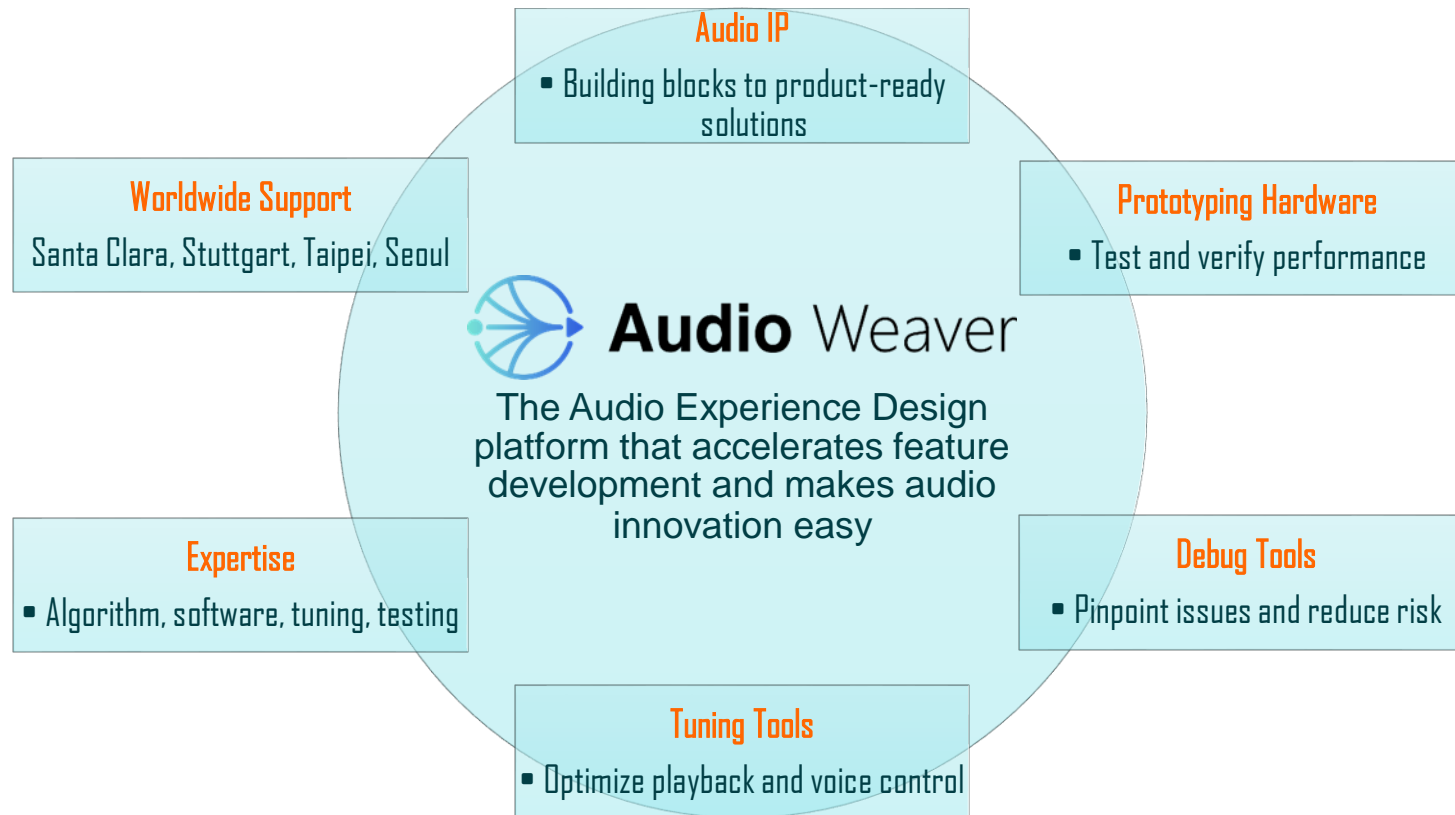# Comparing ML-Based Audio with ML-Based Vision: An Introduction to ML Audio for ML Vision Engineers

Josh Morris

Engineering Manager, Machine Learning

DSP Concepts

# The Audio of Things Approach

DSP Concepts helps product makers deliver remarkable Audio Experience through a flexible and modular approach within a design platform environment. This system makes the entire workflow faster and easier across prototyping, design, debugging, tuning, production, and even over-the-air updates.

