

# **TTIC 31230, Fundamentals of Deep Learning**

David McAllester, Autumn 2022

## **Vector Quantization for Autoregressive Modeling**

# Autoregressive Image and Voice Modeling

Strong VAE image modeling was first achieved with autoregressive token modeling.

van den Oord, Vinyals and Kavukcuoglu,  
Neural Discrete Representation Learning, **2017**



# VQ-VAE-2, June 2019

Generating Diverse High-Fidelity Images with VQ-VAE-2



Figure 1: Class-conditional 256x256 image samples from a two-level model trained on ImageNet.

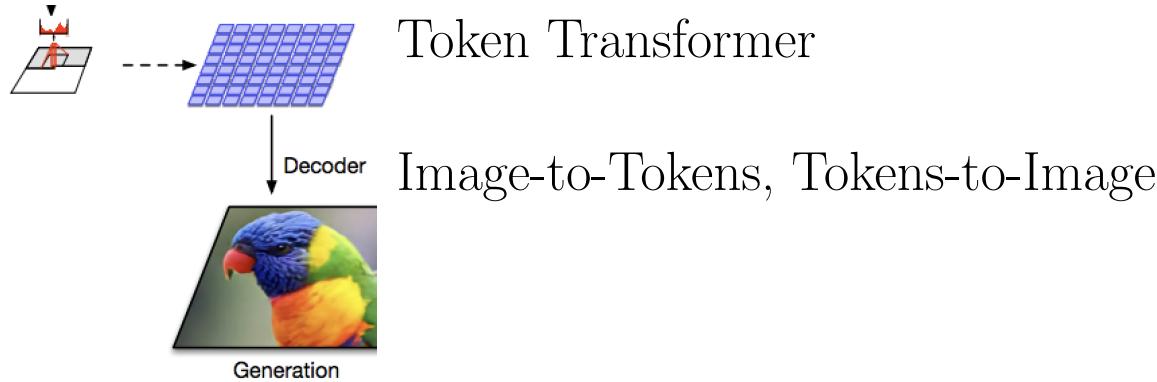
# VAE Tokenization of Images and Voice



Let  $y$  range over a population (such as images or sound waves).

Assume that a given  $y$  is encoded as a tensor denoted  $z_{\text{enc},p}(y)$  where  $p$  is a “position in  $y$ ” (a pixel in an image tensor or time window in a sound tensor) and  $z_{\text{enc},p}(y) \in R^d$  is a vector.

# Vector Quantization (Tokenization)

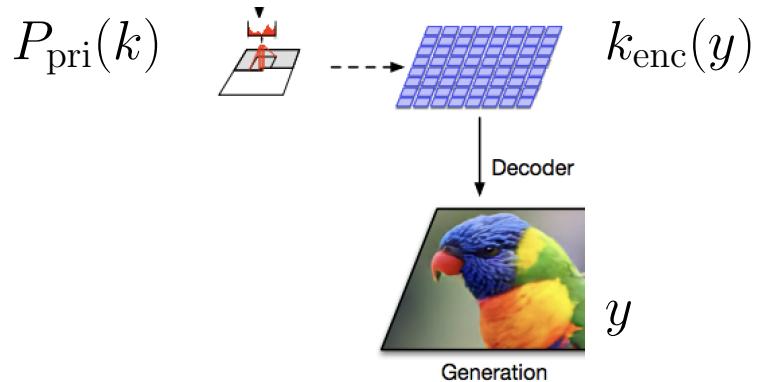


Assume a finite set of  $K$  “tokens” where token  $k$  has an embedding vector  $e(k) \in R^d$ .

