

Generating music in the waveform domain

<https://benanne.github.io/2020/03/24/audio-generation.html>

sanderdieleman@gmail.com



@sediem

Overview

Why audio? Why raw audio?

Generative models

Likelihood-based models of raw audio

Adversarial models of raw audio

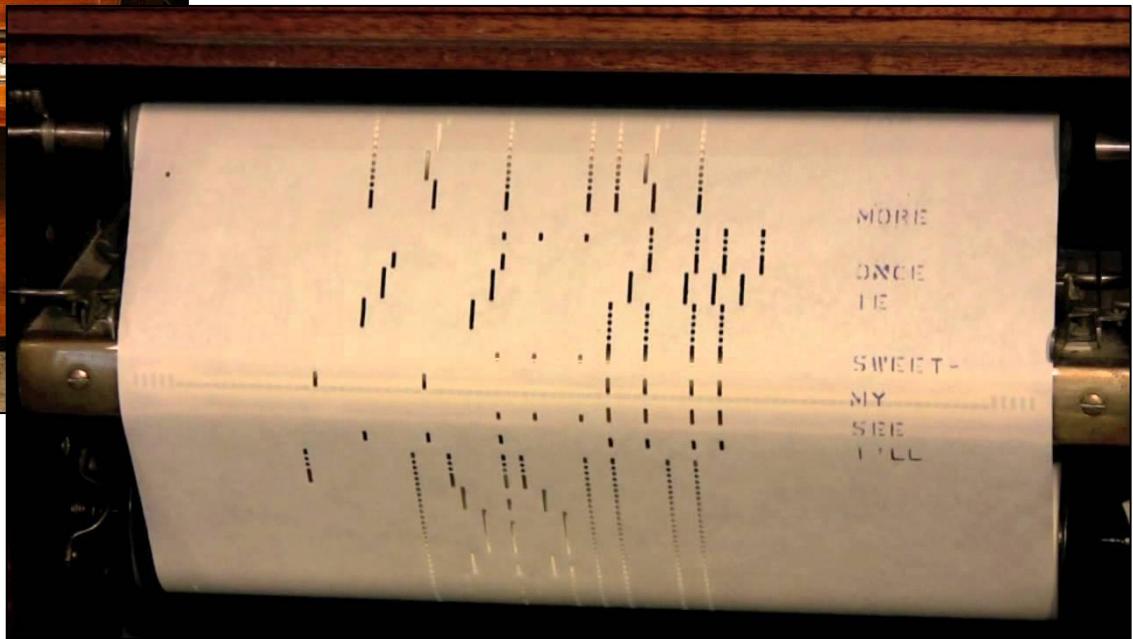
Summary

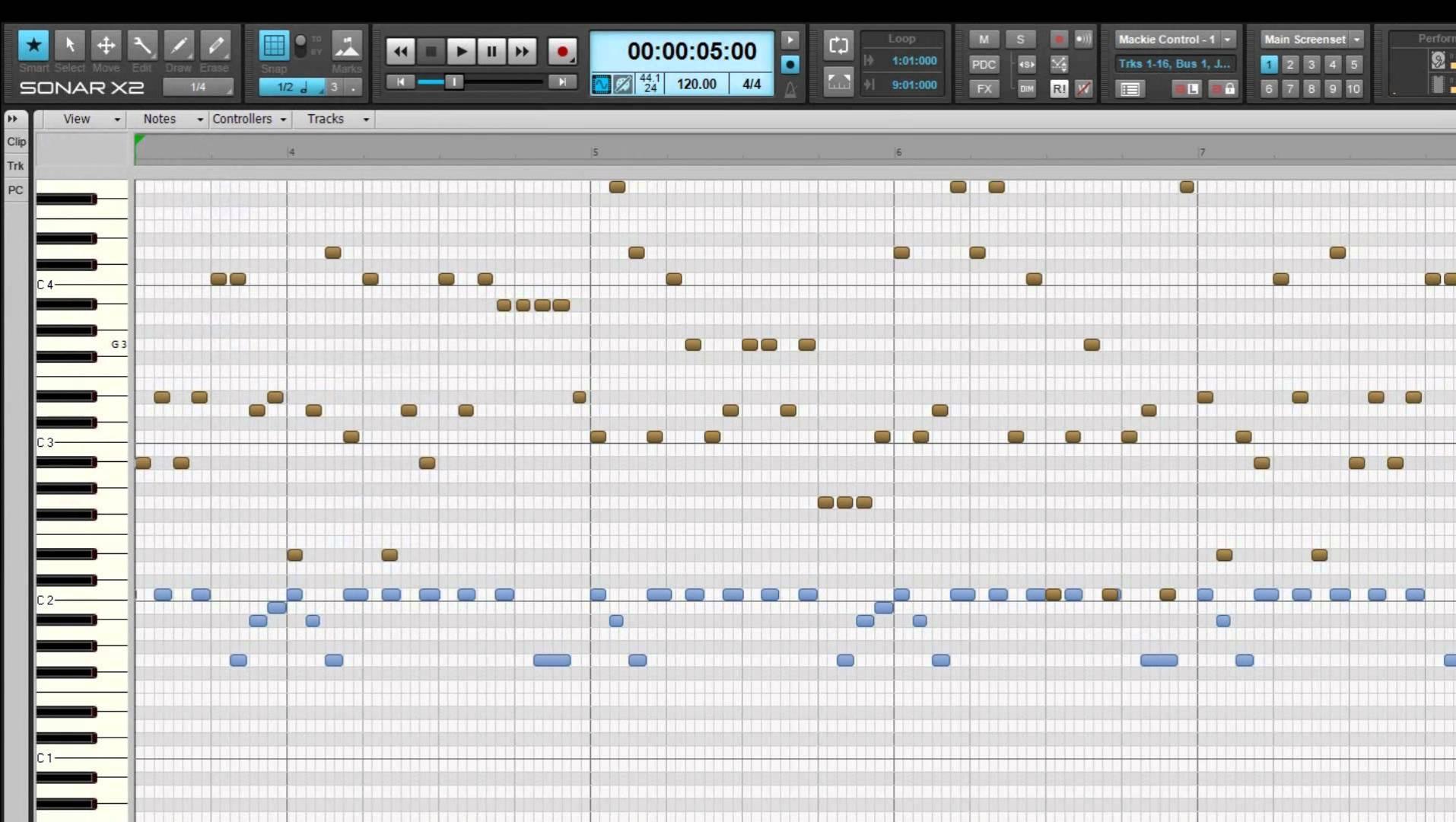
Why audio? Why raw audio?

Why audio?

Music generation is typically studied in the symbolic domain

The image shows two staves of musical notation for piano, page 3. The top staff is in treble clef and the bottom staff is in bass clef. The key signature is A major (three sharps). The tempo is indicated as *Lento, ma non troppo.* (♩ = 100). The dynamics include *p* (piano) and *cresc.* (crescendo). The first staff features sixteenth-note patterns with various fingerings (e.g., 1 2, 3 4, 5 4, 5 4 3, 5 4 2 3, 4 5) and grace notes. The second staff continues the pattern with similar fingerings and dynamics, including *stretto*, *riten.*, and *a tempo* markings. Measure numbers 53 and 54 are visible at the bottom of the page.





Why audio?

Many instruments have complex action spaces

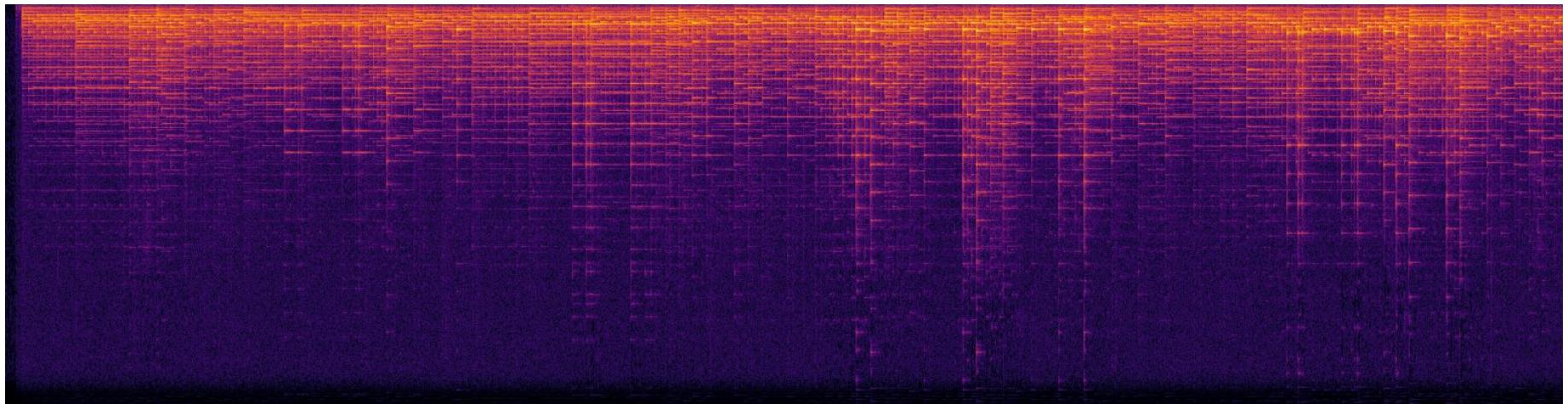
Rich palette of sounds and timbral variations

Guitar

- pick vs. finger
- picking position
- frets
- harmonics
- ...

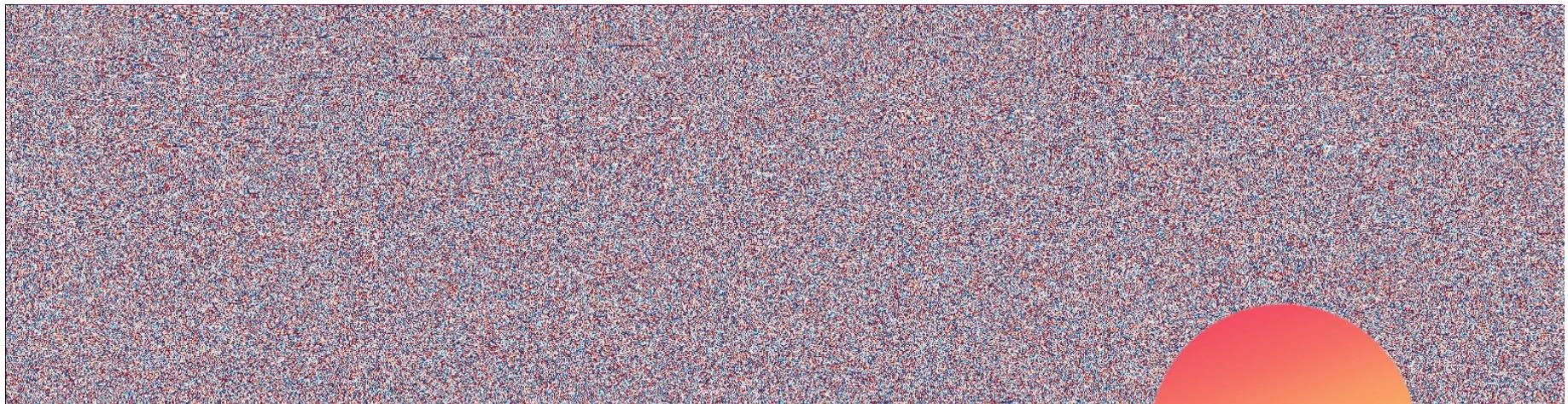


Why raw audio?

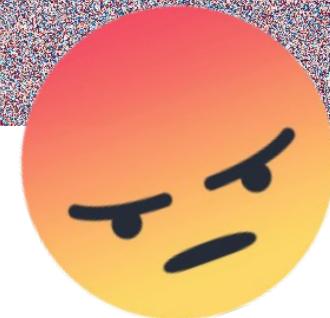


Magnitude spectrogram

Why raw audio?

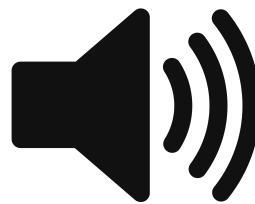


Phase spectrogram

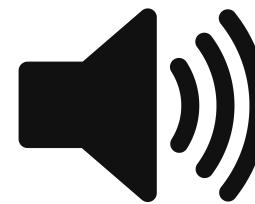


Why raw audio?