# Motivation for Talk

Processing at the edge is getting more and more powerful making it possible to do things that were reserved for the cloud

Audio is becoming increasingly popular

- Standalone applications
  - Smart assistant
  - Voice control
  - Denoising
- Multi-modal applications
  - Industrial sensing
  - Anomaly detection

# Agenda

- **Feature engineering**
  - How is audio different from vision?
  - How is audio like vision?
- **Implementation**
  - Similarities
  - Differences

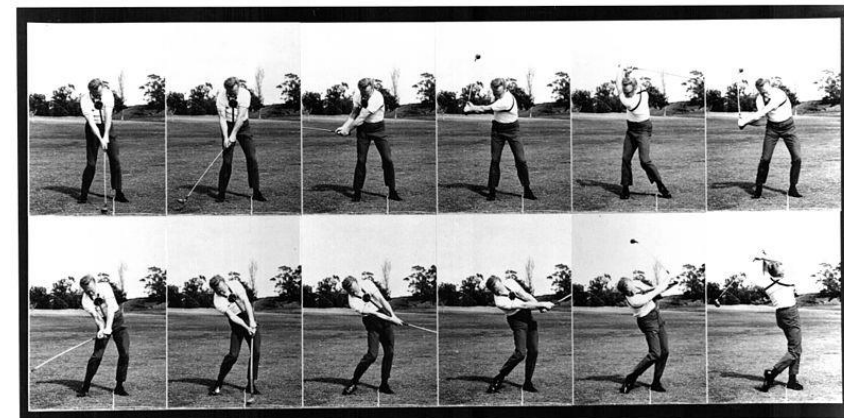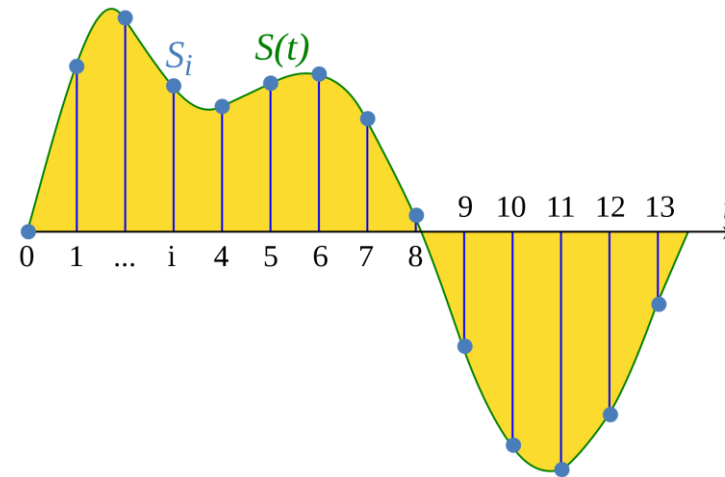- **Common problems in vision and their audio analogues**
  - Classification
  - Sequence decoding
  - Restoration/denoising

# Features

- **Audio is a sequence of samples**
  - Somewhere between video/images
  - Inherent left to right structure to data
  - Sample rate
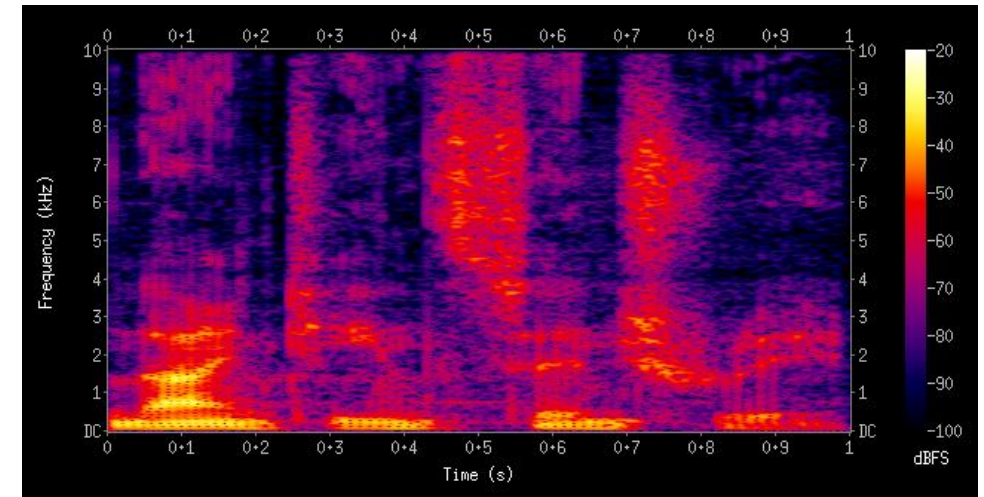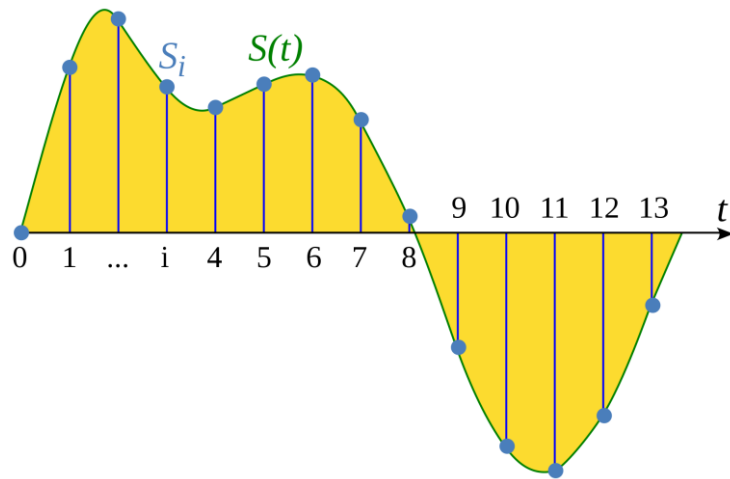  - Bit depth



