# Neural Discrete Representation Learning

**Aaron van den Oord**, Oriol Vinyals, Koray Kavukcuoglu

Google DeepMind

# Generative Models

**Goal**: Estimate the <u>probability distribution</u> of high-dimensional data

Such as images, audio, video, text, ...
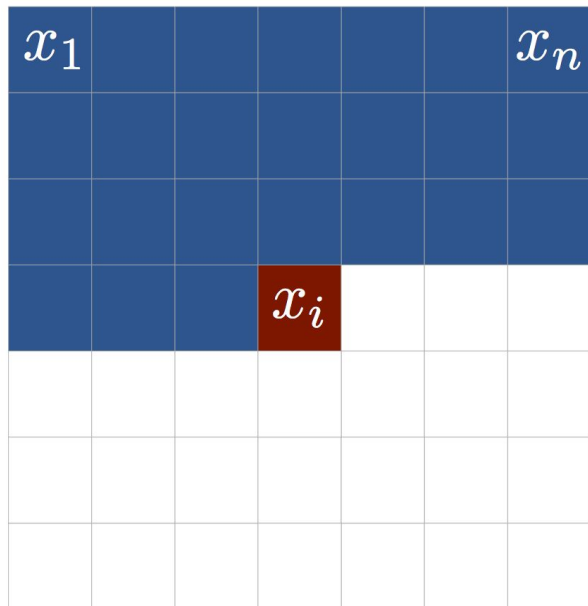
**Motivation:**

Learn the underlying structure in data.

Capture the dependencies between the variables.
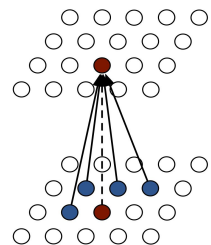
Generate new data with similar properties.

Learn useful features from the data in an unsupervised fashion.

# Autoregressive Models



$$p(\mathbf{x}) = \prod_{i=1}^{n^2} p(x_i | x_1, ..., x_{i-1})$$

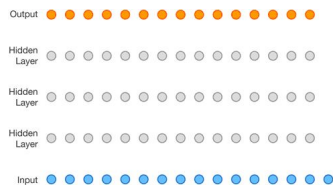# Recent Autoregressive models at DeepMind



PixelRNN

PixelCNN

van den Oord et al, 2016ab

Geyser

White Whale

Hartebeest

Tiger



Video Pixel Networks

Kalchbrenner et al, 2016a



WaveNet

1 Second

van den Oord et al, 2016c



ByteNet

Kalchbrenner et al, 2016b

Google DeepMind

# Modeling Audio



1 Second

# Causal Convolution



Hidden Layer

Input

Google DeepMind
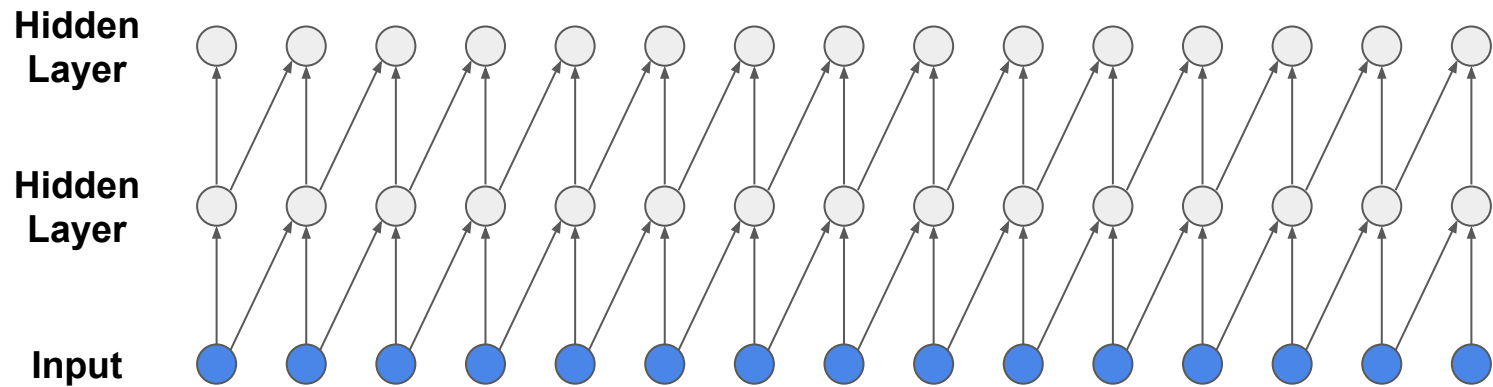
# Causal Convolution



Hidden Layer

Hidden Layer

Input

Google DeepMind

# Causal Convolution

# Causal Convolution

# Causal Convolution



**Output**

**Hidden Layer**

**Hidden Layer**

**Hidden Layer**

**Input**

Google DeepMind

# Causal Dilated Convolution

**Input**  ● ● ● ● ● ● ● ● ● ● ● ● ● ● ●

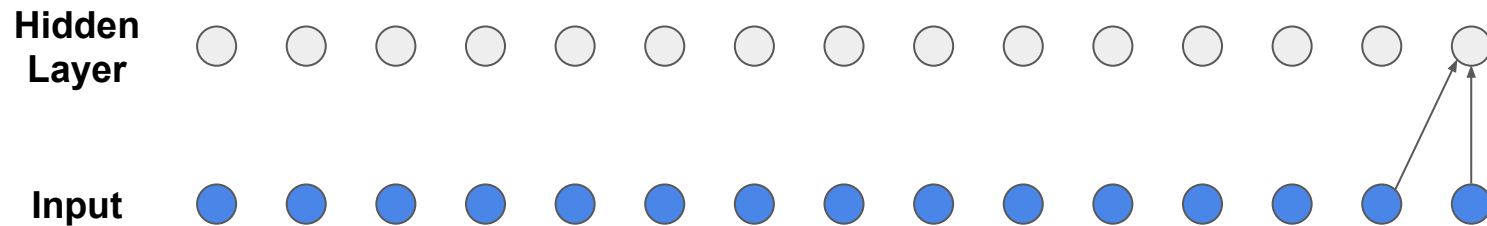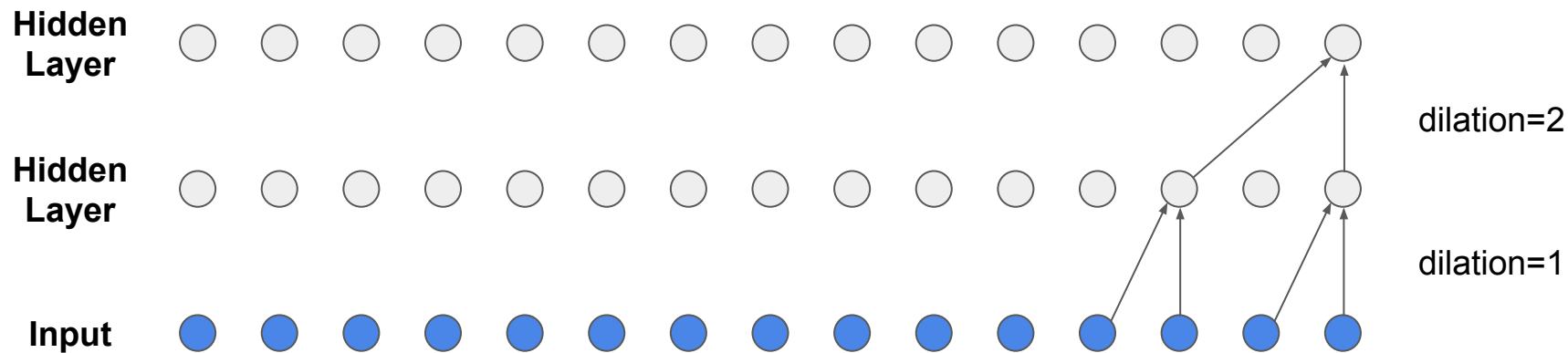Google DeepMind