Phase is often unimportant in discriminative settings,
but is very important perceptually!



original phase

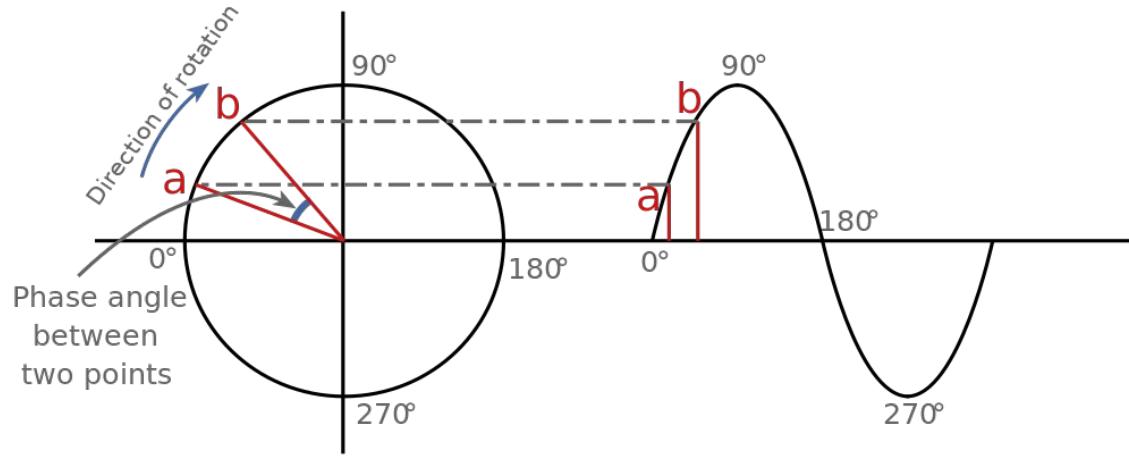
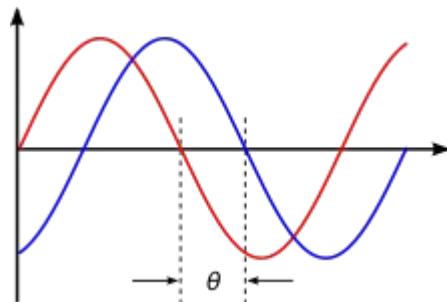


random phase

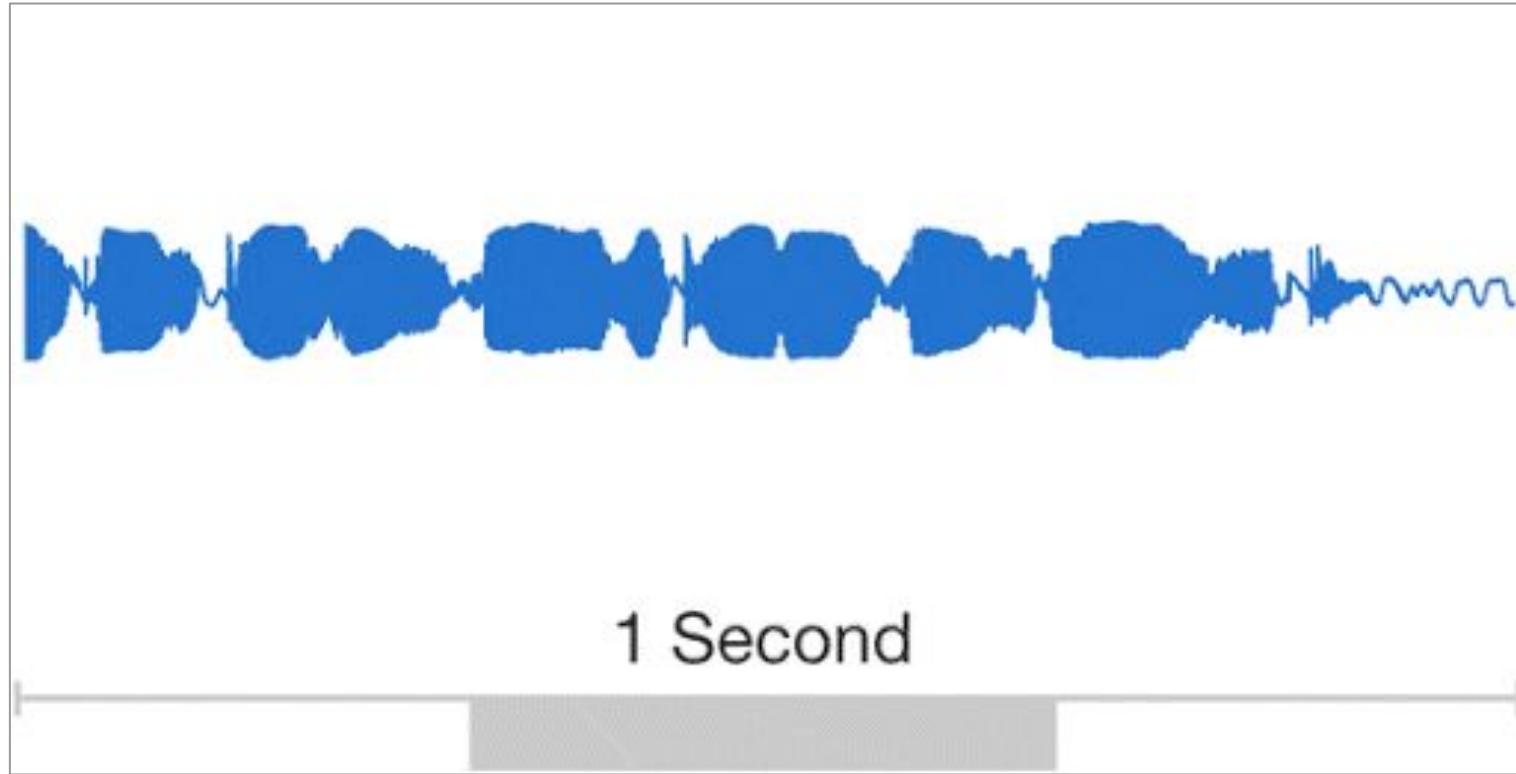
Why raw audio?

Phase is hard to model:

- it is an angle, so it wraps around
- it becomes random as the magnitude tends to 0
- absolute phase is less meaningful, but relative phase differences matter



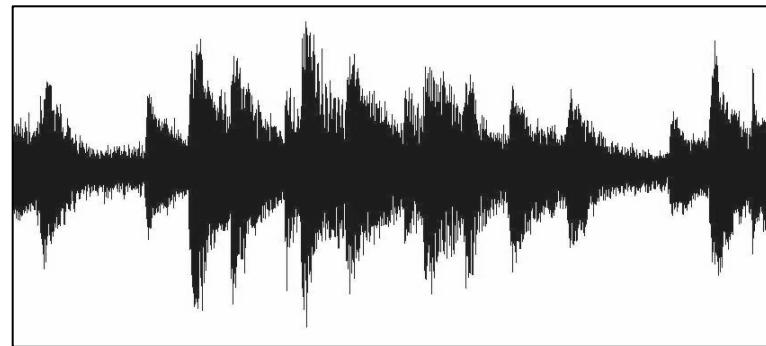
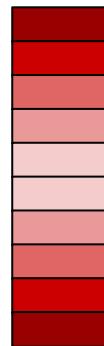
What is “raw audio” anyway?



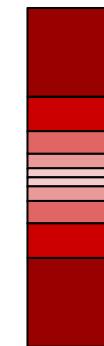
Discretising audio

- Time
- Amplitude

uniform
quantisation



μ -law
quantisation



Generative models

Generative models

Given a dataset of examples \mathbf{X} drawn from $p(\mathbf{X})$:

a generative model estimates $p(\mathbf{X})$

Generative models

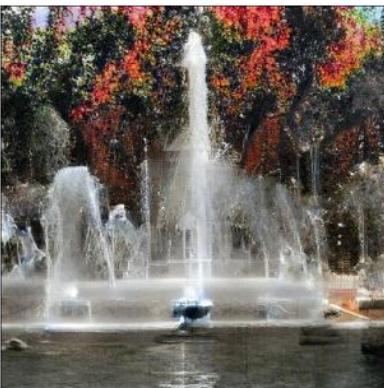
Given a dataset of examples \mathbf{X} drawn from $p(\mathbf{X})$:

a generative model estimates $p(\mathbf{X})$

Explicit: given $x \in \mathbf{X}$, model can infer $p(x)$

Implicit: model can produce new samples $x \sim p(\mathbf{X})$

Generative models



Likelihood-based models

Likelihood-based models parameterise $p(X)$ directly

Objective function: maximise $\sum_{x \in X} \log p(x)$

Autoregressive models

Autoregressive models factorise $p(X)$ into simpler (scalar) distributions

$$x = (x_1, x_2, x_3, \dots, x_n)$$

$$p(x) = \prod_i p(x_i | x_{<i}) \quad \text{chain rule of probability}$$

We can use the same model $p(x_i | x_{<i})$ for all i !

Flow-based models

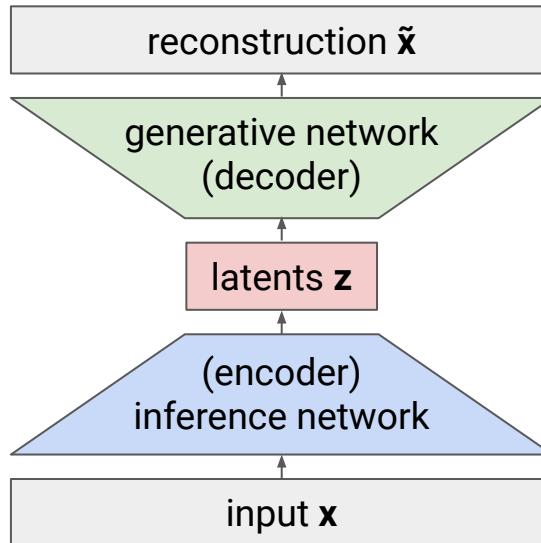
Flow-based models transform $p(\mathbf{X})$ to a simple (factorised) distribution with an invertible mapping

$$p(\mathbf{x}) = p(\mathbf{z}) \cdot |\det \mathbf{J}|^{-1} \quad \text{change of variables theorem}$$
$$\mathbf{J} = d(g(\mathbf{z}))/d\mathbf{z} \quad \mathbf{x} = g(\mathbf{z})$$

Important constraints:

$g(\mathbf{z})$ must be invertible $\det \mathbf{J}$ must be tractable

Variational autoencoders (VAEs)



$$p_\varphi(x|z)$$

$$z \sim \mathcal{N}(0,1)$$

$$q_\theta(z|x)$$

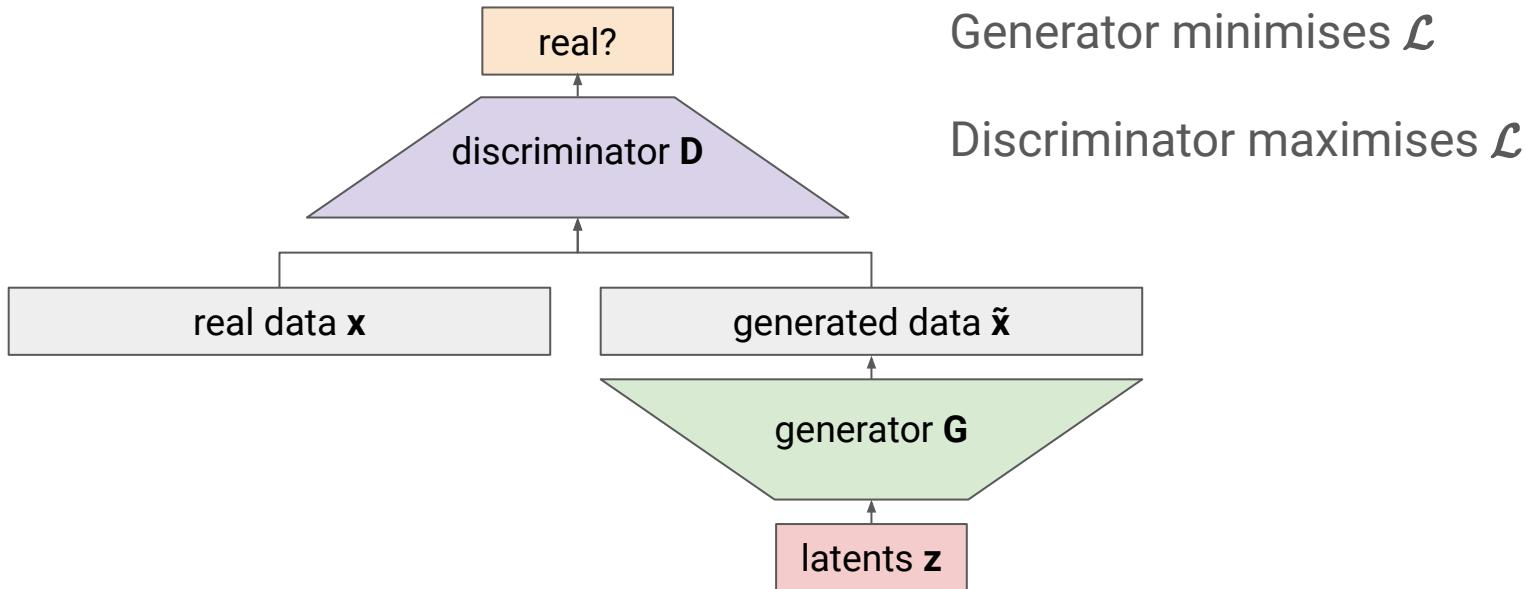
VAE maps latents z from a simple distribution to x with a (non-invertible) generative network.

The inference network approximates the inverse operation.

$p(x)$ cannot be computed exactly, the Evidence Lower BOund (ELBO) is maximised instead.

Adversarial models

$$\mathcal{L} = \mathbb{E}_x[\log(D(x))] + \mathbb{E}_z[\log(1 - D(G(z)))]$$



More exotic flavours

- Implicit quantile networks
- Energy-based models
- Optimal transport (e.g. Wasserstein autoencoders)
- Score-based generative modelling
- Maximum mean discrepancy (energy distance)
- ...