$S_i$ $S(t)$

9  10  11  12  13  $t$

0  1  ...  i  4  5  6  7  8



https://en.wikipedia.org/wiki/Sampling_(signal_processing)
https://en.wikipedia.org/wiki/File:Mike_Austin_Sequence.JPG

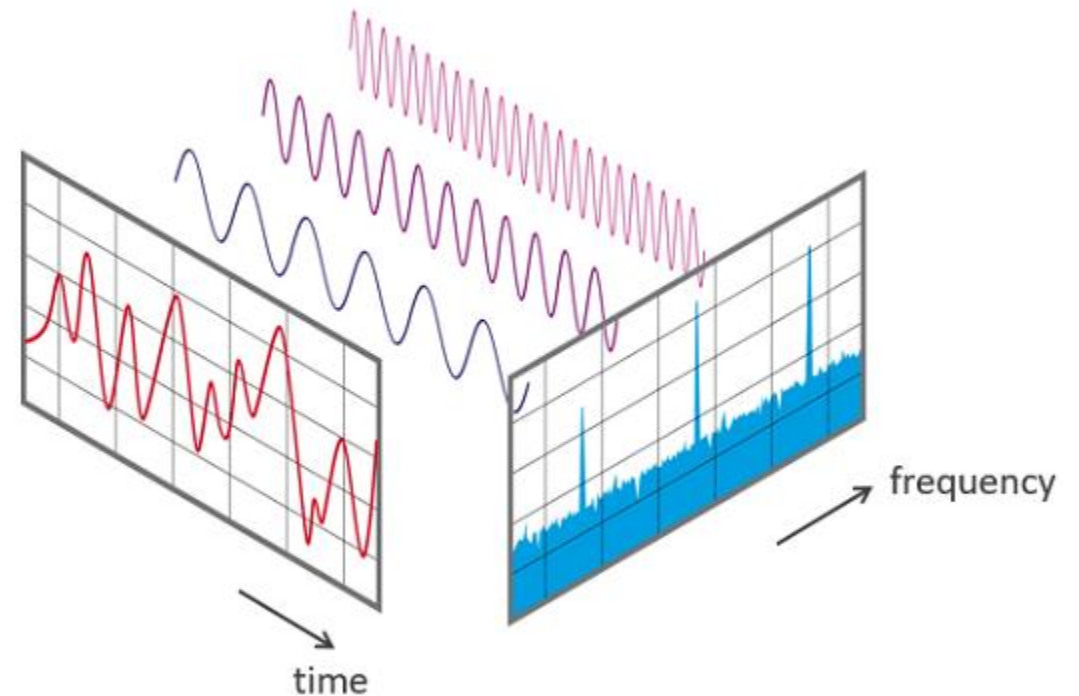# Manipulating Audio from Time Domain to Frequency Domain