# Causal Dilated Convolution
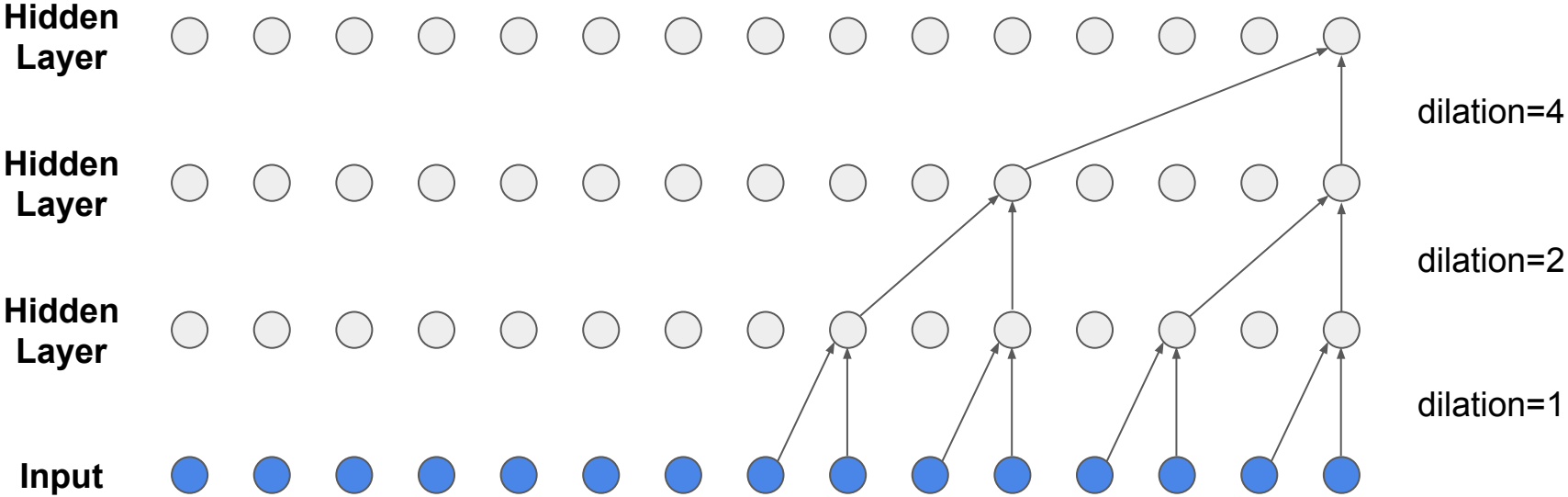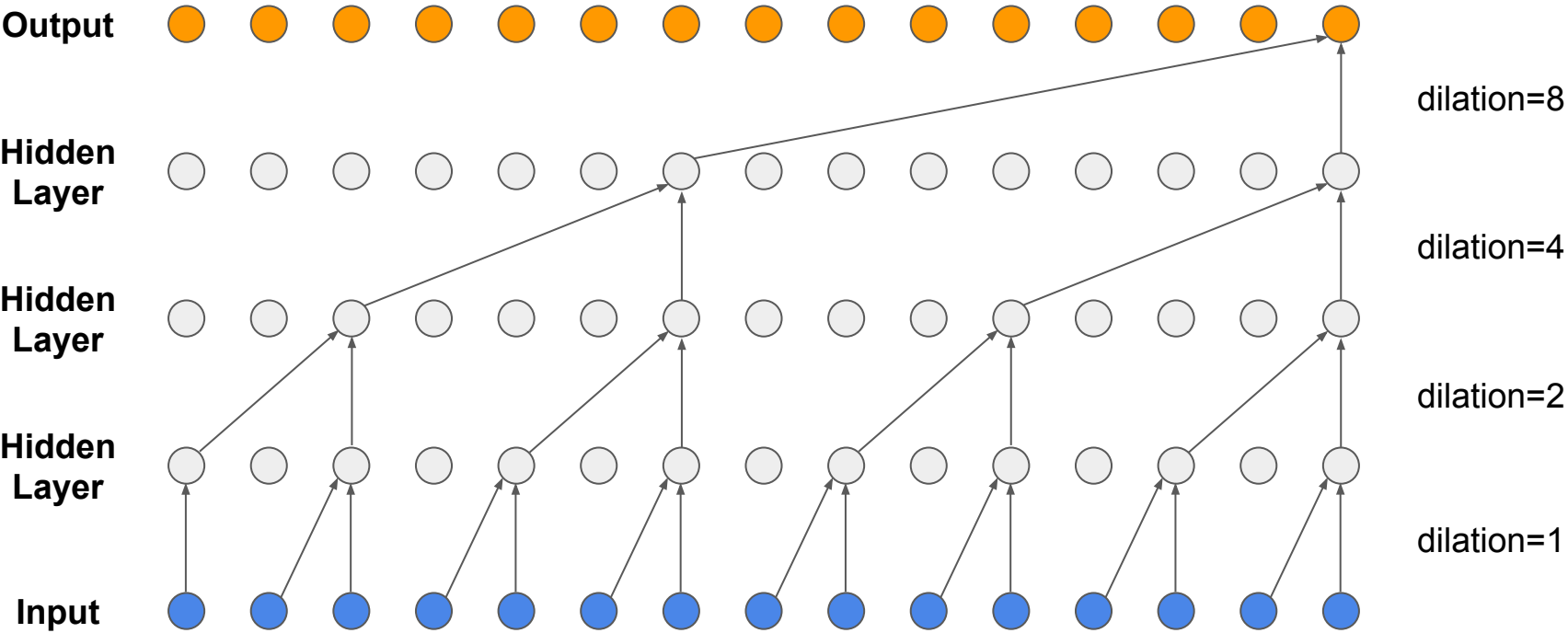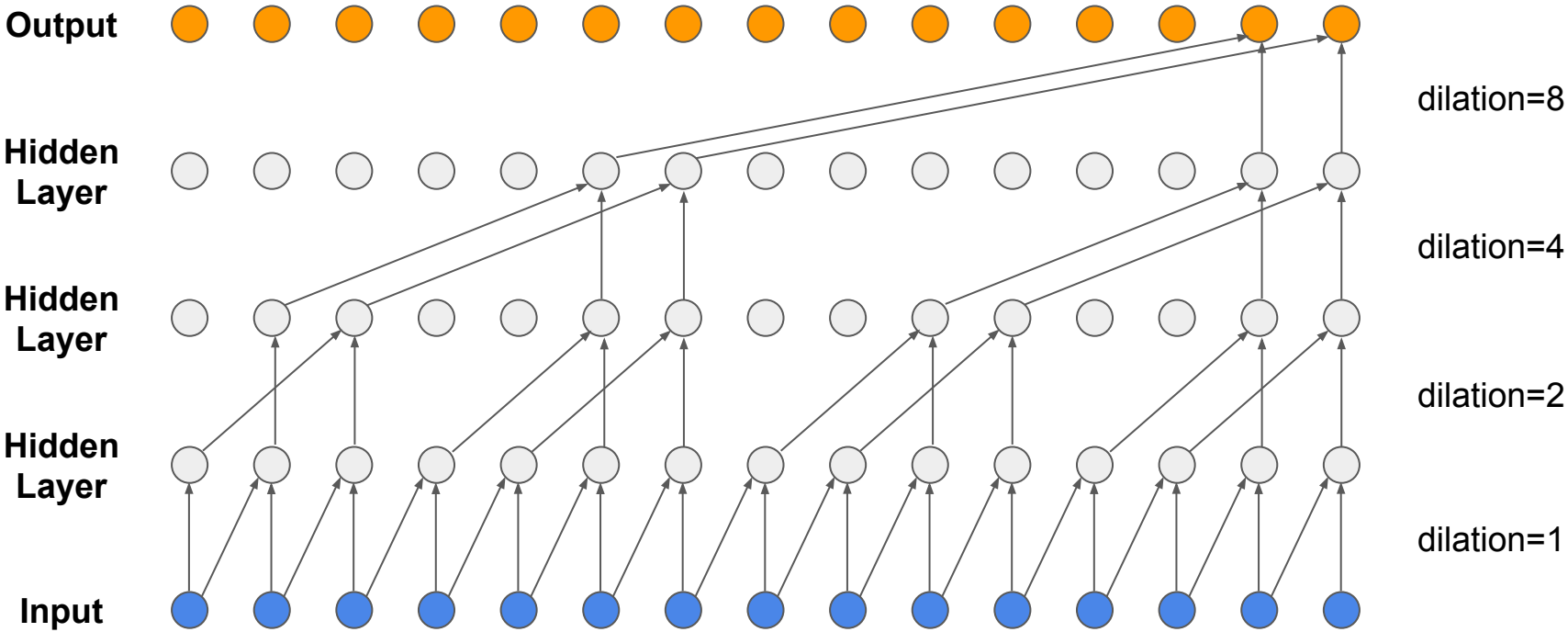
**Hidden Layer**