Conditional generative models

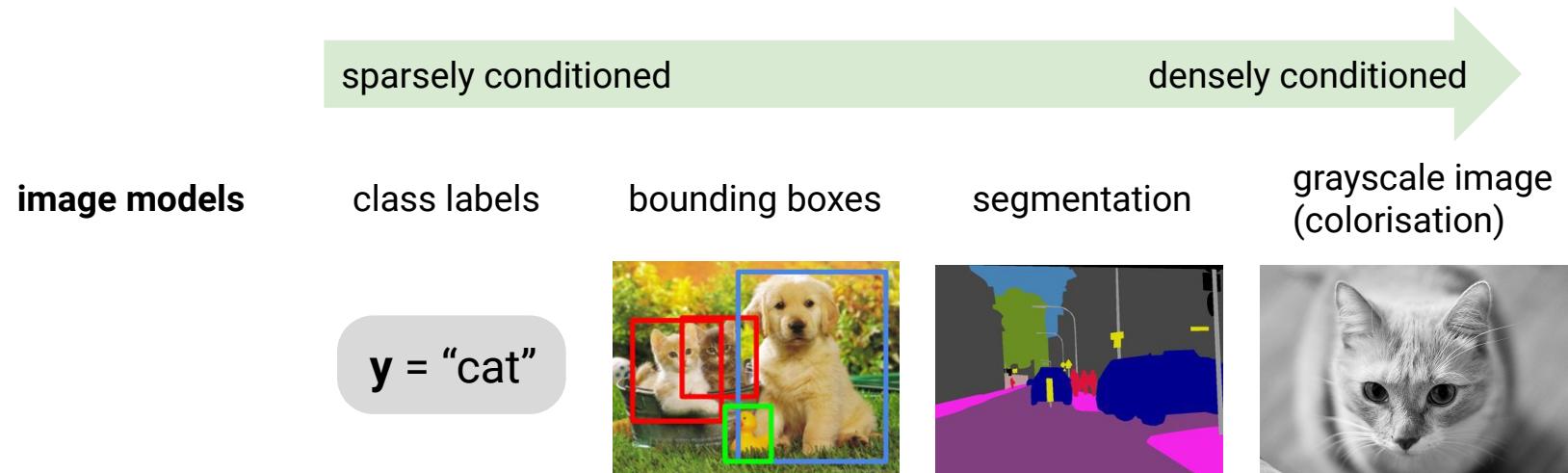
Conditioning is “side information” which allows for control over the model output

$p(x|c)$ vs. $p(x)$

Conditional generative models

Conditioning is “side information” which allows for control over the model output

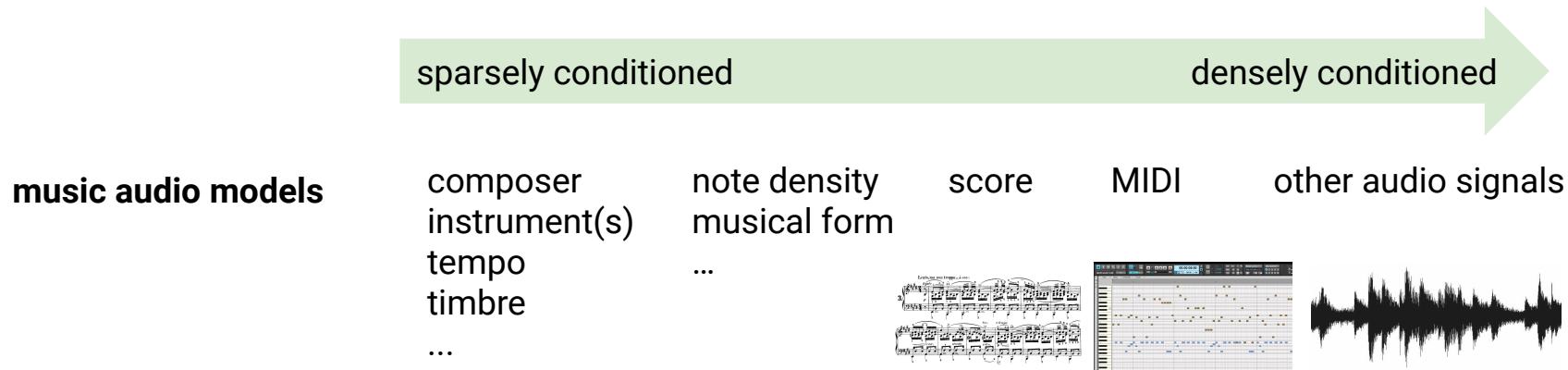
$p(x|c)$ vs. $p(x)$



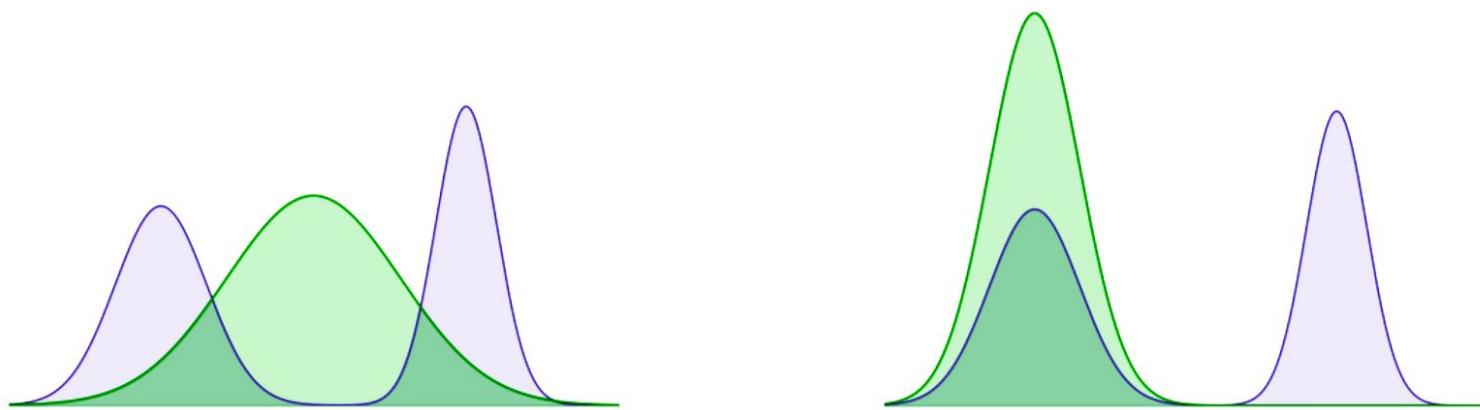
Conditional generative models

Conditioning is “side information” which allows for control over the model output

$p(x|c)$ vs. $p(x)$



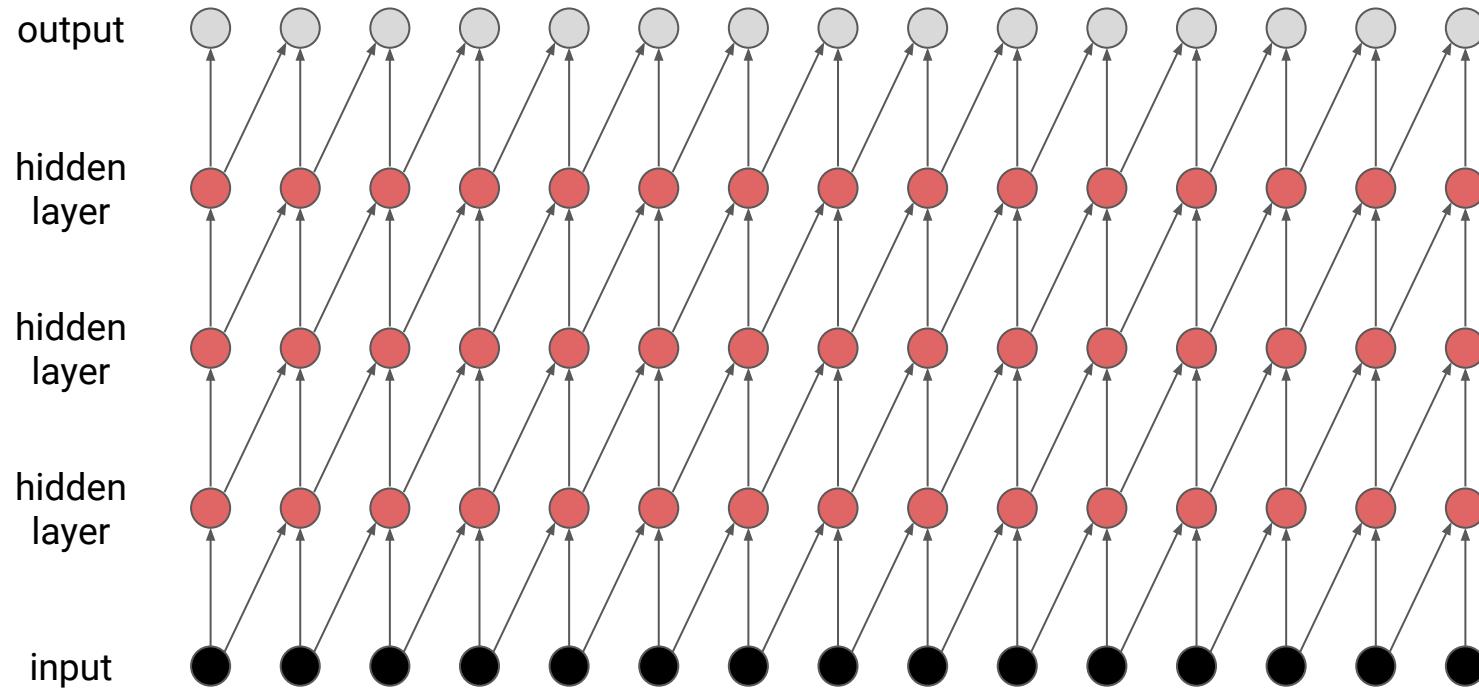
Mode-covering vs. mode-seeking behaviour



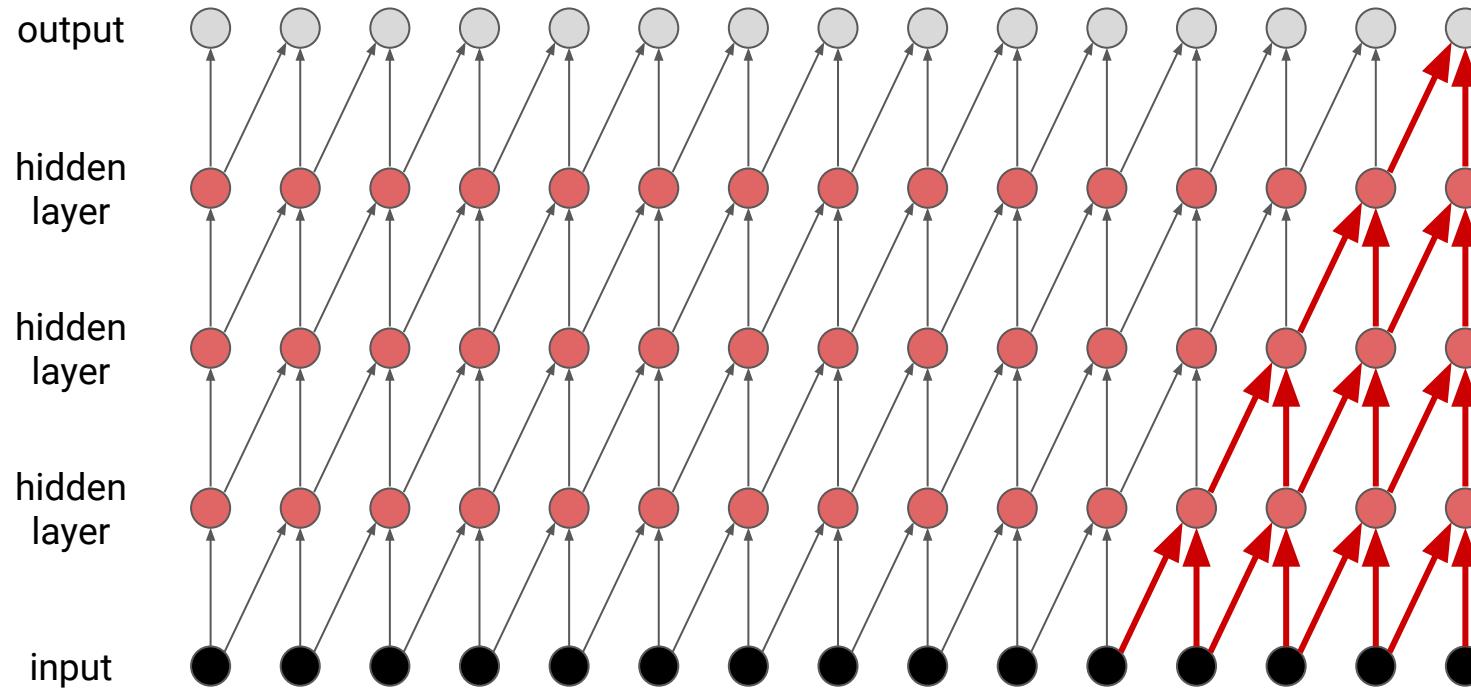
- Likelihood-based models are mode-covering
- Adversarial models are (typically) mode-seeking
- In more densely conditioned settings, we tend to care less about covering all the modes