A lot of techniques employed for ML audio-based solutions borrow techniques from vision. This means we need to take a 1D sequence of samples and make them look like an image.
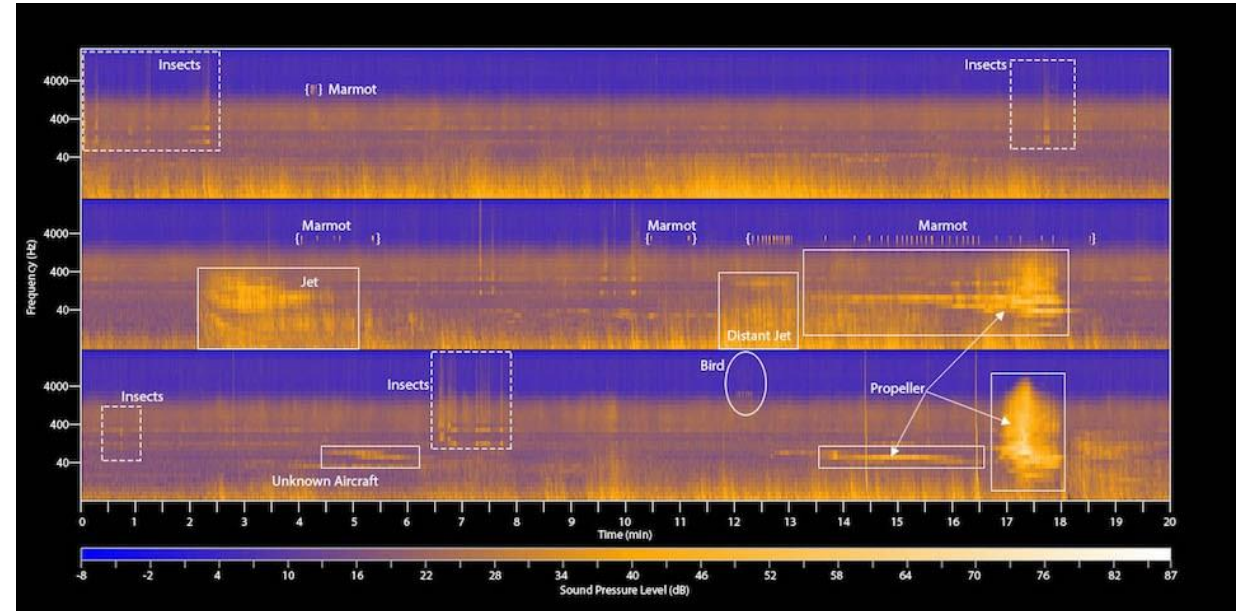
# Fast Fourier Transform

- **Fast Fourier transform (FFT)**

  - Shows frequency over time

  - Linearly spaced frequency bins

  - Can apply processing in the frequency domain and then use an inverse fast Fourier transform (IFFT) to get time domain audio

https://commons.wikimedia.org/wiki/File:FFT-Time-Frequency-View.png
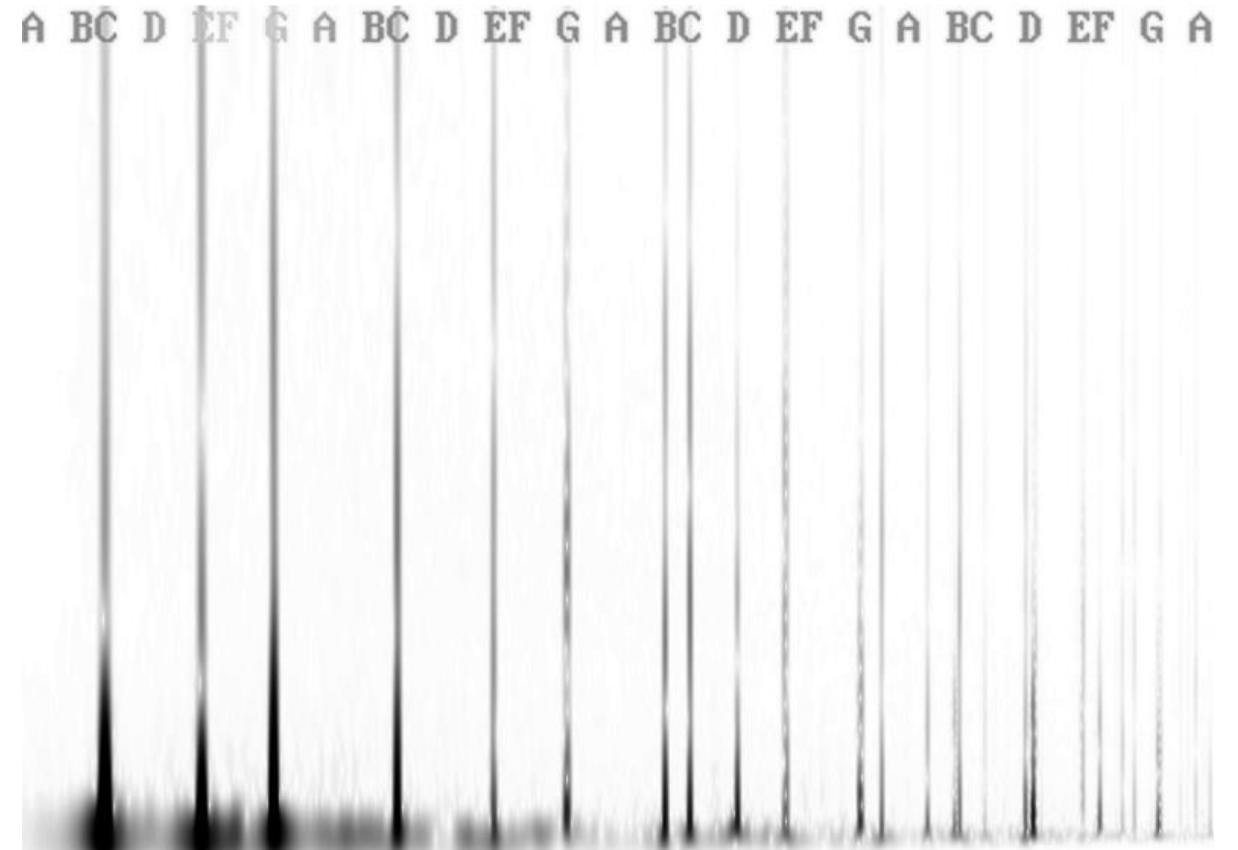
# Power Spectrogram