**Input**

# Causal Dilated Convolution

# Causal Dilated Convolution

# Causal Dilated Convolution



**Output**

**Hidden Layer**

dilation=8

**Hidden Layer**

dilation=4

**Hidden Layer**

dilation=2

**Input**

dilation=1

Google DeepMind

# Causal Dilated Convolution



Output

dilation=8

Hidden Layer

dilation=4

Hidden Layer

dilation=2

Hidden Layer

dilation=1

Input

Google DeepMind

# Multiple Stacks

# Sampling

# Speaker-conditional Generation



Speaker embedding
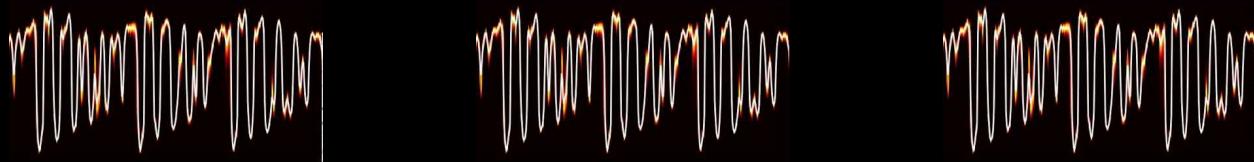
Does **not** depend on timestep
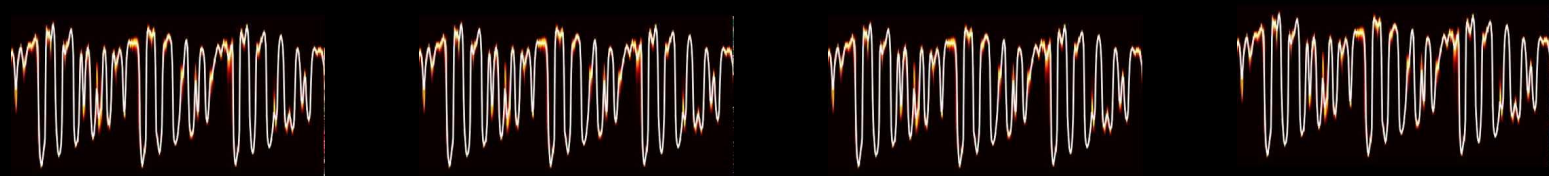
Text-To-Speech samples

# Speaker-conditional samples

(but not conditioned on text)

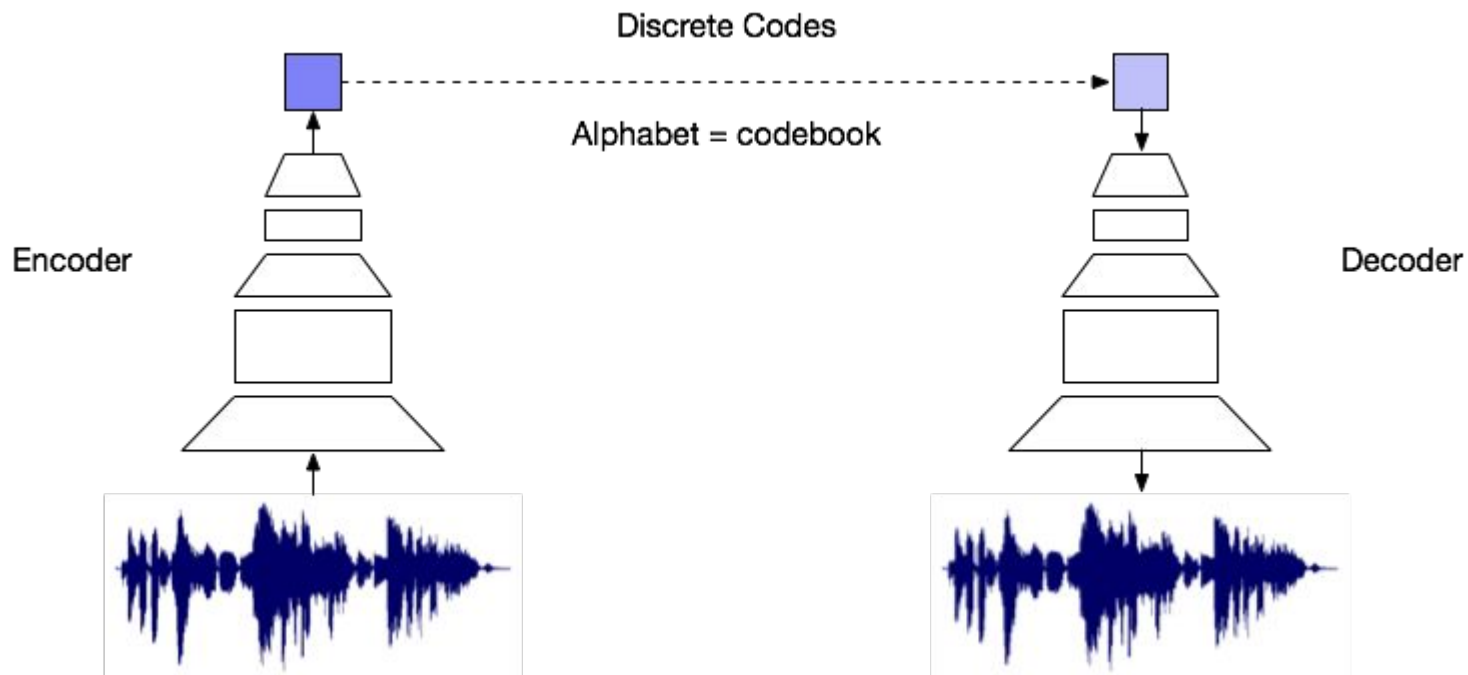Google DeepMind

# Piano Music samples

# VQ-VAE

- Towards modeling a latent space
    - Learn meaningful representations.
    - Abstract away noise and details.
    - Model what's important in a compressed latent representation.