Likelihood-based models

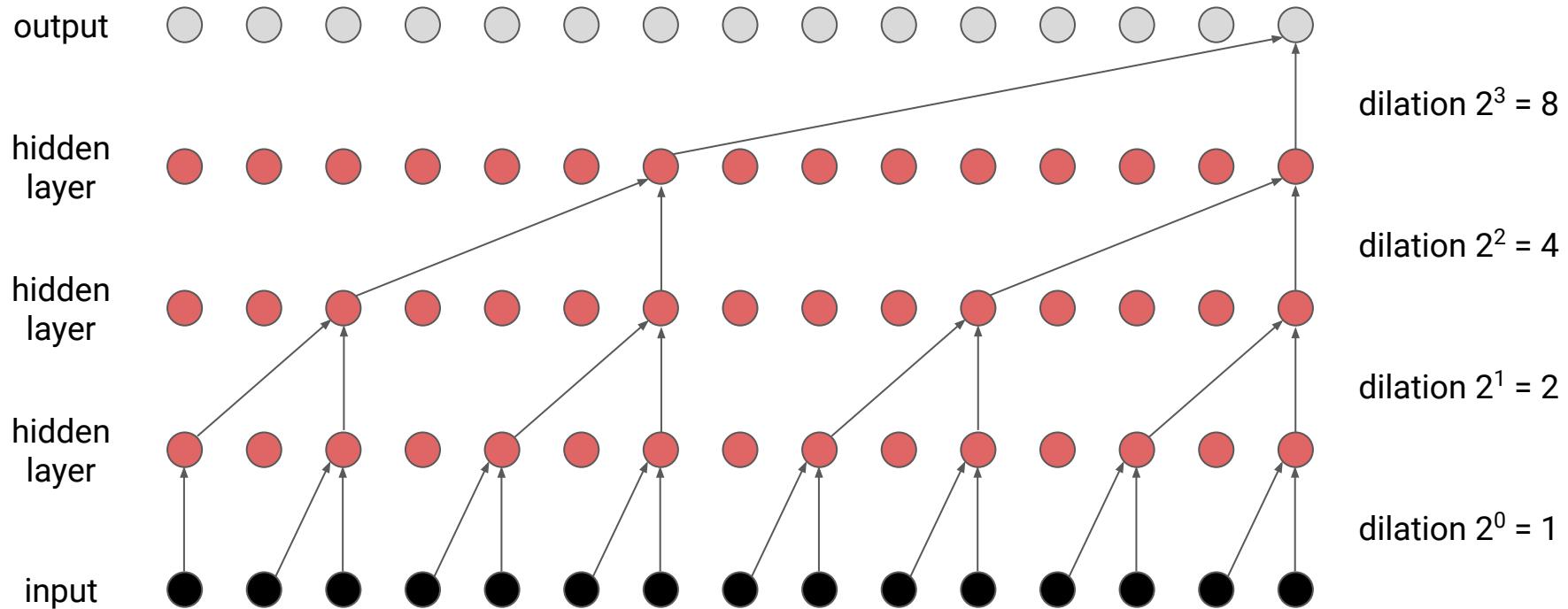
WaveNet



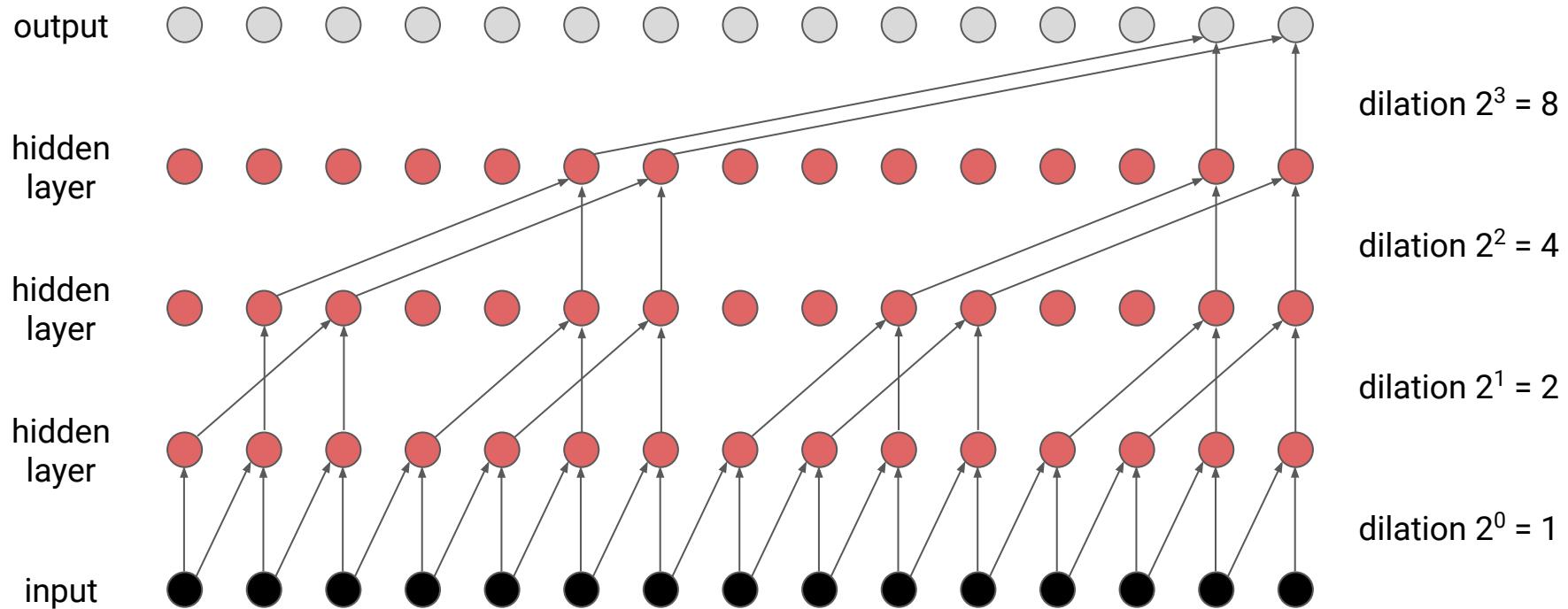
WaveNet

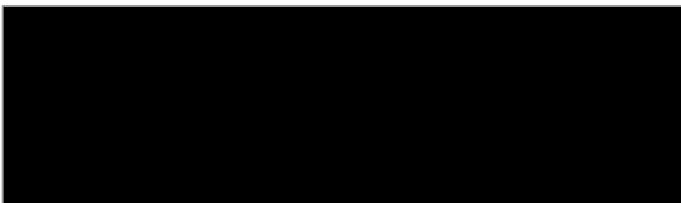
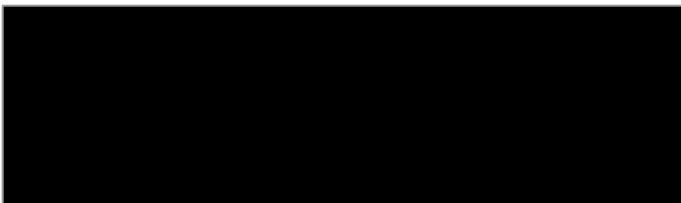