- **Built from short time Fourier transform**

  - Repeat Fourier transform with set window and hop size

    - short-time Fourier transform (STFT)

  - Take magnitude squared of frequency bins

  - Common for classification tasks



https://upload.wikimedia.org/wikipedia/commons/9/99/Mount_Rainier_soundscape.jpg
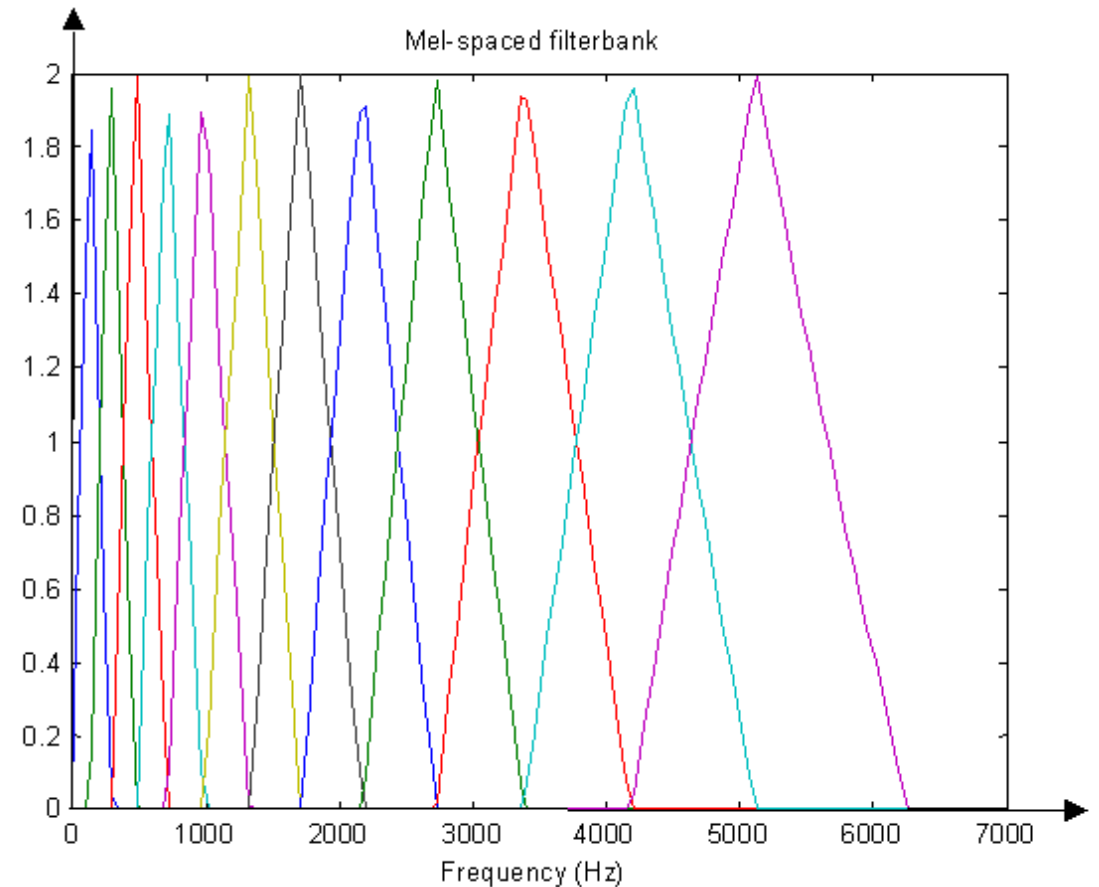
# Constant Q Transform

- **Constant Q transform (CQT)**
  - Logarithmically spaced frequency bins
  - Popular for musical applications
  - More computationally efficient since fewer bins are needed to cover a frequency range



https://en.wikipedia.org/wiki/Constant-Q_transform#/media/File:CQT-piano-chord.png