- Why discrete?
    - Many important real-world things are discrete.
    - Arguably easier to model for the prior (e.g., softmax vs RNADE)
    - Continuous representations are often inherently discretized by encoder/decoder.
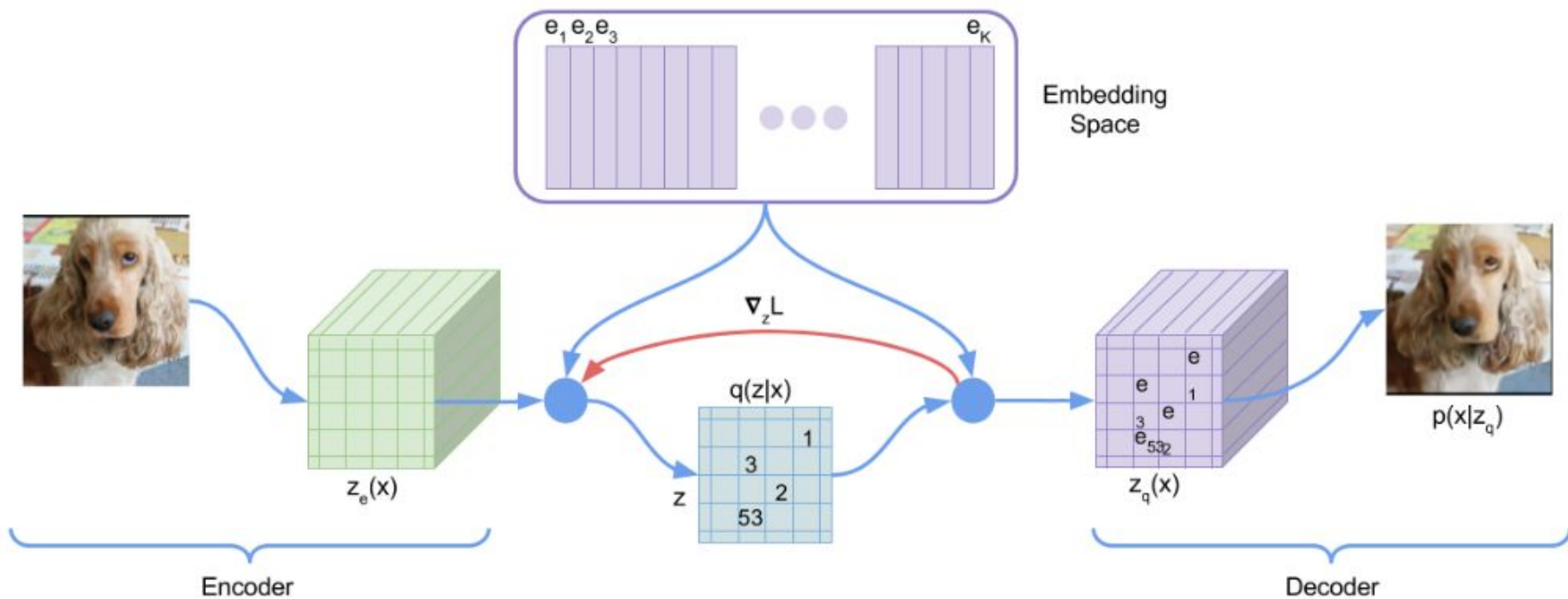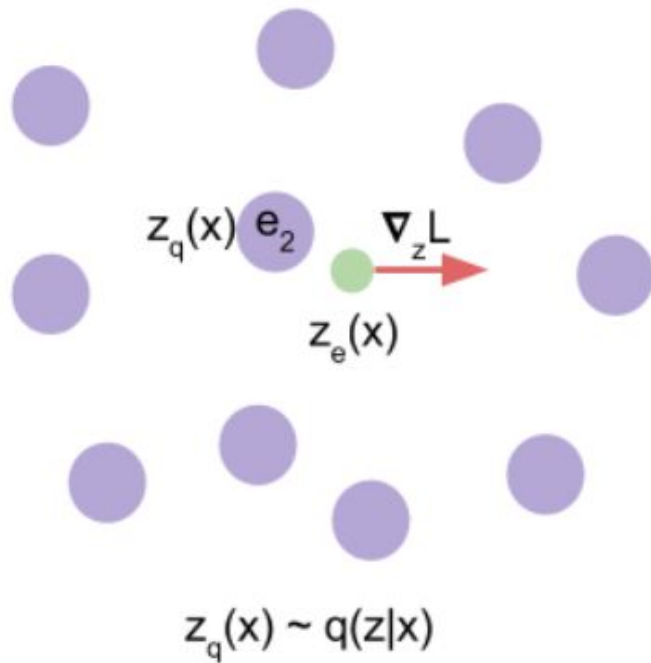
# VQ-VAE

**Related work:**
PixelVAE (Gulrajani et al, 2016)
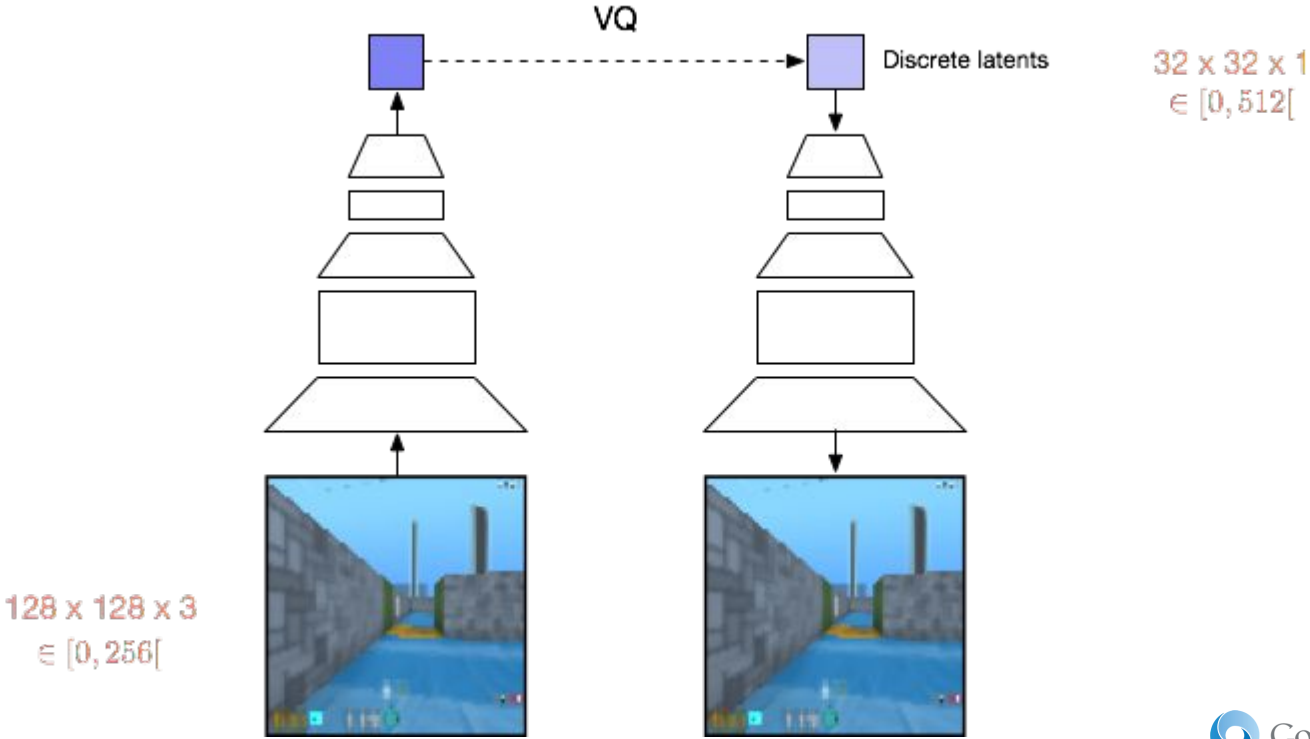Variational Lossy AutoEncoder (Chen et al, 2016)