Define  $k_{\text{enc},p}(y)$  by

$$k_{\text{enc},p}(y) = \operatorname{argmin}_k \frac{1}{2} \|z_{\text{enc},p}(y) - e(k)\|^2$$

## Reconstruction Loss

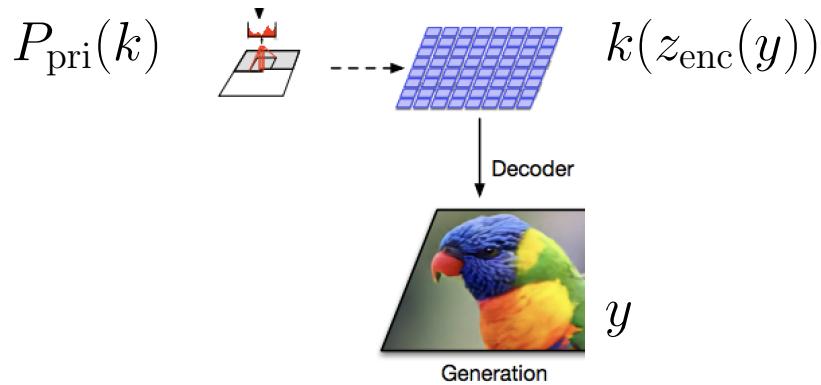


We now have a VAE where the tensor  $k_{\text{enc}}(y)$  is the latent variable.

The encoder and decoder are trained jointly.

The prior is a transformer trained after the encoder and decoder are fully trained.

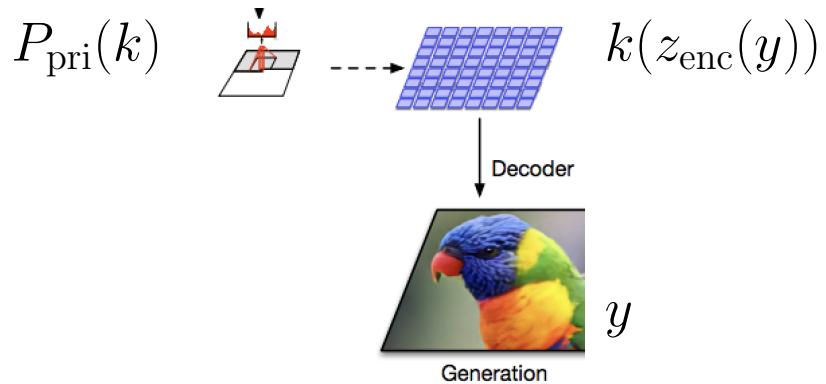
# Training the Encoder and Decoder



Taking  $P_{\text{dec}}(y|k)$  to be  $\mathcal{N}(\hat{y}_{\text{dec}}(k), I)$  we get an  $L_2$  reconstruction loss.

$$\mathcal{L}_{\text{Rec}} = \frac{1}{2} \| y - \hat{y}_{\text{dec}}(e(k(z_{\text{enc}}(y)))) \| ^2$$

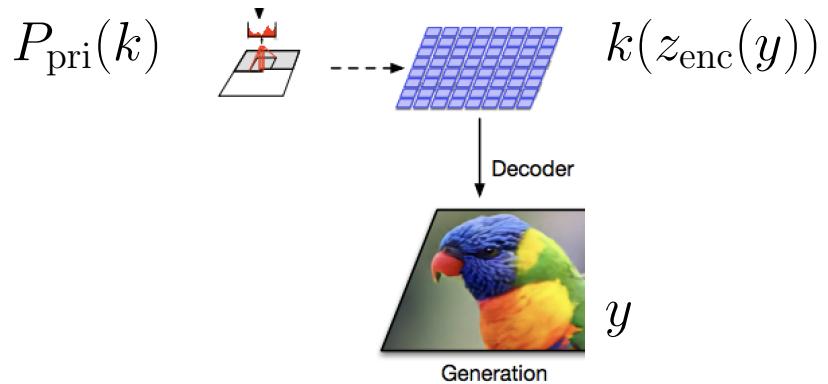
# Training the Encoder and Decoder



$$\mathcal{L}_{\text{Rec}} = \frac{1}{2} \| y - \hat{y}_{\text{dec}}(e(k(z_{\text{enc}}(y)))) \| ^2$$

Because the tokens are discrete we do not get any gradient on  $z_{\text{enc}}(y)$ .