- **Mel spectrogram**

  - Triangular frequency windows

  - Filter banks that attempt to approximate human hearing

    - Humans struggle to hear frequencies that are close together. This anomaly is known as masking.
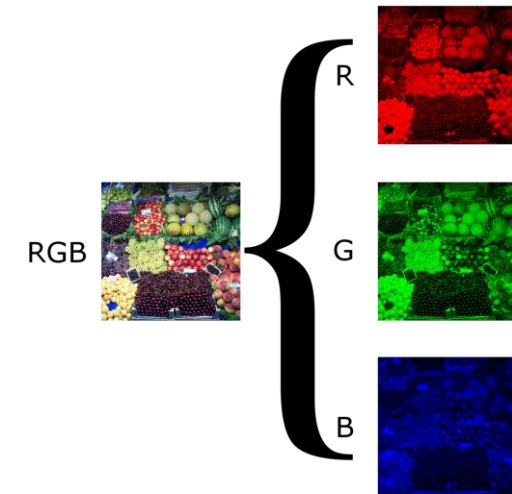


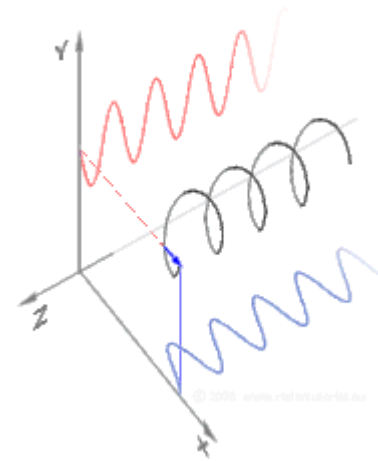http://www.ifp.illinois.edu/~minhdo/teaching/speaker_recognition/speaker_recognition.html

- **Feature normalization**

- **Multi-channel features for complex audio**

  - Color channels in audio

  - Real and imaginary components in frequency time

# Implementation

- **Take audio and create an image-like input with fixed dimensions.**

  - Take the input signal into the frequency domain

  - Take a window of feature vectors

  - Slide window with a hop, usually smaller than the input dimension of the model



https://www.jonnor.com/2021/12/audio-classification-with-machine-learning-europython-2019/

- **Use the same layers**
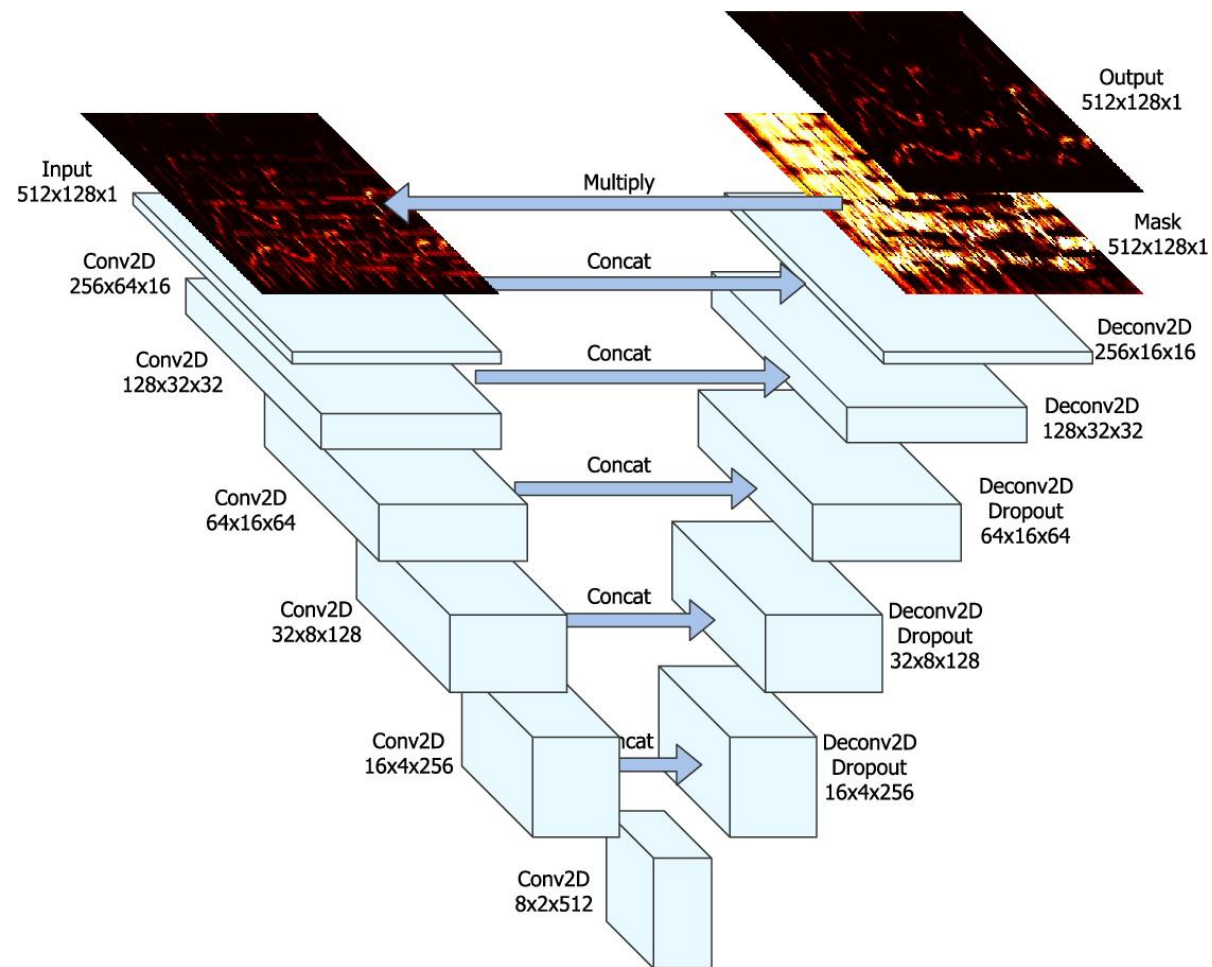  - Convolutional layers
    - Conventional
    - Depth-wise separable
  - Residual blocks
  - RNN