# VQ-VAE



Embedding Space

$e_1 e_2 e_3$ ... $e_K$

$\nabla_z L$

$q(z|x)$

$z_e(x)$

$z$

$z_q(x)$

$p(x|z_q)$

Encoder

Decoder

# VQ-VAE



$z_q(x)$ $e_2$     $\nabla_z L$

$z_e(x)$

$z_q(x) \sim q(z|x)$

Google DeepMind

# Images



VQ

Discrete latents

$32 \times 32 \times 1$
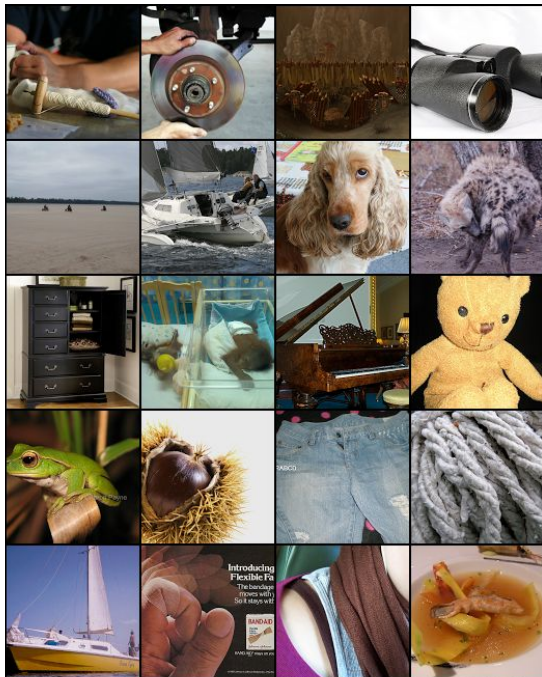$\in [0, 512[$
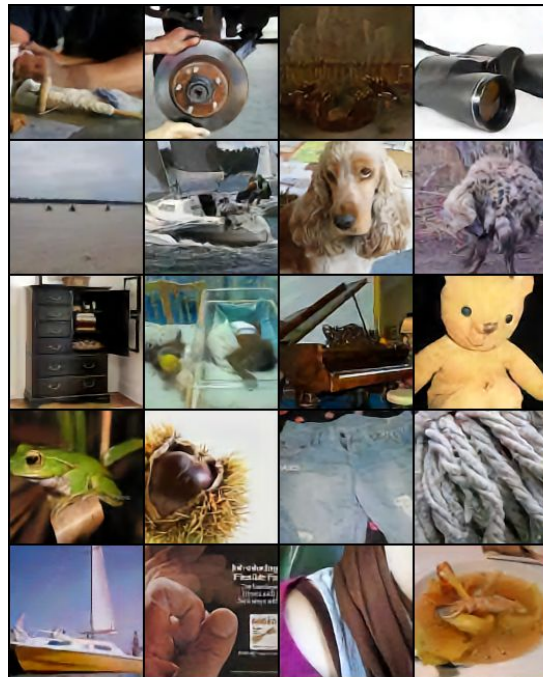
$128 \times 128 \times 3$
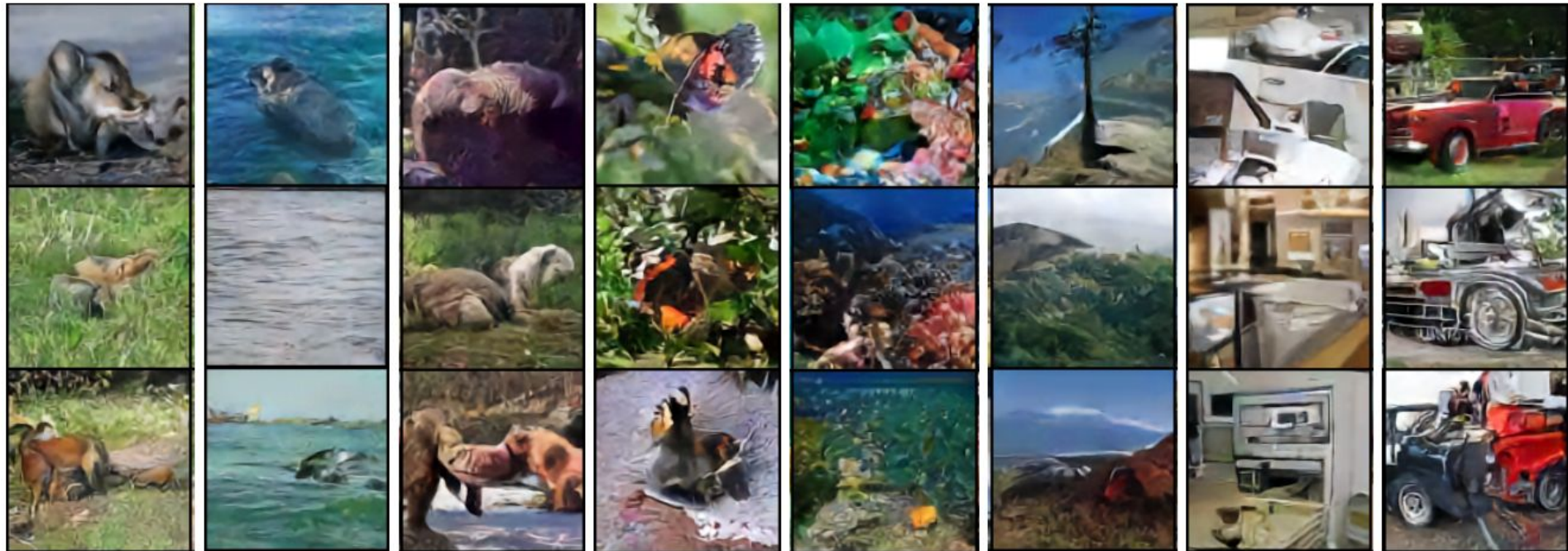$\in [0, 256[$

