pm 0.02$	<b>99.74</b> $\pm 0.00$	99.61 $\pm 0.02$	99.63 $\pm 0.03$	—
VTAB (19 tasks)	<b>77.63</b> $\pm 0.23$	76.28 $\pm 0.46$	72.72 $\pm 0.21$	76.29 $\pm 1.70$	—
TPUv3-core-days	2.5k	0.68k	0.23k	9.9k	12.3k

# Vision transformer

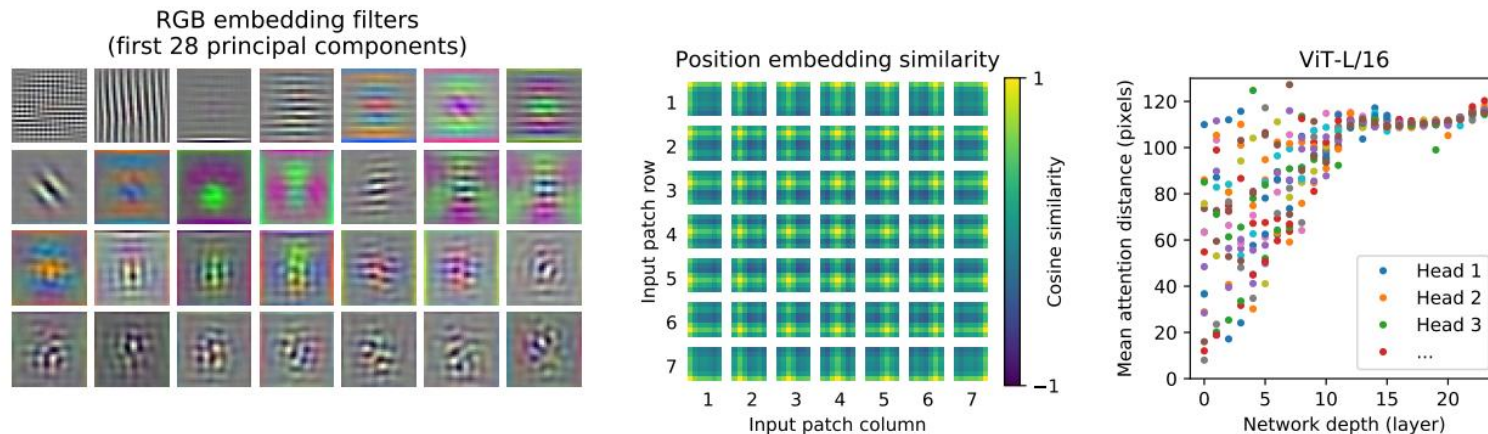


Figure 7: **Left:** Filters of the initial linear embedding of RGB values of ViT-L/32. **Center:** Similarity of position embeddings of ViT-L/32. Tiles show the cosine similarity between the position embedding of the patch with the indicated row and column and the position embeddings of all other patches. **Right:** Size of attended area by head and network depth. Each dot shows the mean attention distance across images for one of 16 heads at one layer. See Appendix D.7 for details.