# Straight-Through Gradients



$$\mathcal{L}_{\text{Rec}} = \frac{1}{2} \| y - \hat{y}_{\text{dec}}(e(k(z_{\text{enc}}(y)))) \|^2$$

$$z_{\text{enc},p}(y).\text{grad} = \nabla_{e(k(z_{\text{enc},p}(y)))} \mathcal{L}_{\text{Rec}}$$

## K-Means Gradients

We train  $z_{\text{enc}}(y)$  and the token embeddings  $e(k)$ .

$$\mathcal{L}_{\text{Rec}} = \frac{1}{2} \| y - \hat{y}_{\text{dec}}(e(k(z_{\text{enc}}(y)))) \|^2$$

$$z_{\text{enc},p}(y).\text{grad} = \nabla_{e(k(z_{\text{enc},p}(y)))} \mathcal{L}_{\text{Rec}}$$

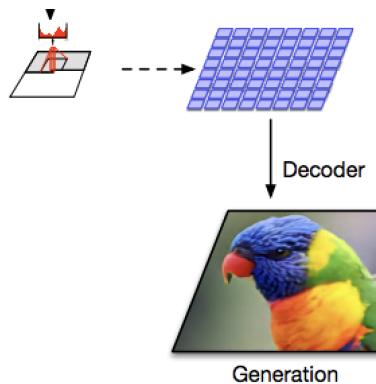
$$\mathcal{L}_{\text{KM}} = \frac{1}{2} \| z_{\text{enc},p}(y) - e(k(z_{\text{enc},p})) \|^2$$

$$e(k(z_{\text{enc},p}(y))).\text{grad} += \beta \nabla_{e(k(z_{\text{enc},p}(y)))} \mathcal{L}_{\text{KM}}$$

$\beta$  is a hyper-parameter that adjust the relative learning rates.

# Transformer Training

Finally we hold the encoder fixed and train the prior  $P_{\text{pri}}(z)$  to be an auto-regressive model of the symbolic image  $k_{\text{enc}}(s)[X, Y]$ .



## Tokenization and Gaussian Mixture Models (GMMs)

Consider modeling  $P(y|x)$  with  $y \in R^d$

A Gaussian model has the form

$$y = \hat{y}(x) + \sigma\epsilon, \quad \epsilon \sim \mathcal{N}(0, I)$$

$$\hat{y} = \underset{\hat{y}}{\operatorname{argmin}} E_{x,y} ||\hat{y}(x) - y||^2$$

## Tokenization and Gaussian Mixture Models (GMMs)

Now consider a tokenizing decoder

$$y = E_{k \sim P_{\text{dec}}(k|x)} [ e(k) + \sigma \epsilon ], \quad \epsilon \sim \mathcal{N}(0, I)$$

We get that  $P_{\text{dec}}(y|x)$  is a Gaussian mixture model (GMM).

GMMs are significantly more expressive than single Gaussians.

## **Wav2Vec 2.0, June 2020, Facebook**

Trained on 53k hours of unlabeled audio (no text) they convert speech to a sequence of discrete quantized vectors they call “pseudo-text units”.

By training on only one hour of human-transcribed audio, and using the Wav2Vec transcription into pseudo-text, they outperform the previous state of the art in word error rate for 100 hours of human-transcribed text.

# DALLE-1, January 2021

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

AI-GENERATED  
IMAGES



[Edit prompt or view more images↓](#)

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

AI-GENERATED  
IMAGES



[Edit prompt or view more images↓](#)

## GLSM, February 2021, Facebook

Generative Spoken Language Model (GSLM)

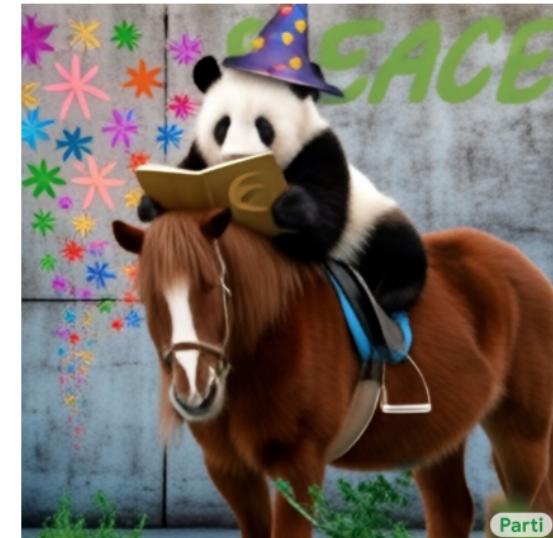
They then train a generative model of the sequences of pseudo-text units learned from unlabeled audio.