Google DeepMind

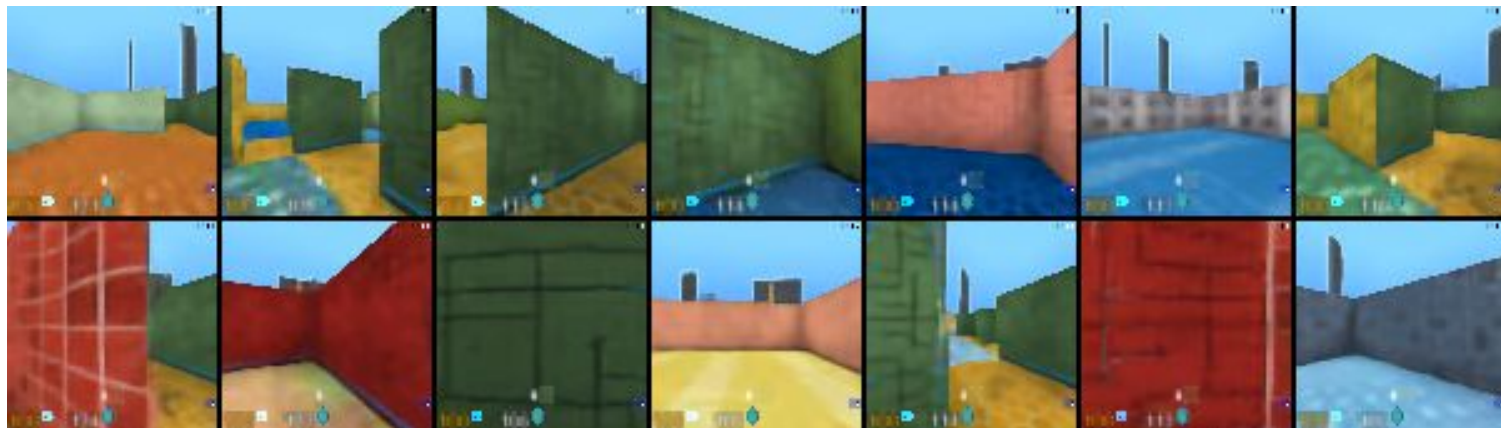# ImageNet reconstructions

Original 128x128 images

Reconstructions



Google DeepMind
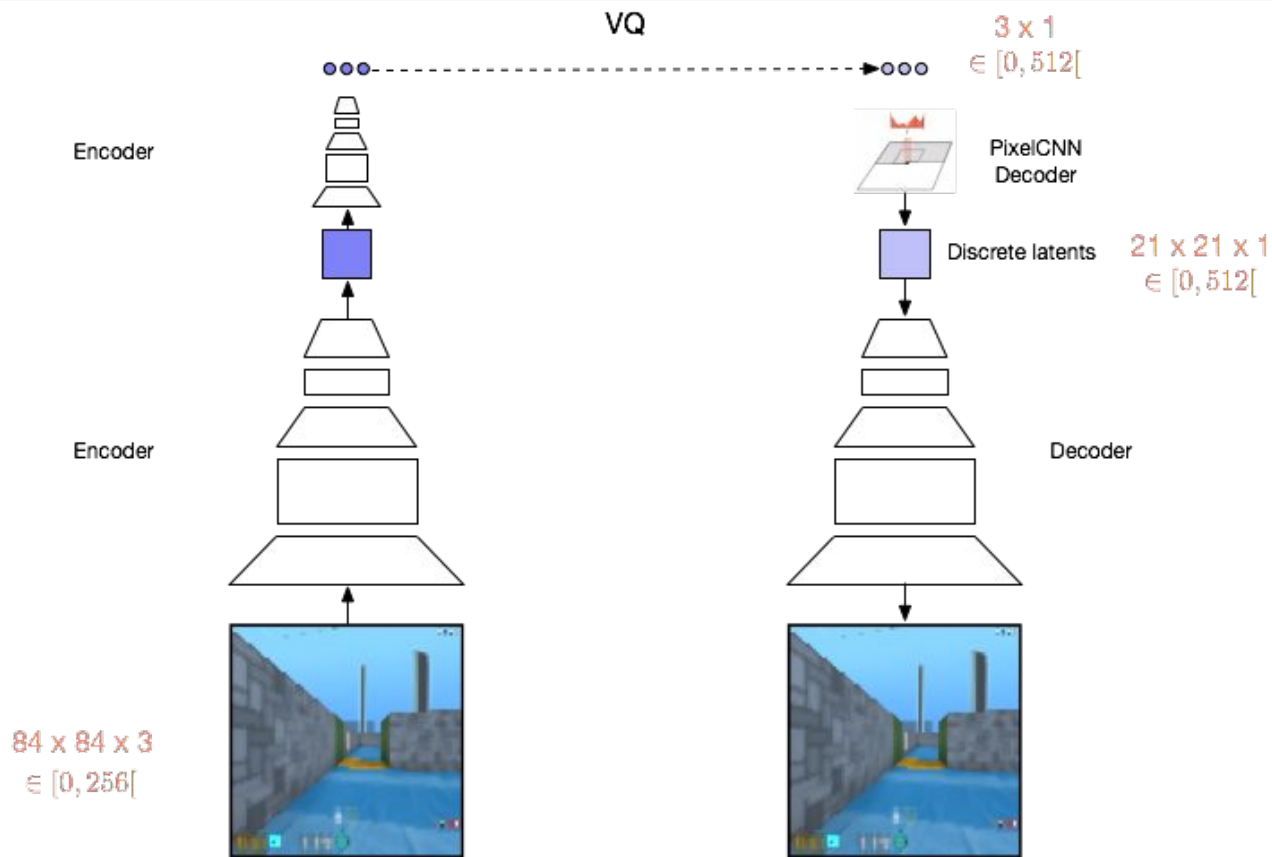
# VQ-VAE - Sample



PixelCNN Prior

Sampled Discrete latents

32 x 32 x 1

Decoder

Google DeepMind

# ImageNet samples

# DM-Lab Samples

# 3 Global Latents Reconstruction



VQ

$3 \times 1$
$\in [0, 512[$

Encoder

PixelCNN
Decoder

Discrete latents

$21 \times 21 \times 1$
$\in [0, 512[$

Encoder

Decoder

$84 \times 84 \times 3$
$\in [0, 256[$

Google DeepMind

# 3 Global Latents Reconstruction
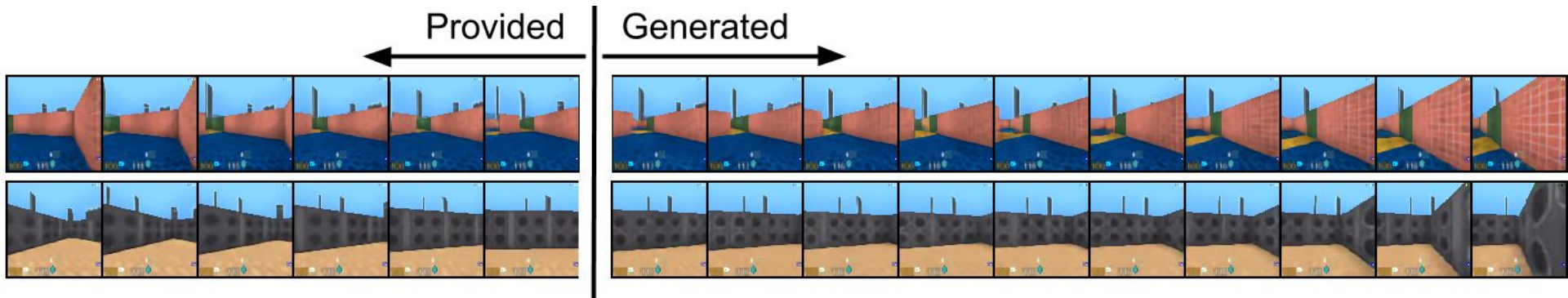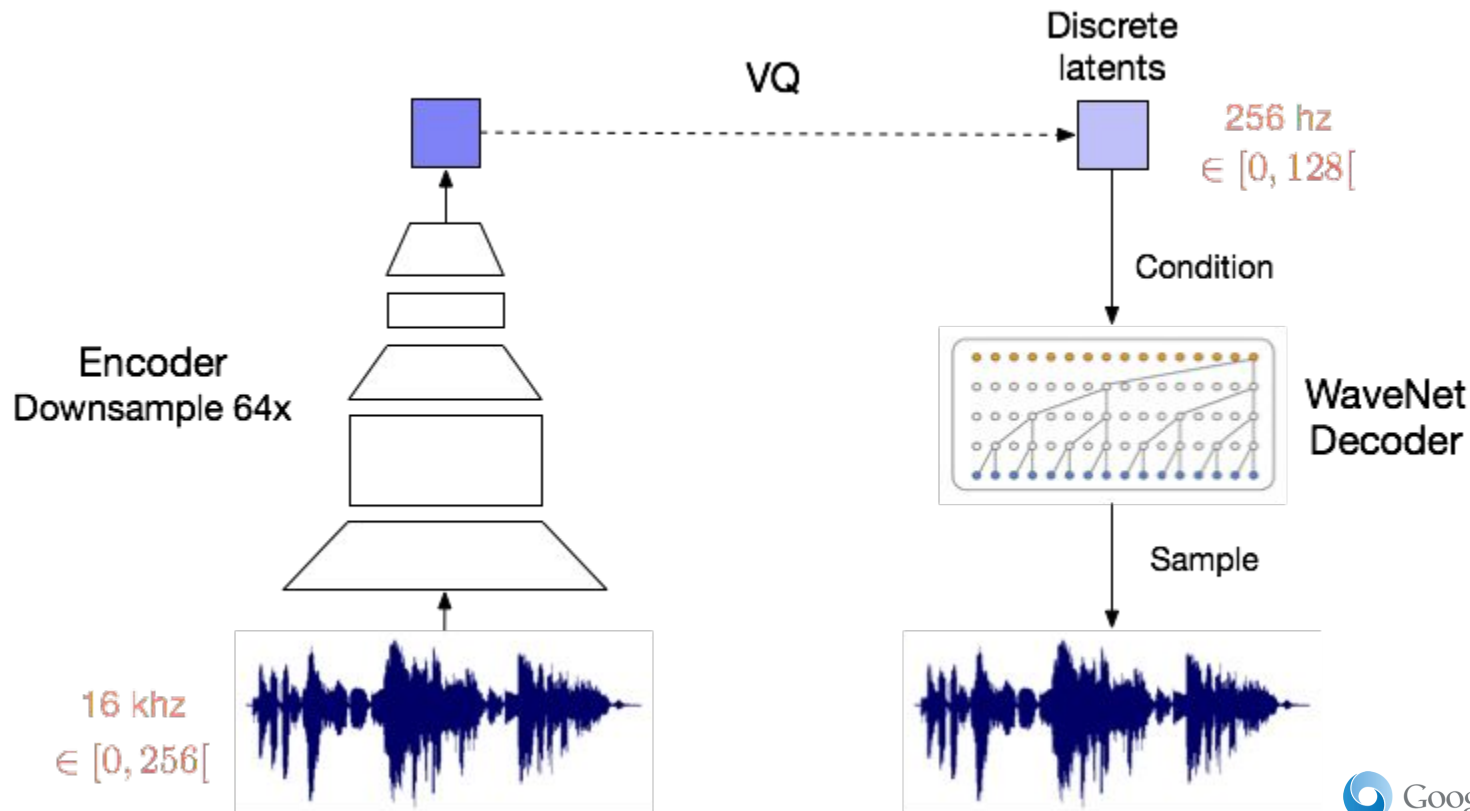
Originals



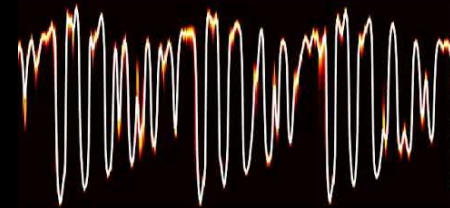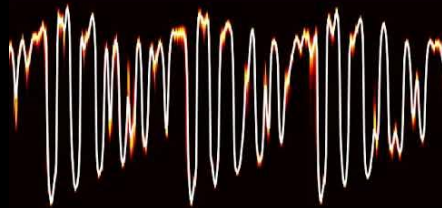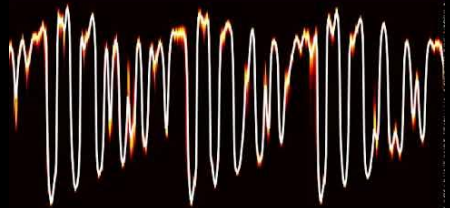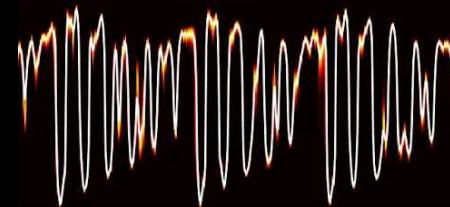Reconstructions from compressed representations (27 bits per image).