# Vision transformer. Задачи

кодирование изображения в последовательность для LLM

классические задачи с изображениями – классификация, детекция изображений





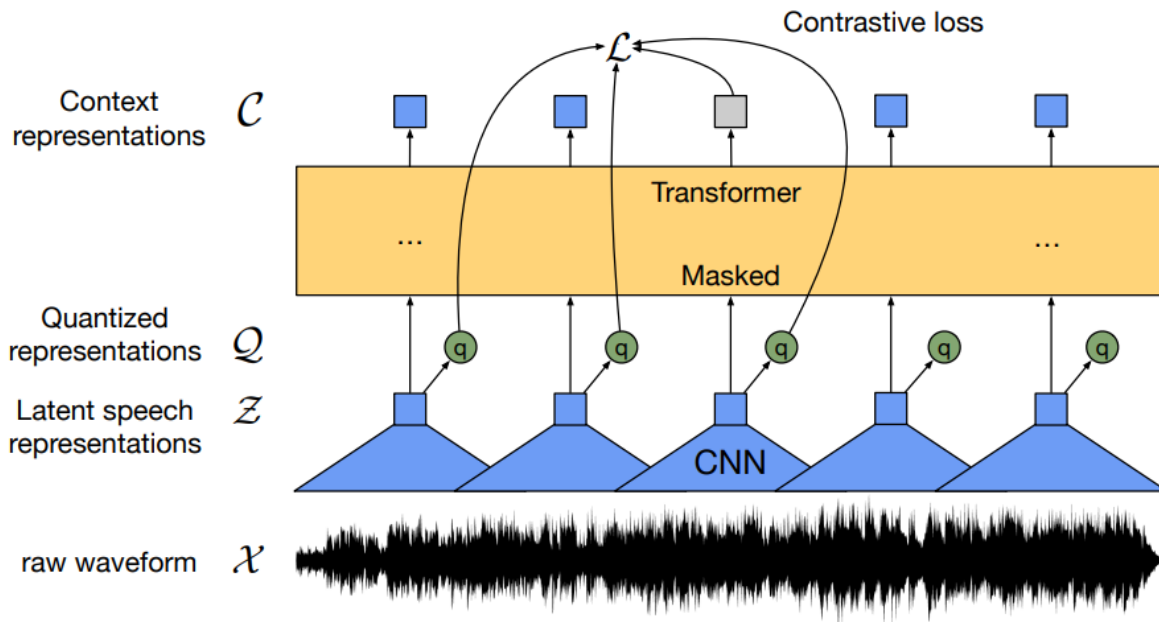
# Аудио энкодеры

энкодеры – спектрограмма, wav2vec, HuBERT

токенизаторы – RVQ, Encodec, Mimi, Xcodec

большие модели - whisper

# Wav2Vec2





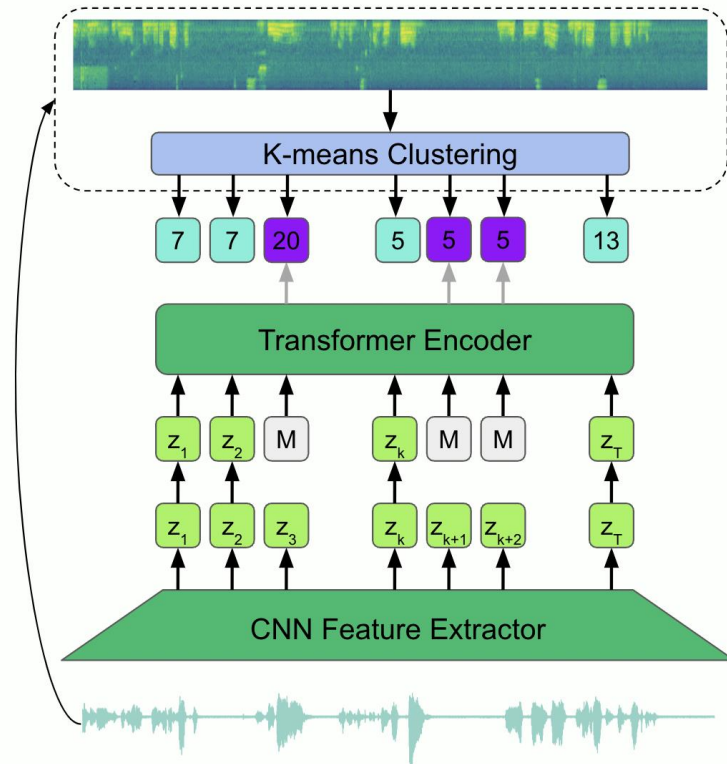
# Wav2Vec2

кодирует общую информацию об аудио

подходит для универсальных задач (captioning, classification, etc)



# HuBERT





# HuBERT

кодирует псевдофонетическую информацию  
подходит для речевых задач (ASR, translation, etc)



# Whisper

## Multitask training data (680k hours)

### English transcription

- 🗣️ "Ask not what your country can do for ..."
- 📄 Ask not what your country can do for ...

### Any-to-English speech translation

- 🗣️ "El rápido zorro marrón salta sobre ..."
- 📄 The quick brown fox jumps over ...

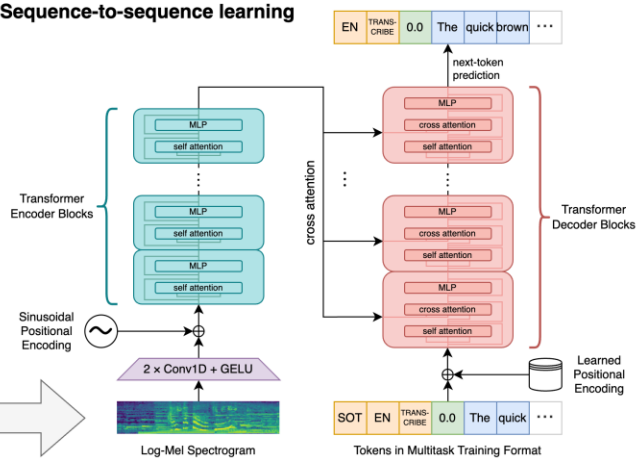
### Non-English transcription

- 🗣️ "언덕 위에 올라 내려다보면 너무나 넓고 넓은 ..."
- 📄 언덕 위에 올라 내려다보면 너무나 넓고 넓은 ...