WaveNet: dilated convolutions

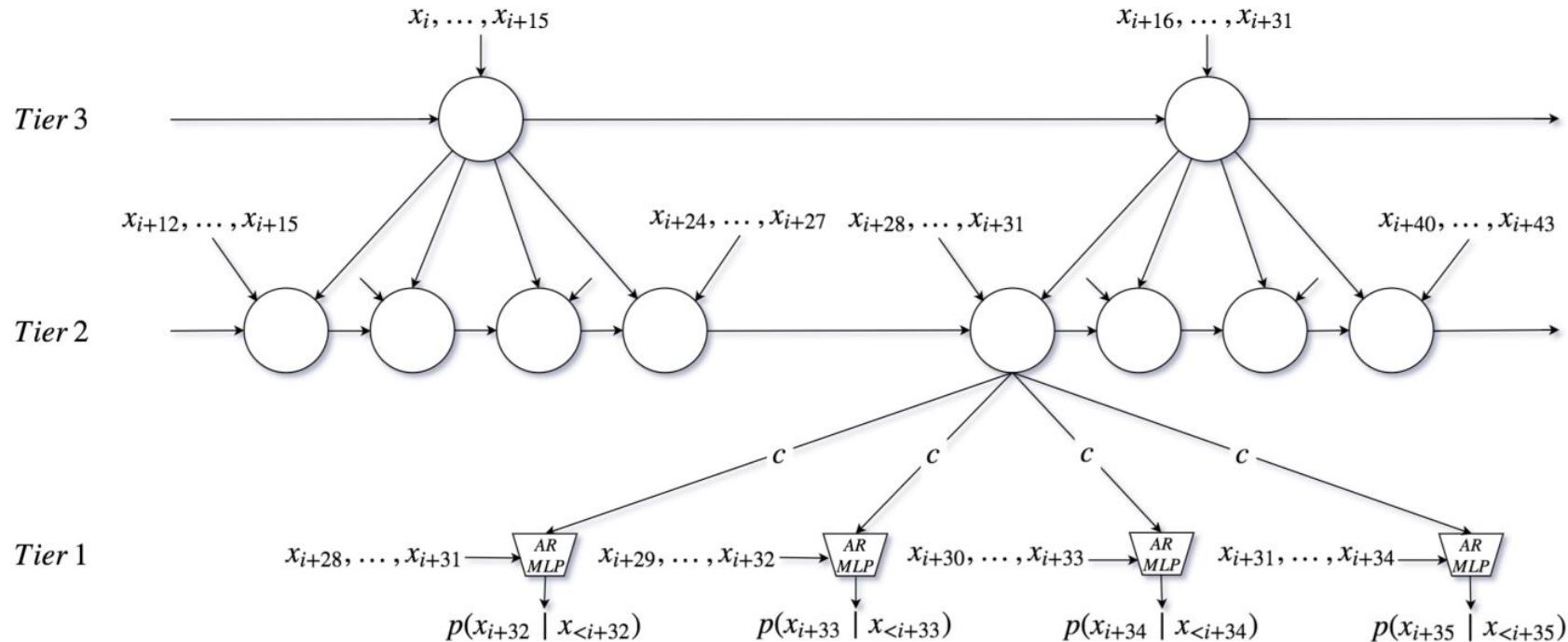


WaveNet: dilated convolutions

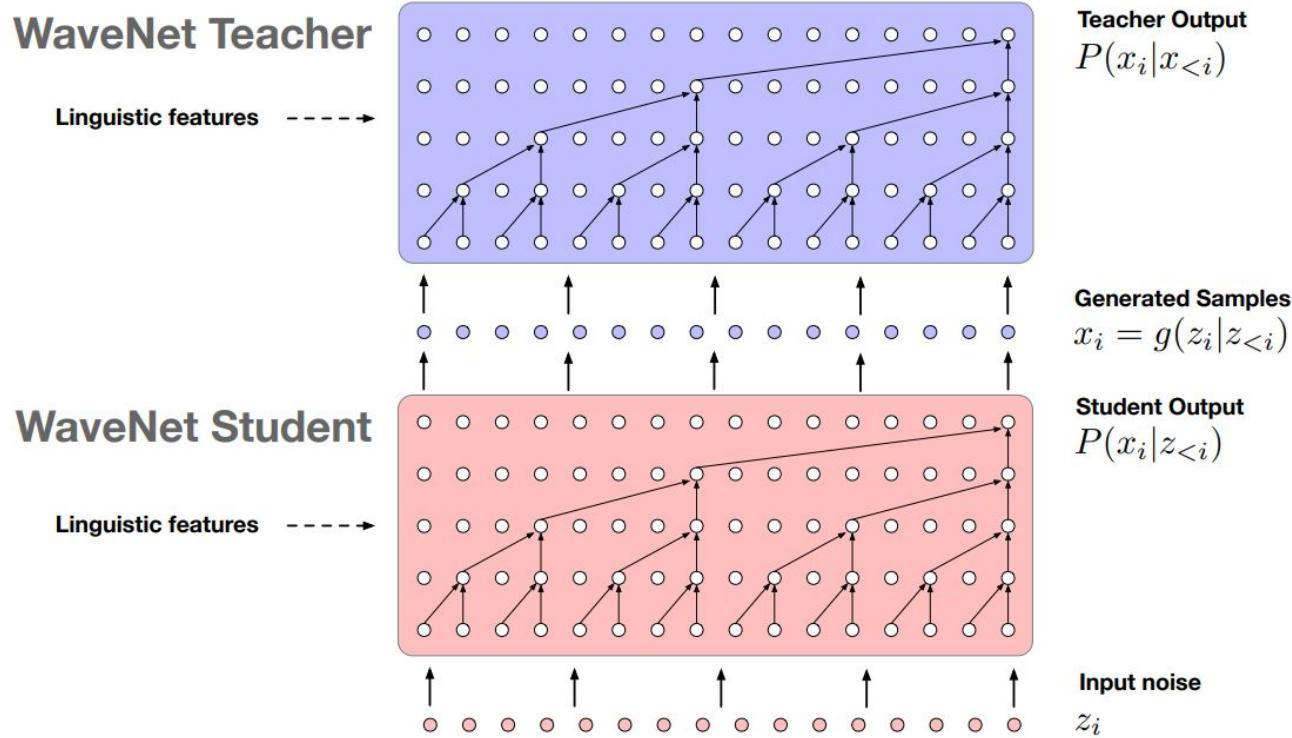




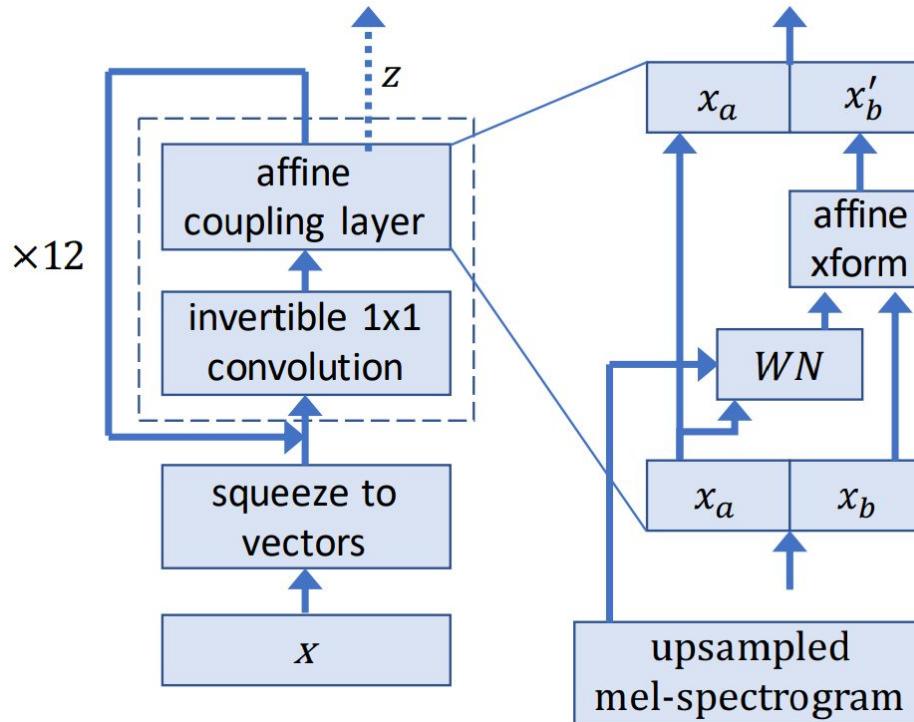
SampleRNN



Parallel WaveNet, ClariNet



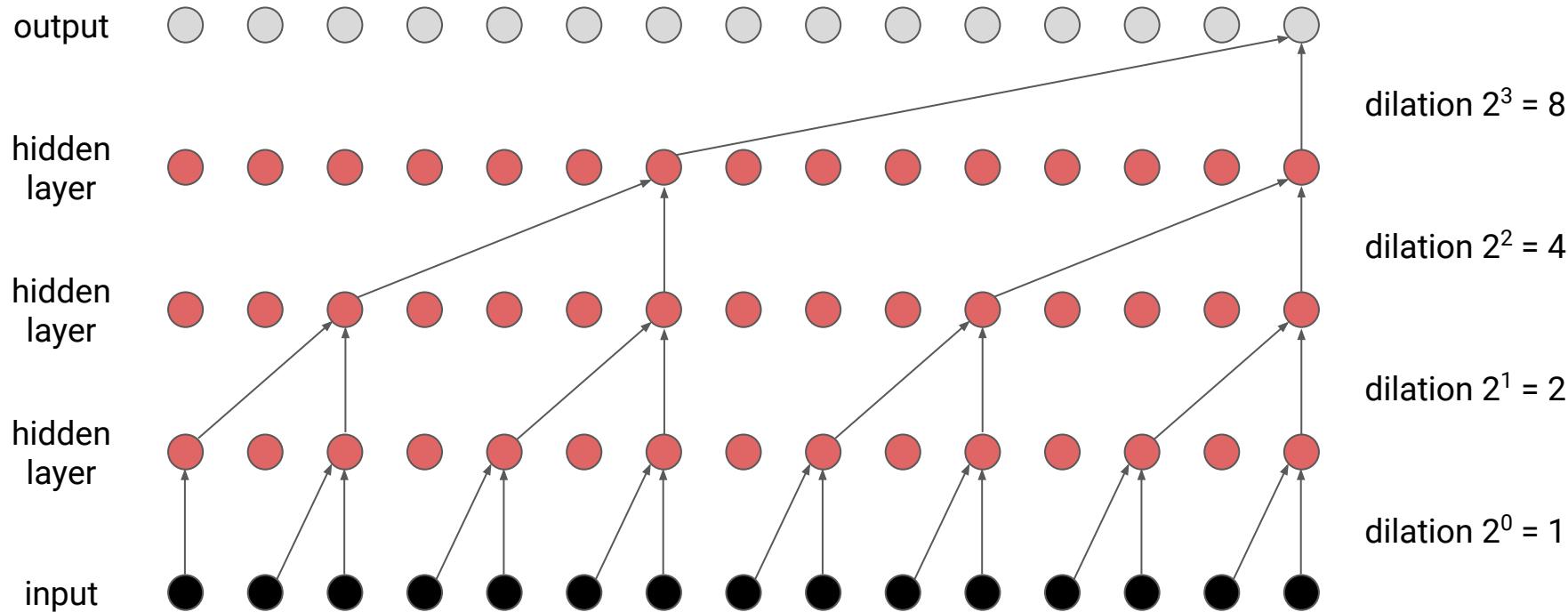
WaveGlow, FloWaveNet



Prenger et al., 2019. "Waveglow: A flow-based generative network for speech synthesis", ICASSP.

Kim et al., 2019. "FloWaveNet: A generative flow for raw audio", ICML.

WaveNet: #layers ~ log(receptive field length)



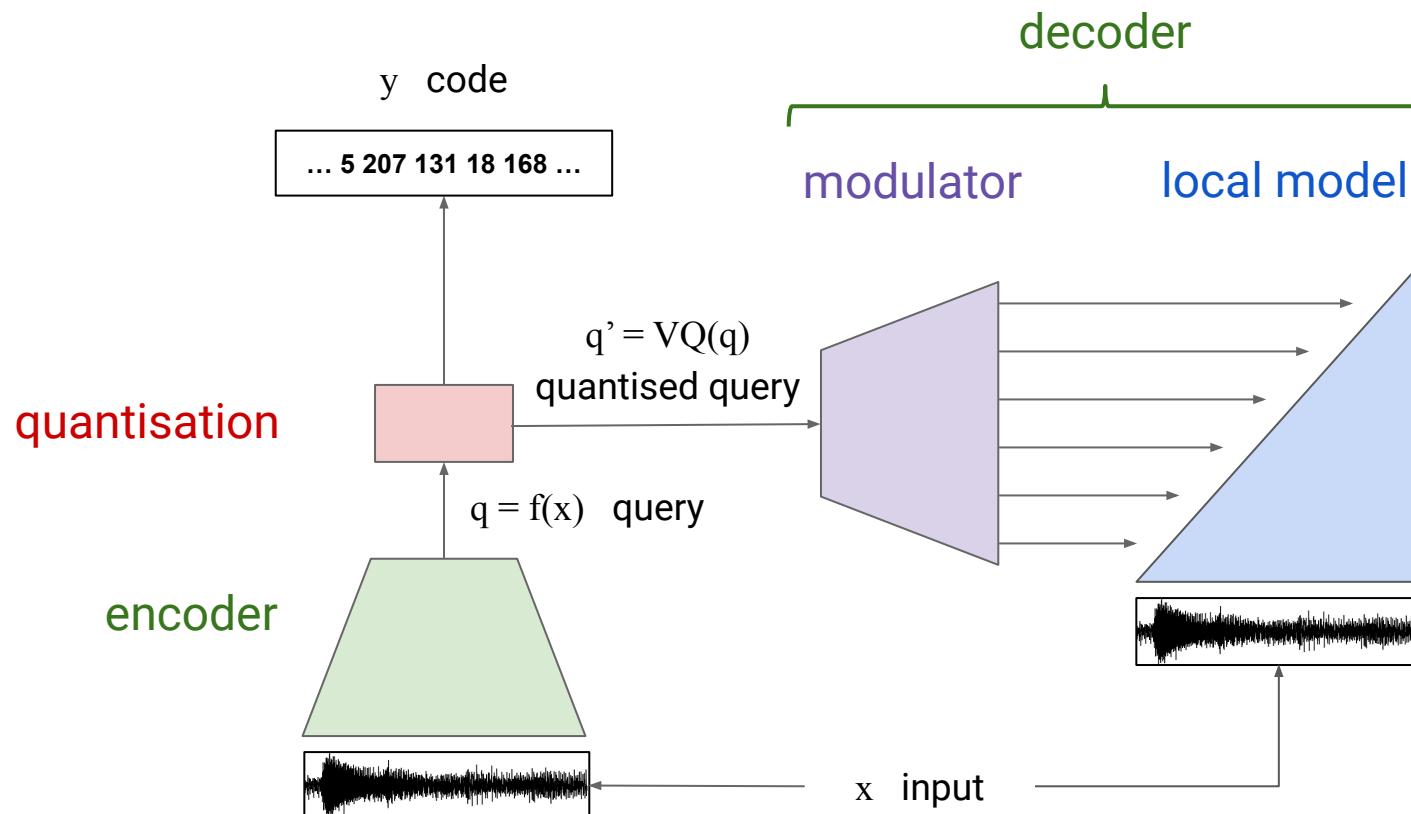
... but memory usage \sim receptive field length

Required model depth is **logarithmic** in the desired receptive field length

Required memory usage during training is still **linear** in the desired receptive field length!