- **Use the same training tricks**



Input 512x128x1
Conv2D 256x64x16
Conv2D 128x32x32
Conv2D 64x16x64
Conv2D 32x8x128
Conv2D 16x4x256
Conv2D 8x2x512
Multiply
Concat
Concat
Concat
Concat
ncat
Output 512x128x1
Mask 512x128x1
Deconv2D 256x16x16
Deconv2D 128x32x32
Deconv2D Dropout 64x16x64
Deconv2D Dropout 32x8x128
Deconv2D Dropout 16x4x256

https://www.semanticscholar.org/paper/Singing-Voice-Separation-with-Deep-U-Net-Networks-Jansson-Humphrey/83ea11b45cba0fc7ee5d60f608edae9c1443861d

# How Are Implementations of Audio Solutions Like Vision Solutions?

- **Some popular vision model architectures show up**
  - EfficientNet
- **Similar concepts**
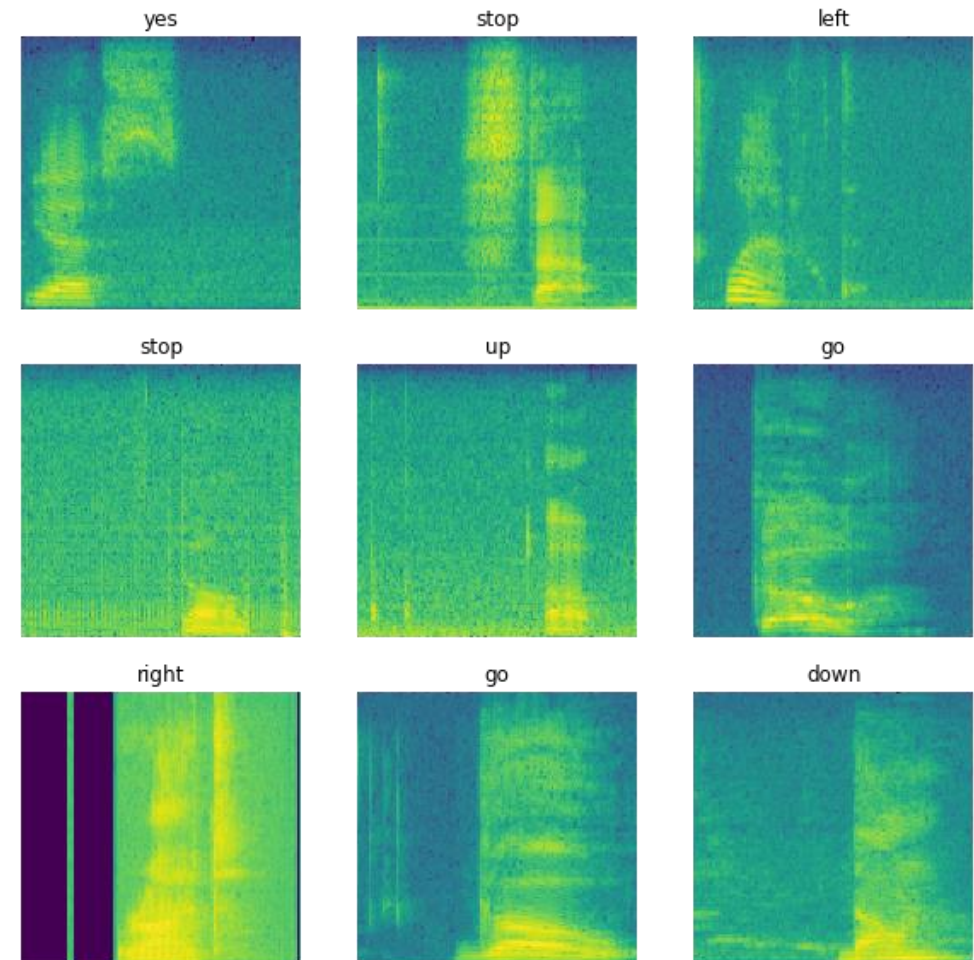  - Encoder-decoder network
  - Transfer learning