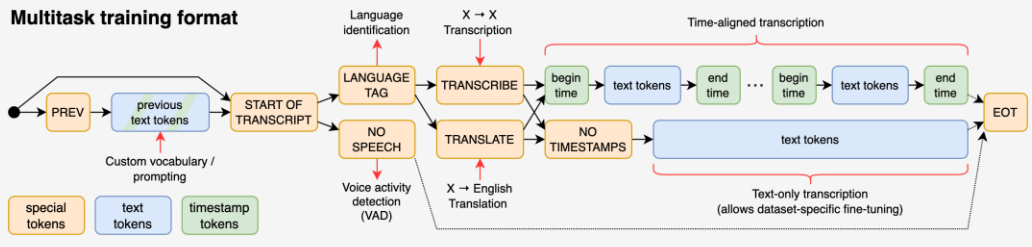
### No speech

- 🎧 (background music playing)
- 📄 ∅

## Sequence-to-sequence learning



## Multitask training format



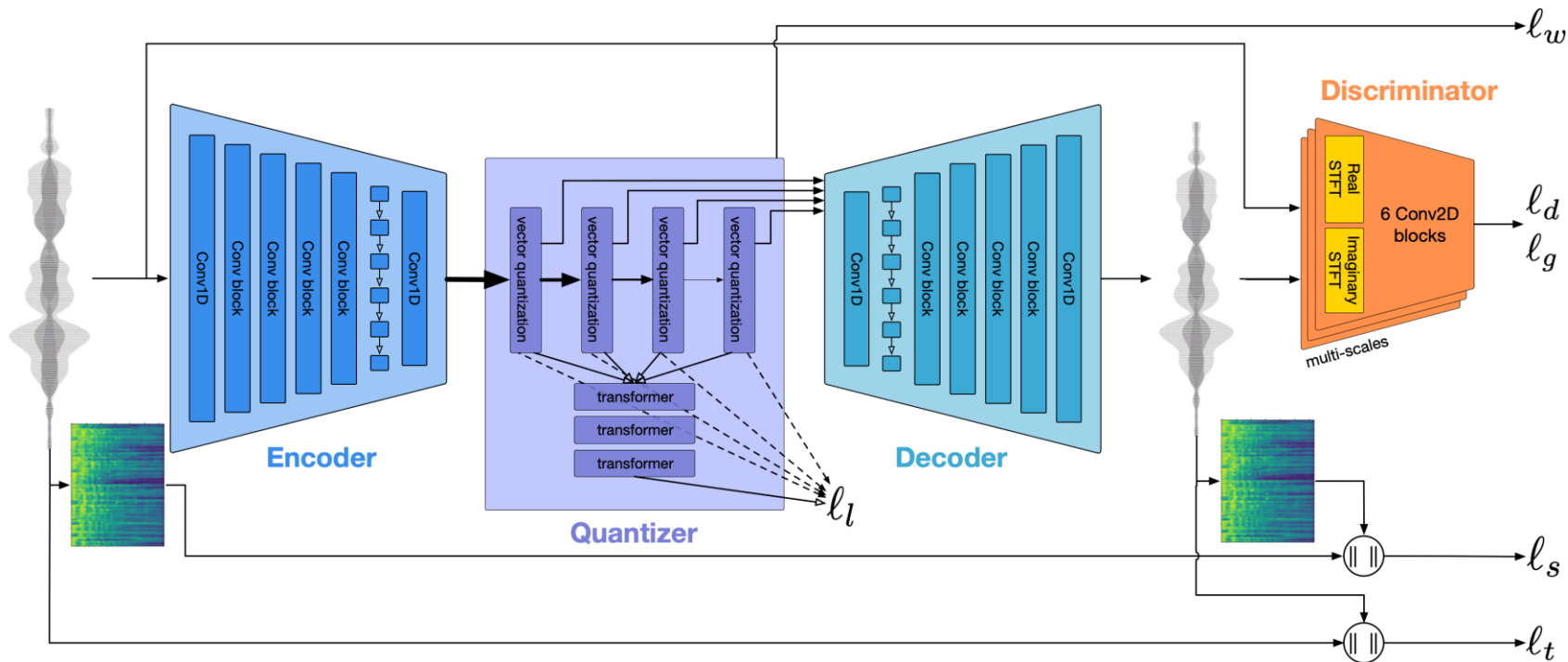


# Whisper

понимание языка

мультязычные задачи

# Encodec







# Encodec

сжатие аудио  
генерация аудио



# Видео

кодирование изображения и звука отдельно и объединение  