⇒ We cannot scale indefinitely using dilation

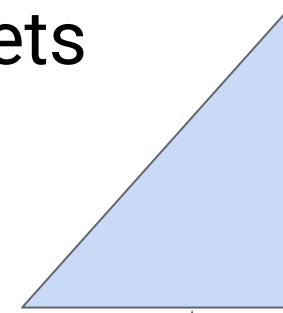
Autoregressive discrete autoencoders



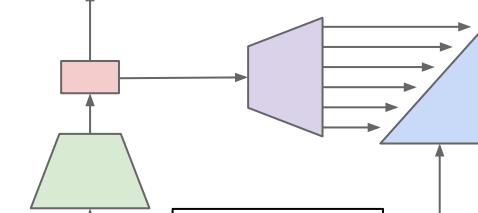


Hierarchical WaveNets

level 3
unconditional model

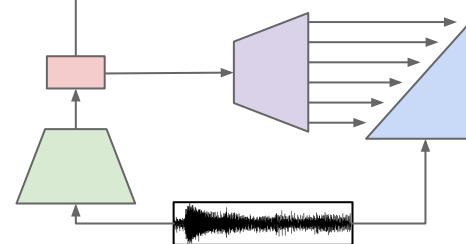


level 2
ADA



250 Hz

level 1
ADA

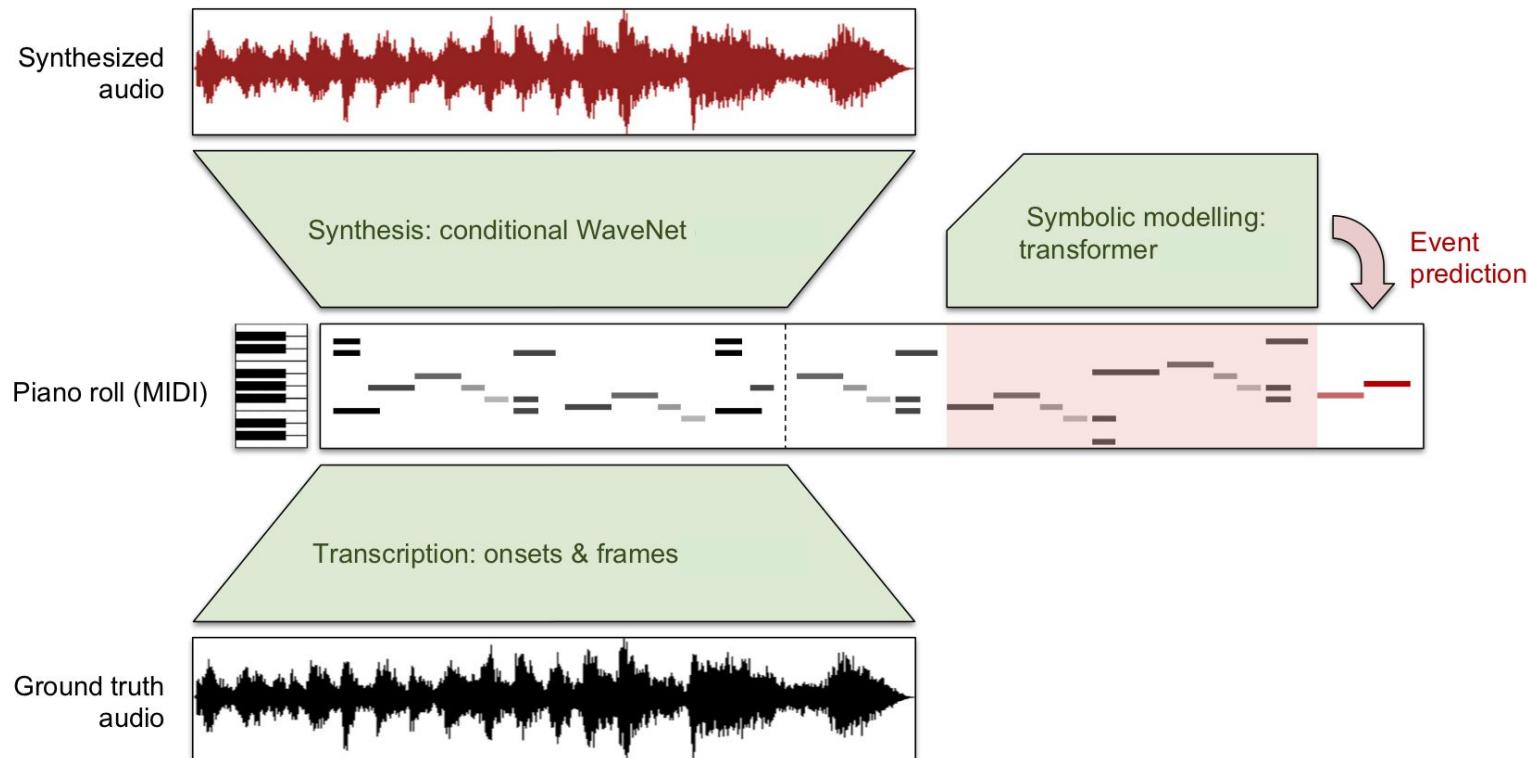


2 kHz

16 kHz

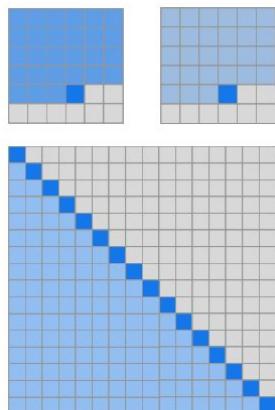


Wave2Midi2Wave and the MAESTRO dataset

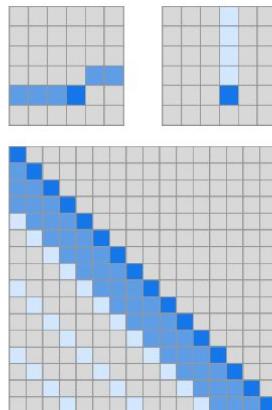


Sparse transformers

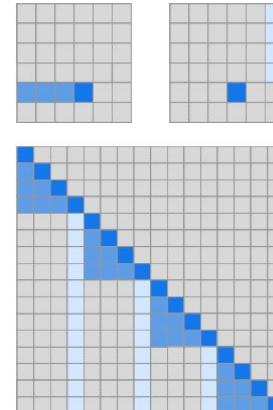
Attention (with sparse masks) instead of recurrence / convolutions.



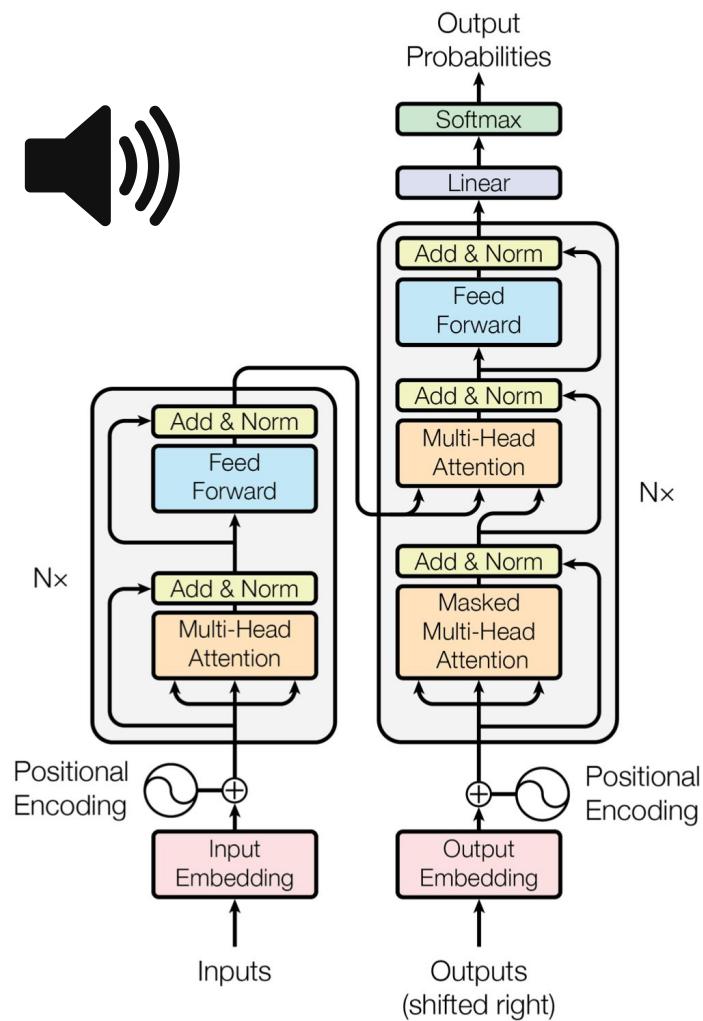
(a) Transformer



(b) Sparse Transformer (strided)

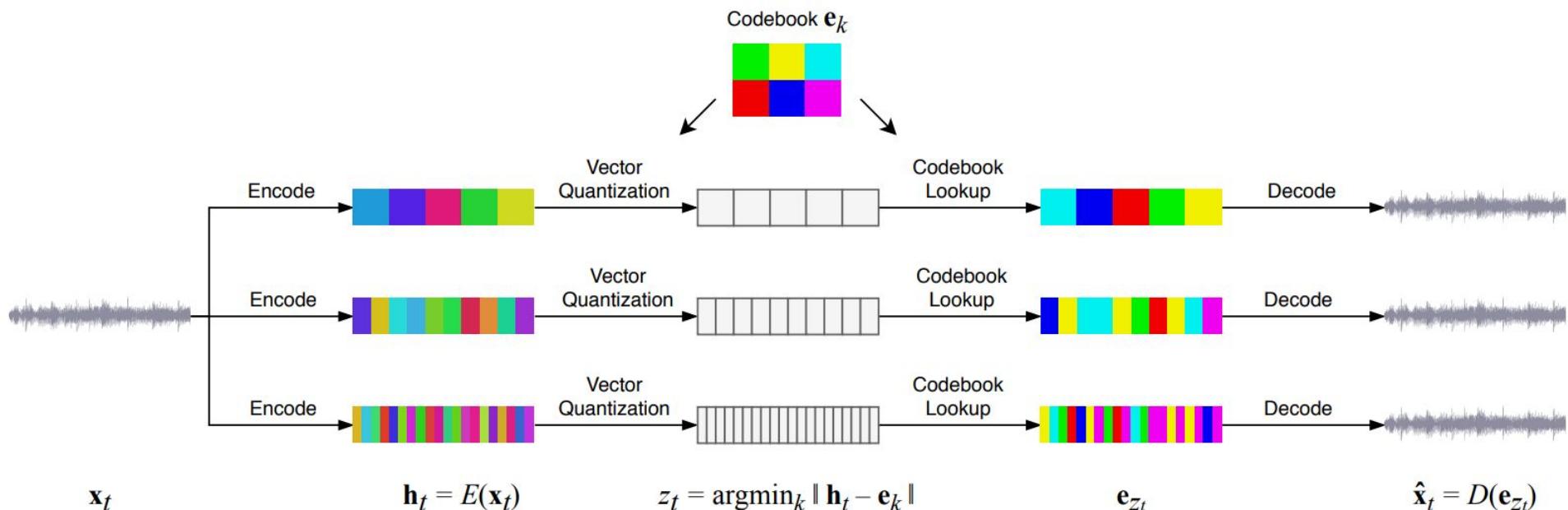


(c) Sparse Transformer (fixed)