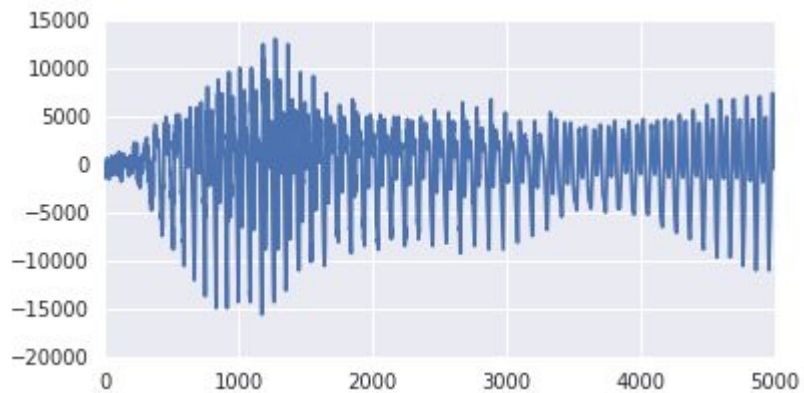Google DeepMind
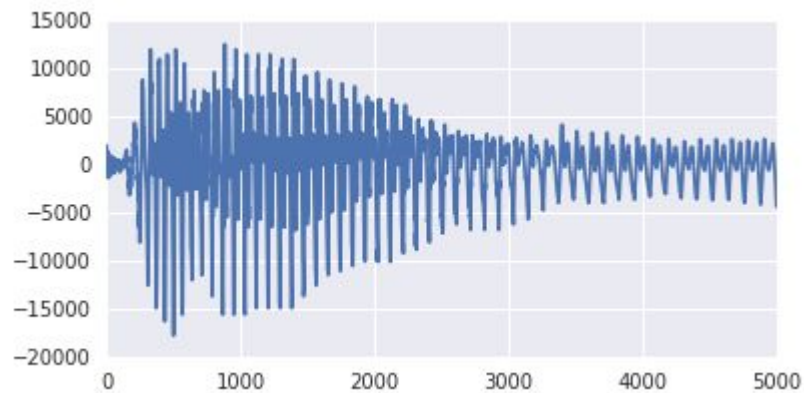
# Video Generation in the latent space



Provided ← → Generated

# Speech

# Speech - reconstruction



Original

Reconstruction

Google DeepMind

# Speech - Sample from prior

# Speech - speaker conditional

# Unsupervised Learning of phonemes

# Unsupervised Learning of phonemes

41-way classification

49.3% accuracy **fully unsupervised**
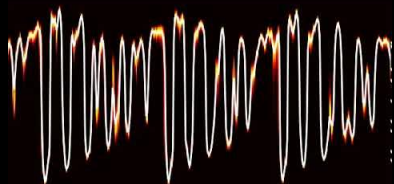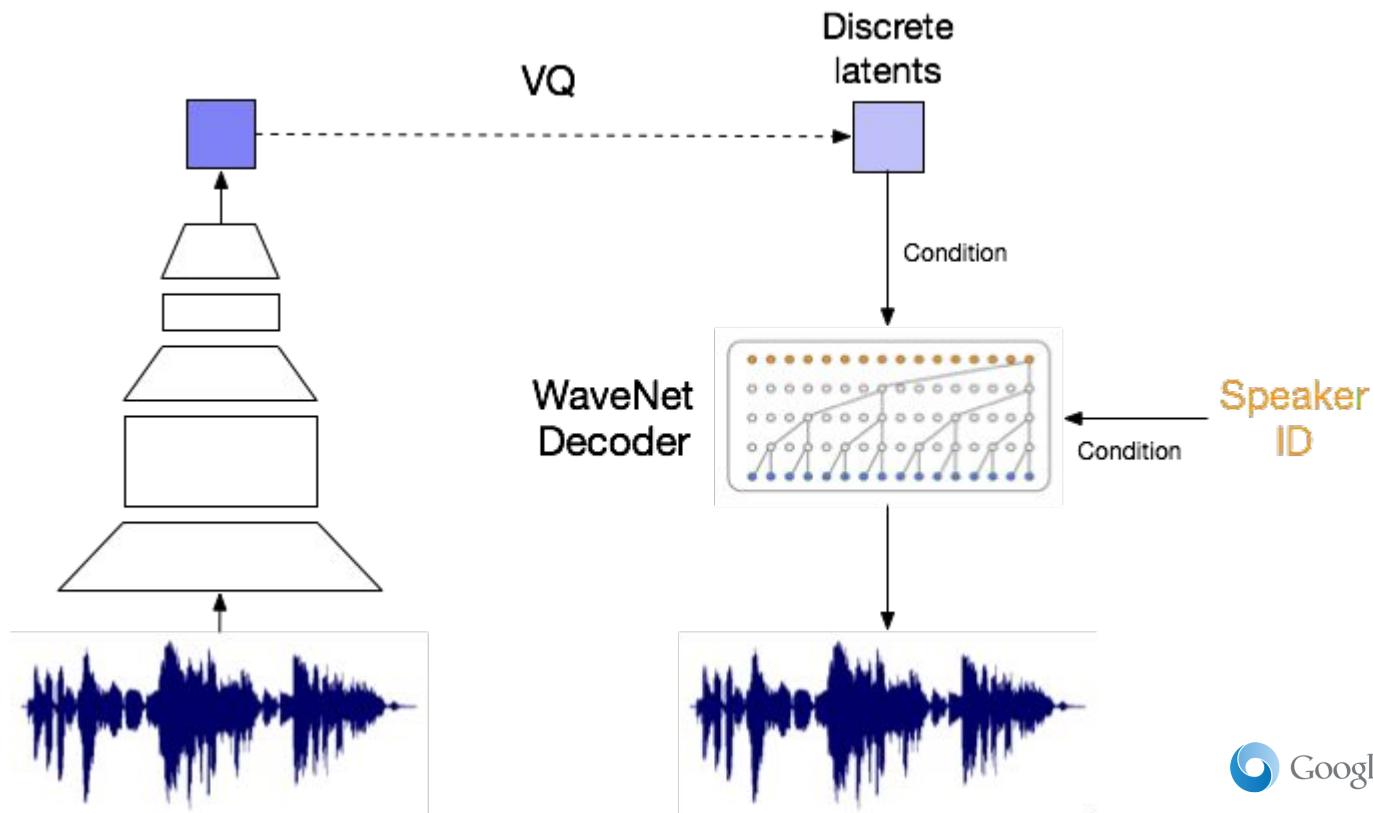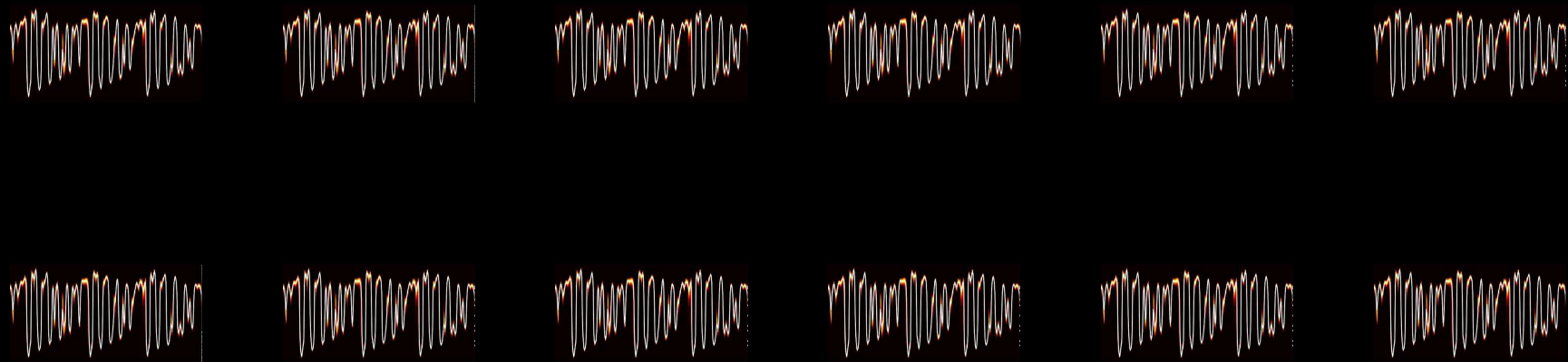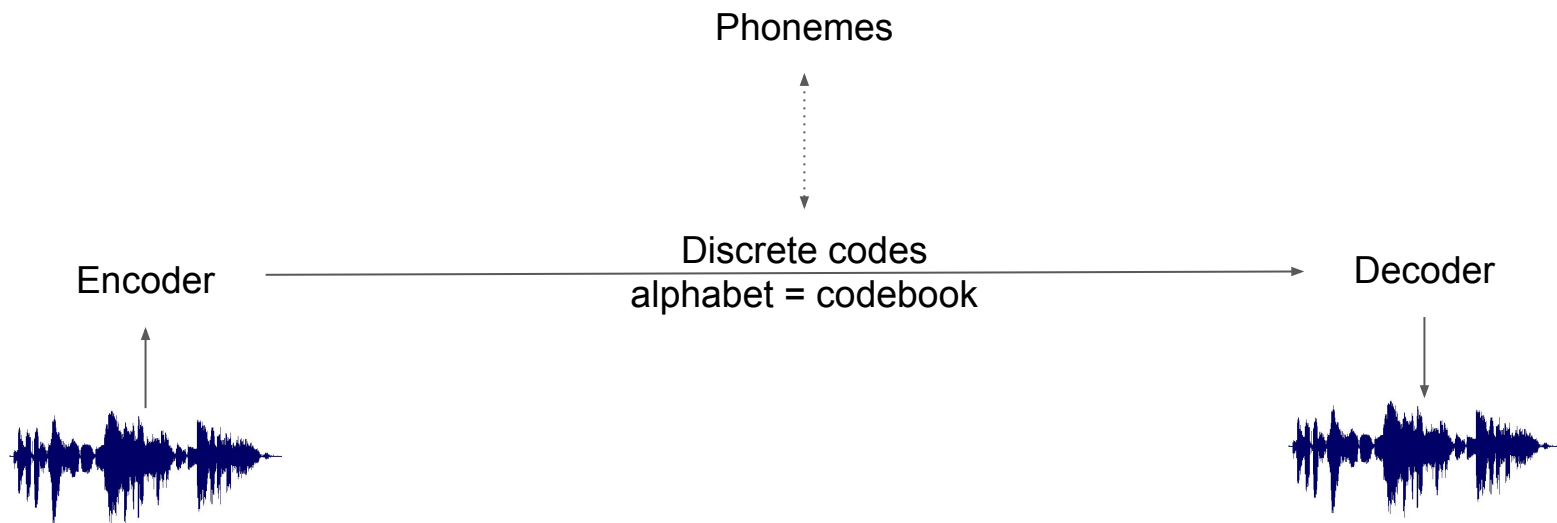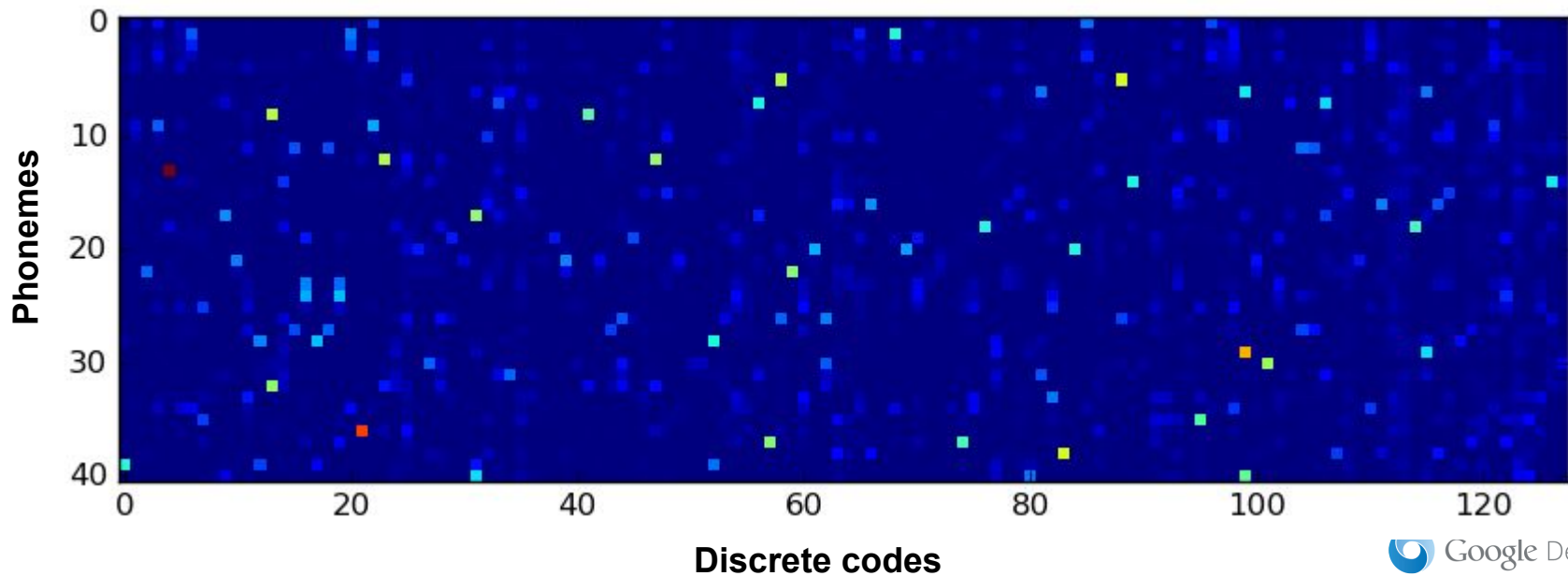