(SALMONN, VideoLLaVa etc)

Специальные энкодеры (Cosmoss)



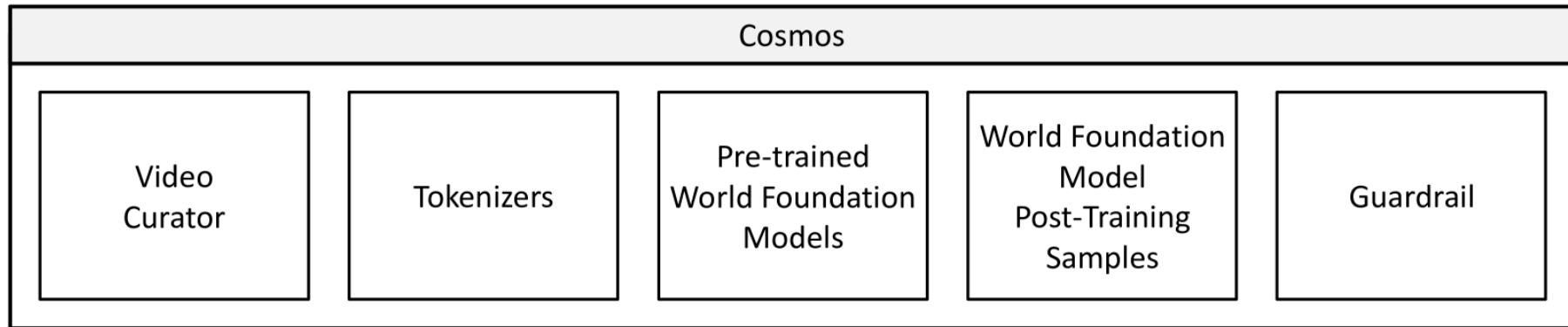
# Видео

кодирование изображения и звука отдельно и объединение  
(SALMONN, VideoLLaVa etc)

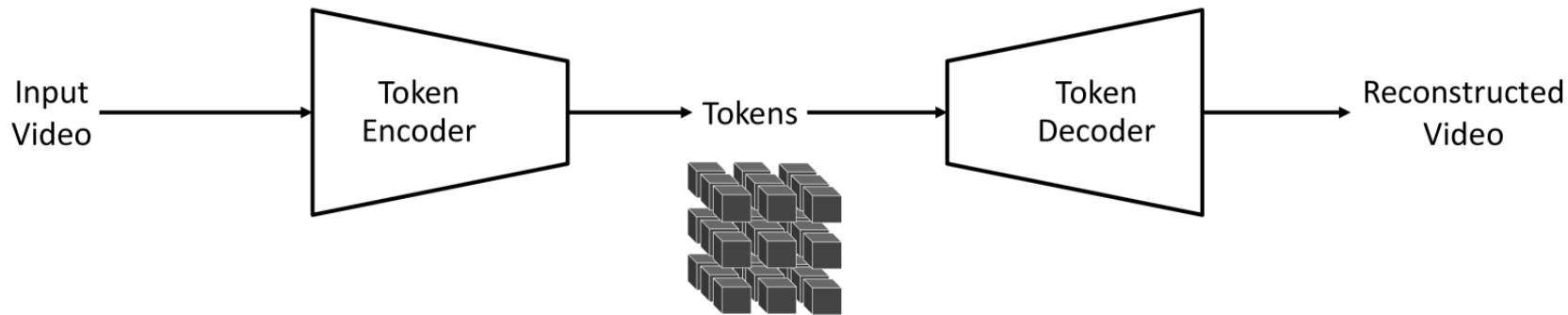
Специальные энкодеры (Cosmoss)



## Видео. Cosmos



# Видео. Cosmos





**ВСЁ!**