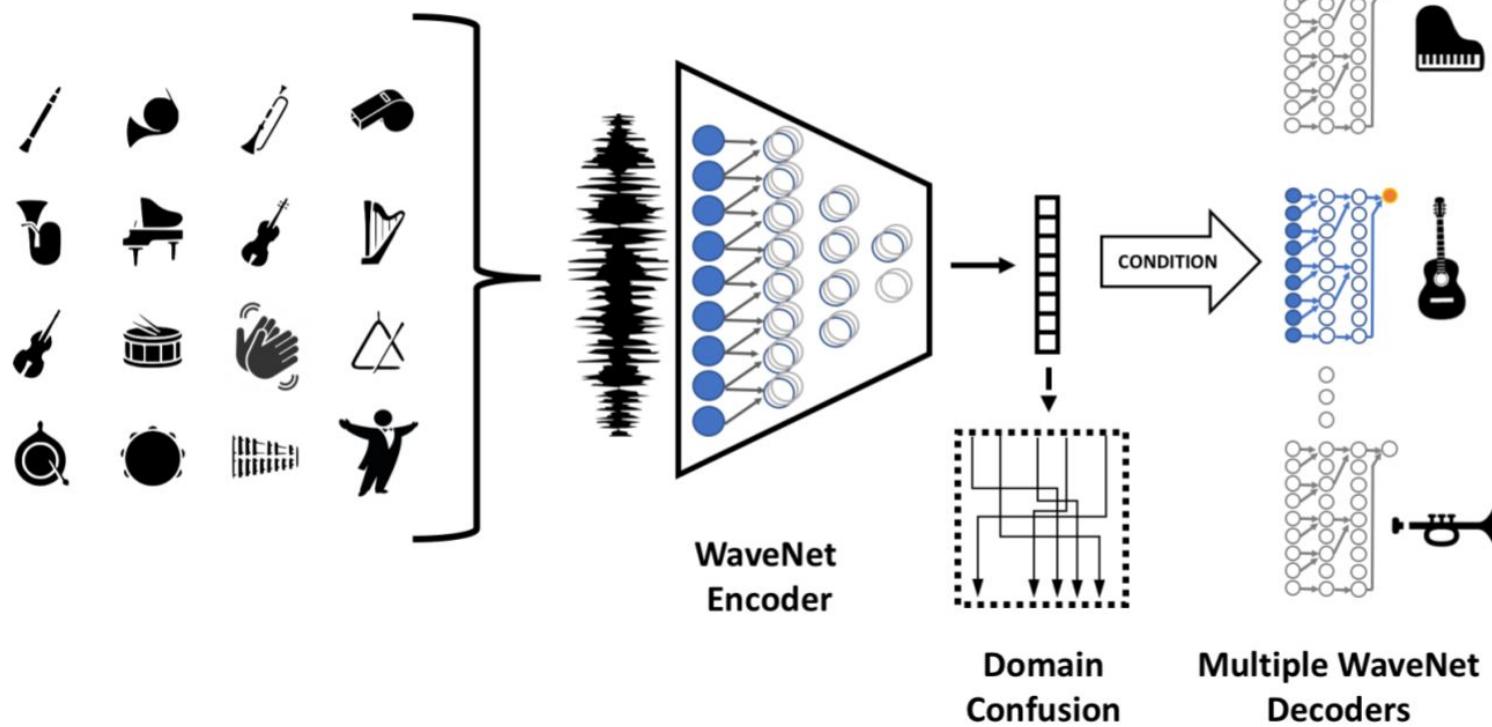


Jukebox

VQ-VAE + Transformer at scale



Universal music translation network



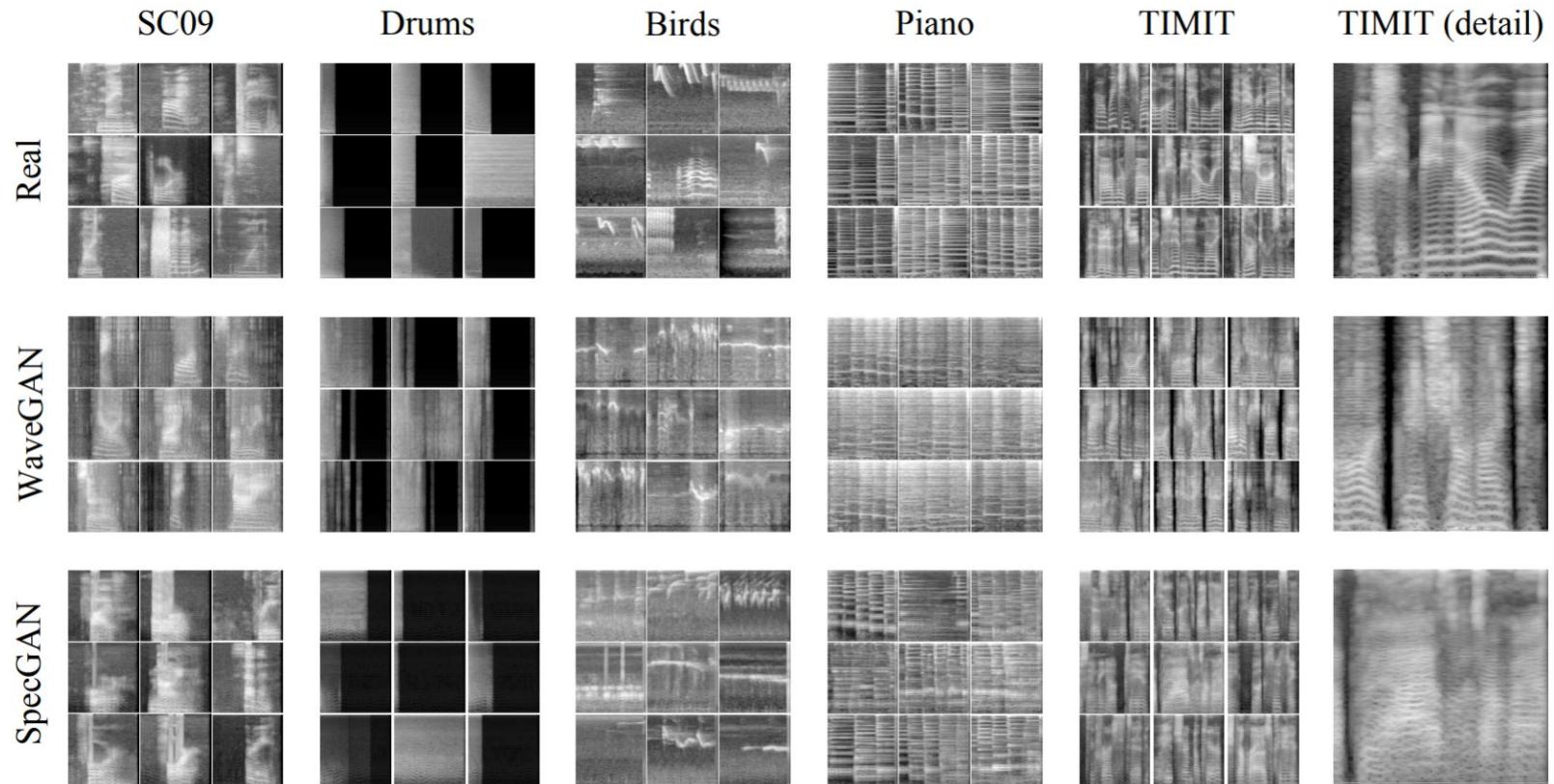
Neural networks

generating death metal
via livestream 24/7 to infinity

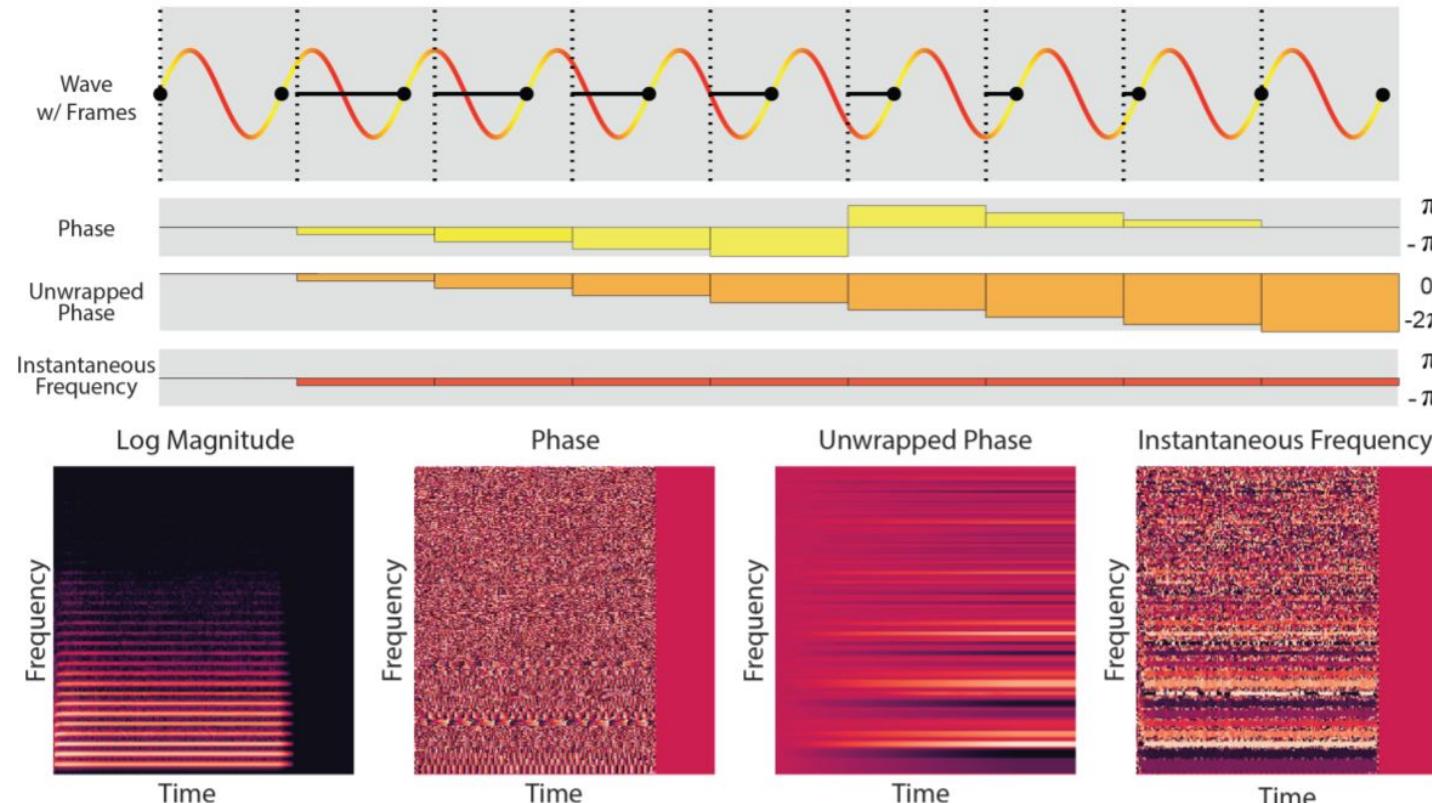
We make raw audio neural networks
that can imitate bands

Adversarial models

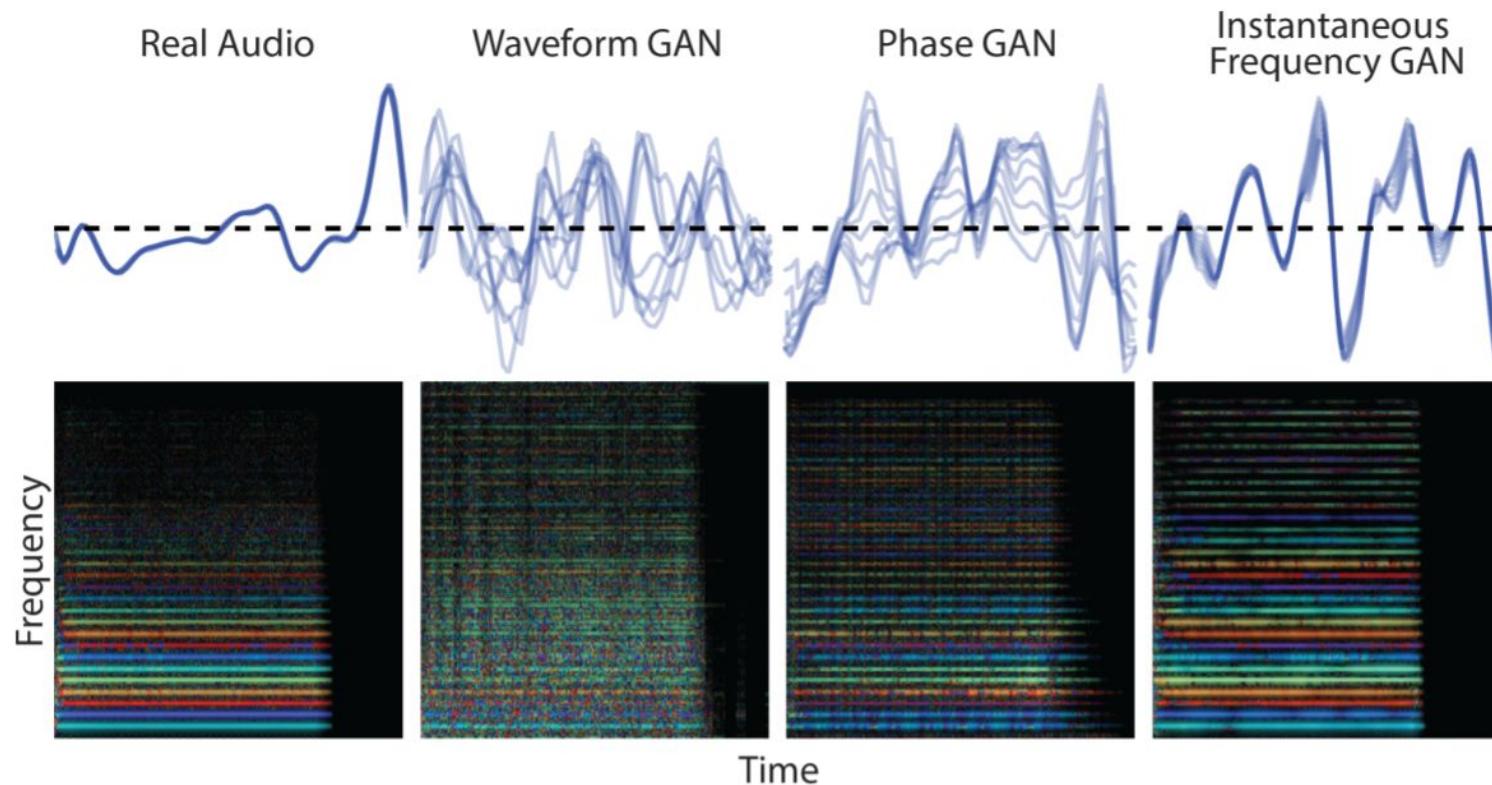
WaveGAN



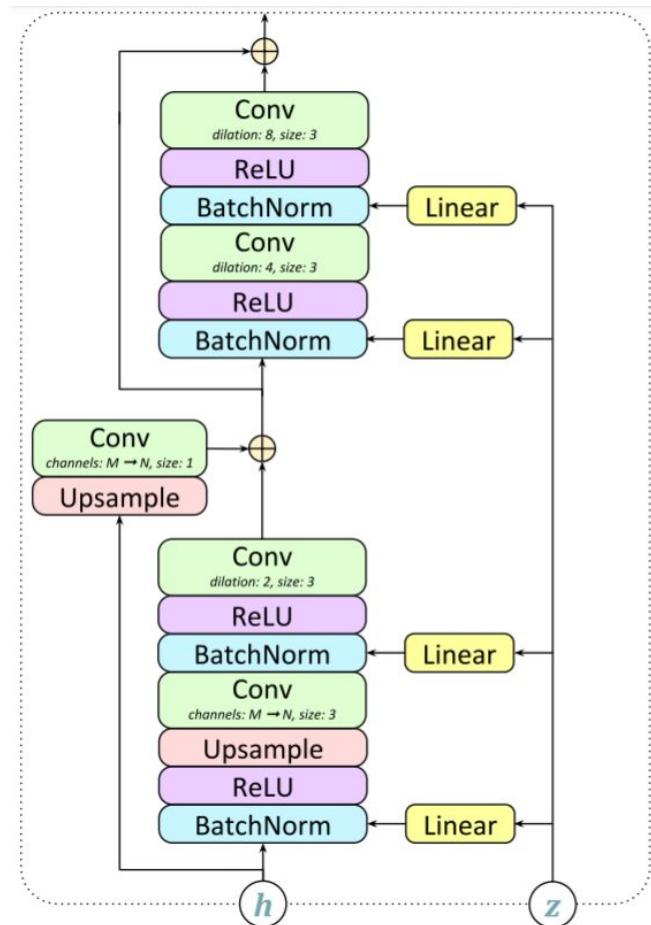
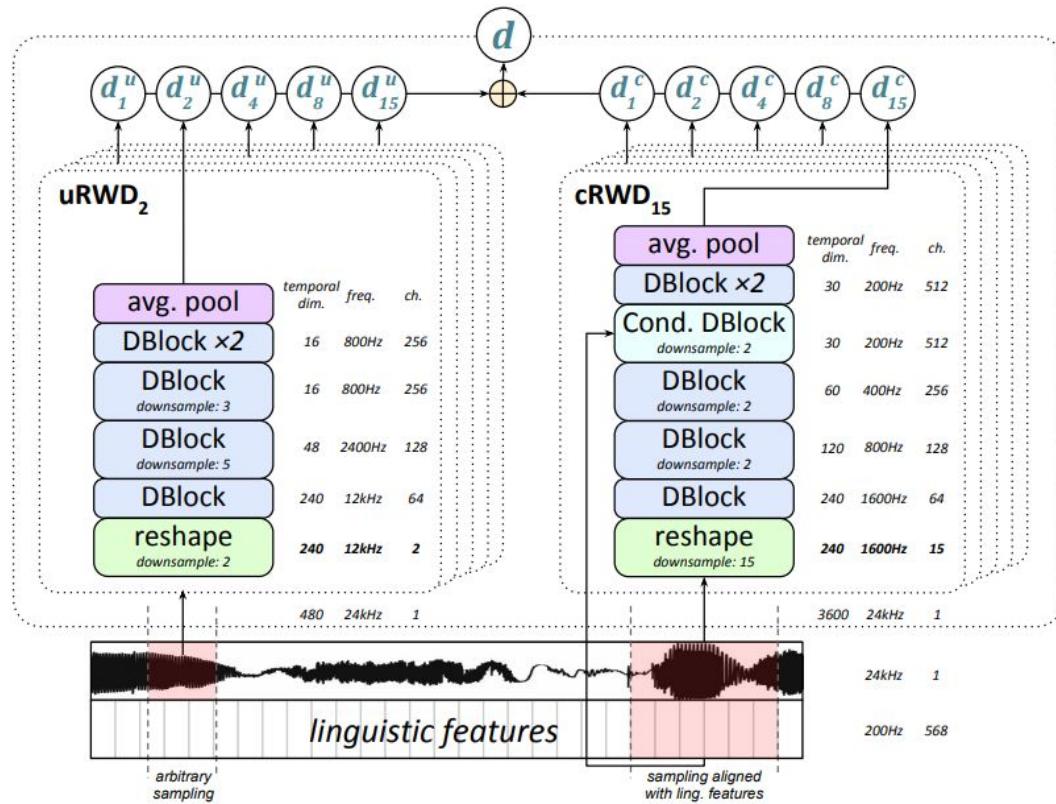
GANSynth