This model can continue speech from a speech prompt in much the same way that GPT-3 continues text from a text prompt.

Semantic and grammatical structure in a “unit language model” is recovered from speech alone.

# Parti, June 2022

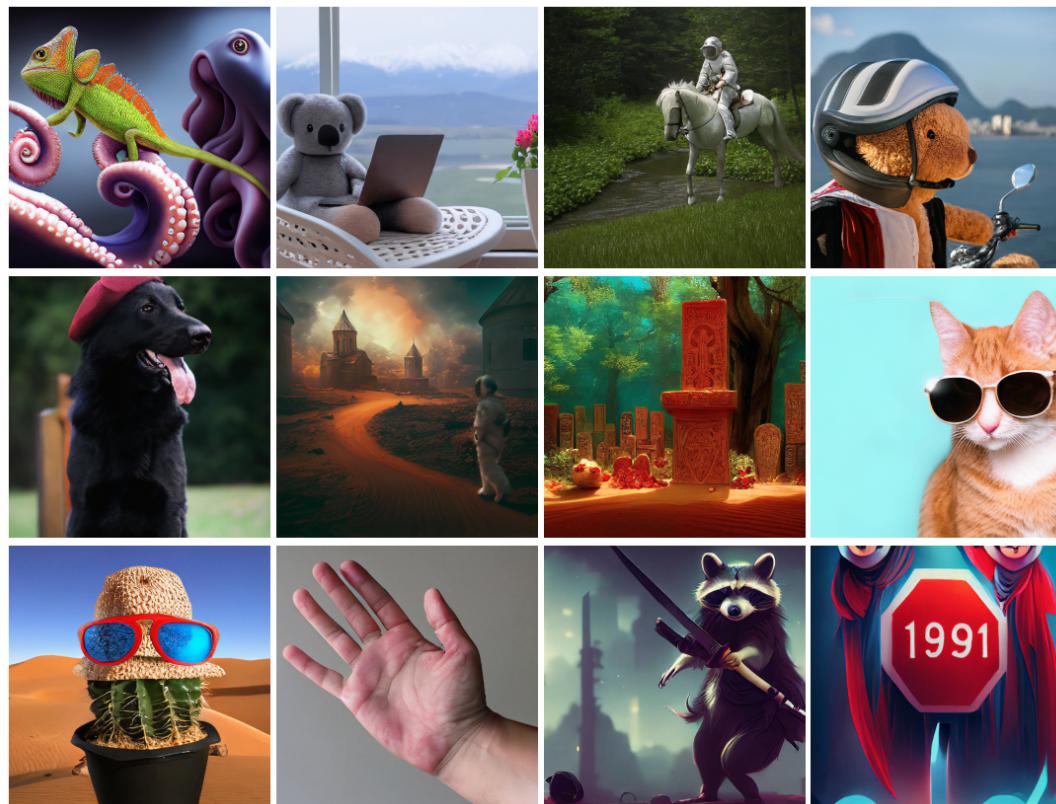
Scaling Autoregressive Models for Content-Rich Text-to-Image Generation  
Yu et al.



# CM3Leon, September 2023

Scaling Autoregressive Multi-Modal Models: Pretraining and Instruction Tuning

Yu et al.



## **Voice-Text Language Model (VoxtLM), September 2023**

This is similar to CM3Leon but for voice and text rather than images and text.

Voice is tokenized and then a transformer is used to model sequences that alternate voice and text.

**END**