Google DeepMind

# References and related work

**Pixel Recurrent Neural Networks** - van den Oord et al, ICML 2016

**Conditional Image Generation with PixelCNN Decoders** - van den Oord et al, NIPS 2016

**<u>WaveNet: A Generative Model For Raw Audio</u>** - van den Oord et al, Arxiv 2016

**Neural Machine Translation in Linear Time** - Kalchbrenner et al, Arxiv 2016

**Video Pixel Networks** - Kalchbrenner et al, ICML 2017

**<u>Neural Discrete Representation Learning</u>** - van den Oord et al, NIPS 2017


**Related work:**

**The Neural Autoregressive Distribution Estimator** - Larochelle et al, AISTATS 2011
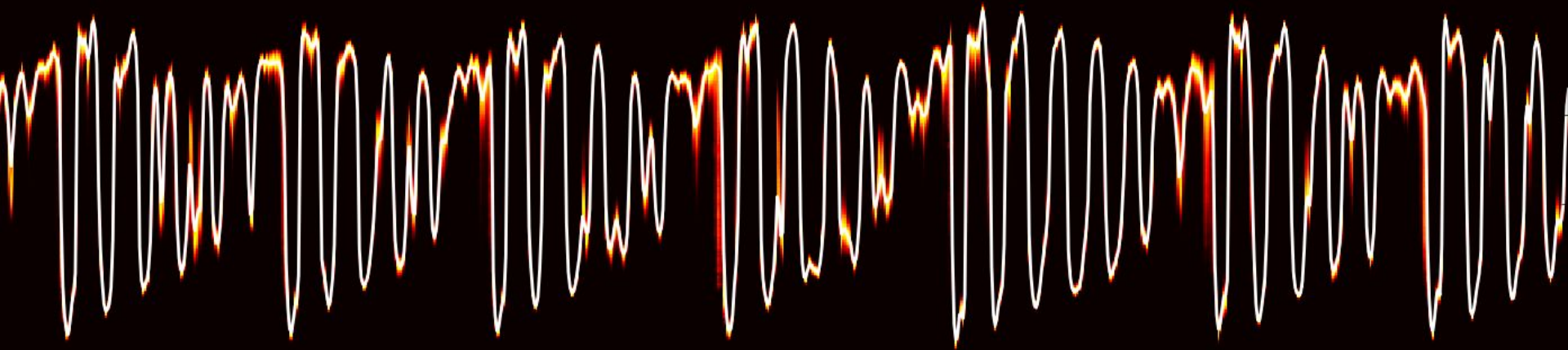
**Generative image modeling using spatial LSTMs** - Theis et al, NIPS 2015

**SampleRNN: An Unconditional End-to-End Neural Audio Generation Model** - Mehri et al, ICLR 2017

**PixelVAE: A Latent Variable Model for Natural Images** - Gulrajani et al, ICLR 2017

**Variational Lossy Autoencoder** - Chen et al, ICLR 2017

**Soft-to-Hard Vector Quantization for End-to-End Learning Compressible Representations** - Agustsson et al, NIPS 2017

Google DeepMind

Thank you!

Google DeepMind