https://ai.googleblog.com/2021/08/soundstream-end-to-end-neural-audio.html

# Classification

- **Very similar**

- **Convolution layers pull out structural information**

- **Trained on frequency domain features**

- **In audio, we usually have a sliding window**
  - Like working with a camera stream

- **Can be important for energy savings in complex systems**
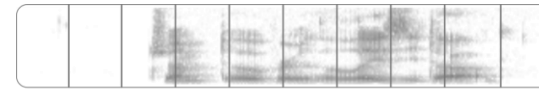  - Motion detection
  - Voice activity detection



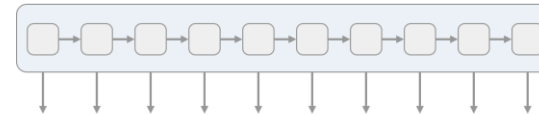https://www.tensorflow.org/tutorials/audio/simple_audo
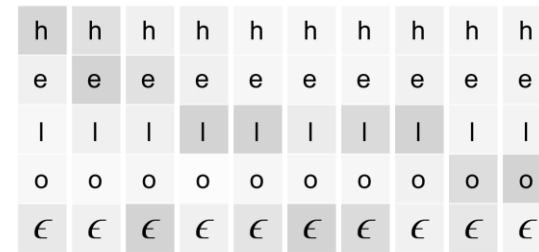
# Sequence Decoding

- **Vision**
  - Optical character recognition
- **Audio**
  - Automatic speech recognition
- **Both look for structural information in their inputs and decode them to a character sequence**
- **Similar architectures**
  - RCNN
  - Transformer

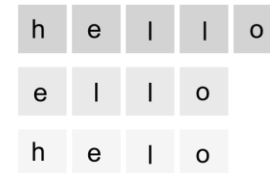We start with an input sequence, like a spectrogram of audio.

The input is fed into an RNN, for example.

The network gives $p_t(a \mid X)$, a distribution over the outputs $\{h, e, l, o, \epsilon\}$ for each input step.

With the per time-step output distribution, we compute the probability of different sequences

By marginalizing over alignments, we get a distribution over outputs.

https://distill.pub/2017/ctc/

- **Vision**

  - Directly regress the image

- **Audio**

  - Regress a gain mask which is applied to audio stream

  - Applied in a streaming fashion

    - Window and hop



https://www.mathworks.com/help/audio/ug/denoise-speech-using-deep-learning-networks.html

# Conclusion

- **Feature engineering**
  - After some preprocessing things are more similar than not

- **Implementation**
  - Windowing with a set stride and hop allows us to deal with streams of data

- **Common problems in vision and their audio analogues**
  - Classification
  - Sequence decoding
  - Restoration/denoising

## Getting Started with Audio

Audio Classification using Transfer Learning

https://www.tensorflow.org/tutorials/audio/transfer_learning_audio

Speech Command Recognition

https://www.tensorflow.org/tutorials/audio/simple_audio

Get a 30 Day Trial of Audio Weaver

https://w.dspconcepts.com/audio-weaver

# The Audio Weaver Framework: Overview

Audio Weaver **accelerates** audio feature development and **enables** collaboration across product teams. With over 550 optimized processing modules, audio designs can be developed and implemented on hardware **without writing any DSP code**.

**AWE Designer**
Windows-based graphical design environment
- ✓ Standard Edition: Design GUI
- ✓ Pro Edition: Works with MathWorks® MATLAB® platform

**AWE Core**
The embedded processing engine
- ✓ Optimized target-specific libraries
- ✓ Available for multiple processors
- ✓ Supports multicore and multi-instance implementation

**Audio IP Modules**
Building blocks for product developers
- ✓ From low-level primitives to complete designs
- ✓ From DSP Concepts and our third-party partners