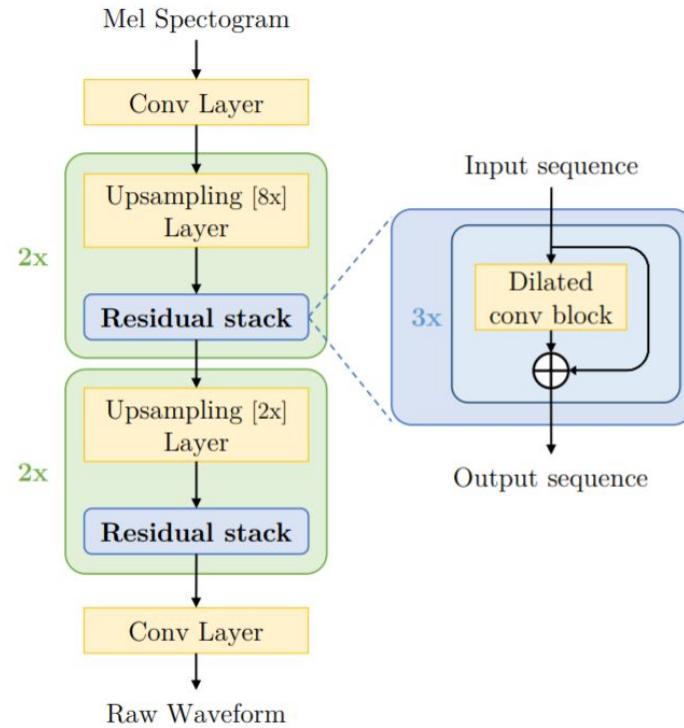
GANSynth



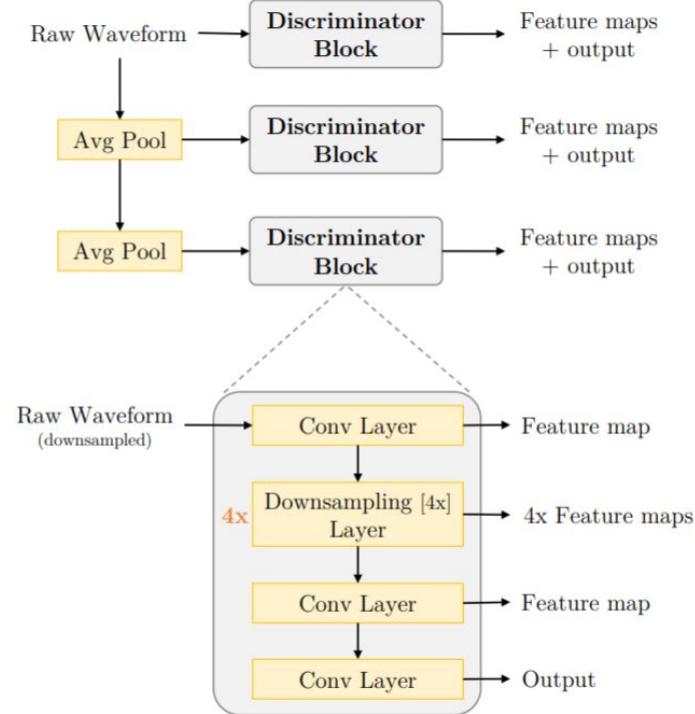
GAN-TTS



MelGAN

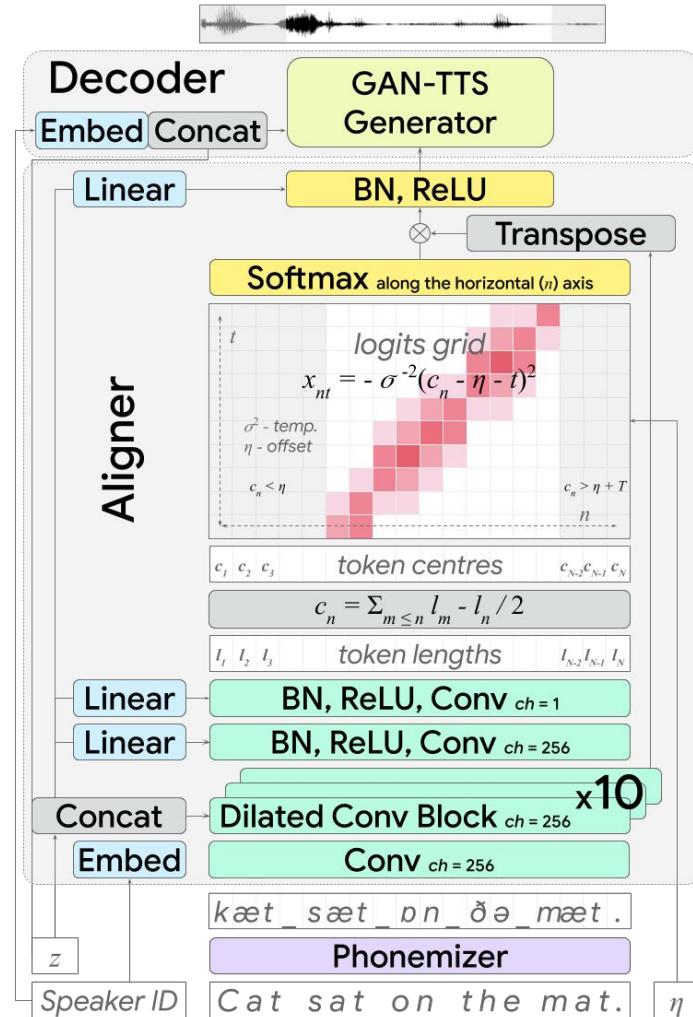
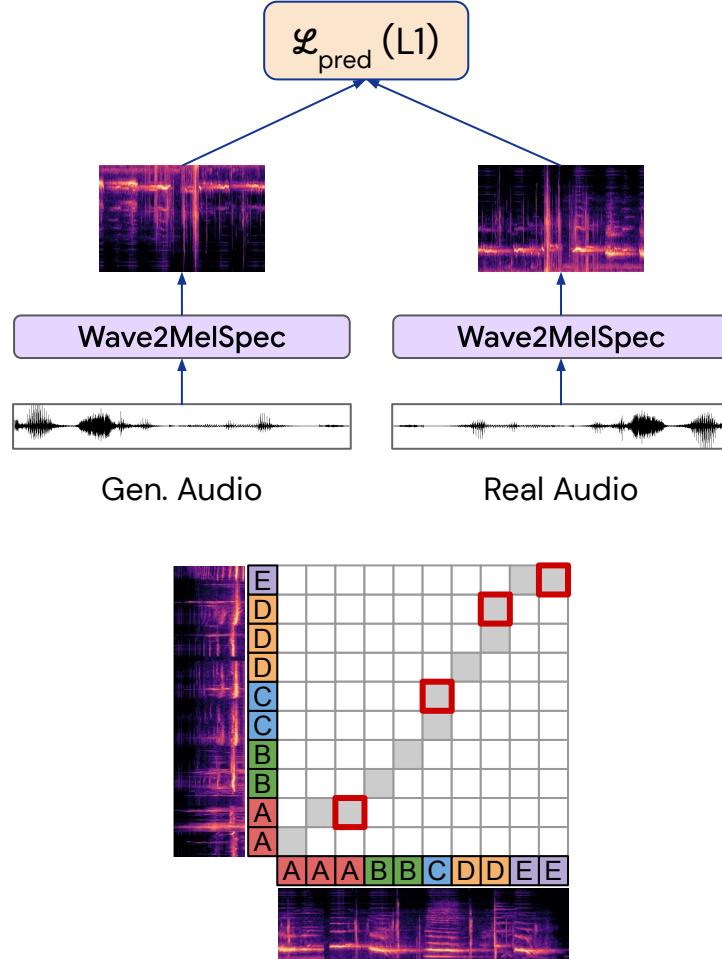


(a) Generator



(b) Discriminator

EATS



Why the emphasis on likelihood in music modelling?

Most popular generative modelling paradigm:

GANs

Why the emphasis on likelihood in music modelling?

Most popular generative modelling paradigm:

GANs

Most popular generative modelling paradigm for music:

likelihood-based (autoregressive)

Why the emphasis on likelihood in music modelling?

- We are still figuring out the right architectural priors for audio discriminators
 - For images, a stack of convolutions is all you need
 - What do we need for audio? Multiresolution? Dilation?
Something phase shift invariant?

Why the emphasis on likelihood in music modelling?

- We are still figuring out the right architectural priors for audio discriminators
 - For images, a stack of convolutions is all you need
 - What do we need for audio? Multiresolution? Dilation?
Something phase shift invariant?
- The sparsely-conditioned setting is dominant
 - We care about “creativity” and capturing diversity
 - GANs are worse at this than likelihood-based models

Alternatives to modelling raw audio directly

- Model complex-valued spectrograms and “deal” with phase (GANSynth)
- Model magnitude spectrograms
Use a vocoder or Griffin-Lim to invert
(Tacotron 1 & 2, MelGAN, MelNet, ...)

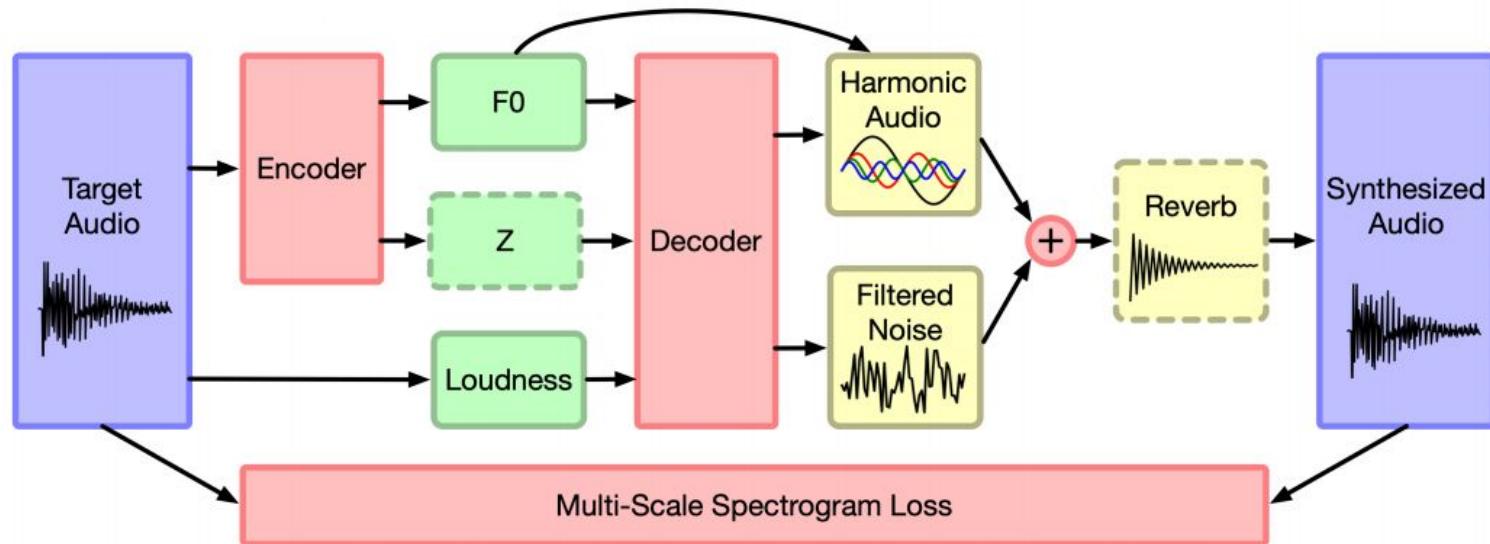
Wang et al., 2017. “Tacotron: Towards end-to-end speech synthesis”, ISCA.

Shen et al., 2018. “Natural TTS Synthesis by Conditioning WaveNet on Mel Spectrogram Predictions”, ICASSP.

Vasquez & Lewis, 2019. “MelNet: A Generative Model for Audio in the Frequency Domain”, arXiv.

Alternatives to modelling raw audio directly

- Differentiable Digital Signal Processing
Use raw audio input, but put DSP components in the model



Summary

- Generative modelling of raw audio is feasible, even in the sparsely conditioned setting
- Likelihood-based models dominate, but GANs are making in-roads in the densely conditioned setting
- Modelling large-scale structure from raw audio is an unsolved problem

Thank you

<https://benanne.github.io/2020/03/24/audio-generation.html>

sanderdieleman@gmail.com



@sediem