

# Neural audio synthesis

WaveNet, SampleRNN, and the waveform domain

Sevag Hanssian

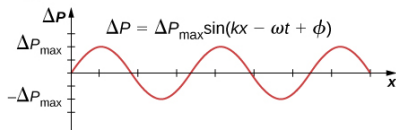
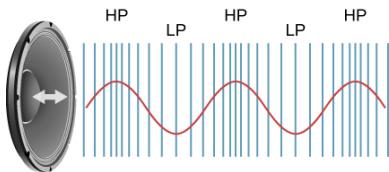
MUMT 618, McGill University

October 14, 2020

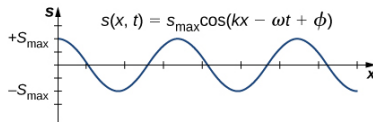
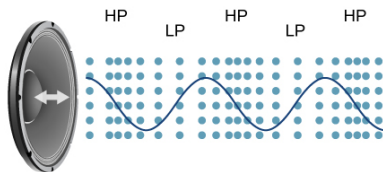
SOUND CHECK

# Physical sound waves

HP = Compression    LP = Rarefaction



(a)

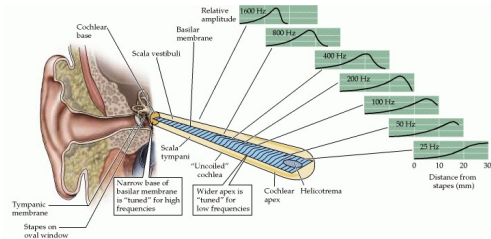


(b)

Figure: a) air pressure vs. distance, b) air molecule displacement<sup>1</sup>

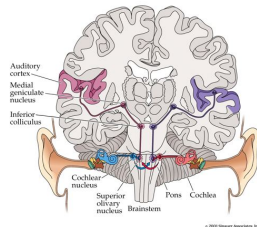
<sup>1</sup>17.1 Sound Waves. URL: <https://courses.lumenlearning.com/suny-osuniversityphysics/chapter/17-1-sound-waves/>.

# Waveforms (human pov)



(a) Tonotopic organization of the inner ear<sup>2</sup>

## Ascending auditory pathways



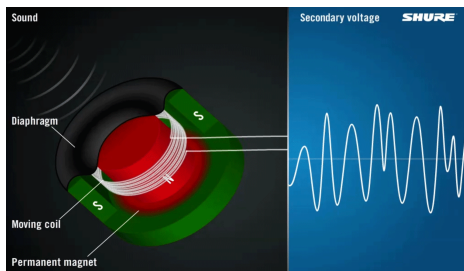
(b) Ascending auditory pathways<sup>3</sup>

Figure: Human auditory neural network

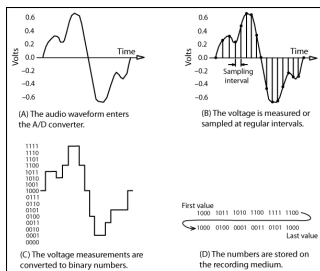
<sup>2</sup>Purves et al. *Neuroscience 2nd ed.* 2001.

<sup>3</sup>Professor David Heeger. *Perception Lecture Notes: Auditory Pathways and Sound Localization.*

# Waveforms (digital pov)



(a) Microphone/transducer<sup>4</sup>



(b) Analog to digital<sup>5</sup>

Figure: Sound pressure wave to analog signal to digital signal

<sup>4</sup>Microphones: Transducer Types (Dynamic, Condenser, Ribbon). URL: <https://www.shure.eu/musicians/discover/educational/transducer-types>.

<sup>5</sup>Bruce Bartlett. Digital Recording Does Not Chop Up Your Music. URL: <https://l2pnet.com/digital-recording-does-not-chop-your-music-11102010/>.

# Waveform domain

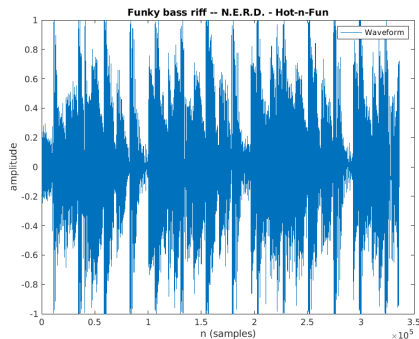


Figure: Bass riff waveform (CLICK2PLAY)<sup>6</sup>

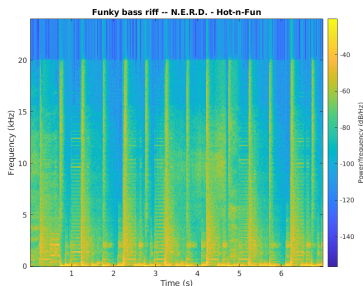
How meaningful is it?

```
>> disp(x(50:128)); 0.0609, 0.0456, 0.0326, 0.0375, ...
```

---

<sup>6</sup>*N.E.R.D. - Hot-n-fun (Official Music Video) ft. Nelly Furtado – YouTube.* URL: <https://www.youtube.com/watch?v=F6ZfA5QZDHY>.

# Frequency domain



**Cons** Often in the spectral domain, phase information is omitted<sup>7</sup>, sacrificing timbre quality in the recreation (good vs. bad)

**Pros** Dimensionality of recognizable audio features (notes, melodies, utterances) is more compact<sup>8</sup>

---

<sup>7</sup>Sander Dieleman. *Generating music in the waveform domain*. URL: <https://benanne.github.io/2020/03/24/audio-generation.html>.

<sup>8</sup>Mehri et al. *SampleRNN: An unconditional end-to-end neural audio generation model*. 2017.

## Symbolic domain

e.g. bass tab<sup>9</sup>

```
G|-----|
D|-----3--3--4-----3--4----|
A|----3--3-----3-----|
E|-----|
```

*Automatic music generation dates back to more than half a century. A prominent approach is to generate music symbolically in the form of a piano roll, which specifies the timing, pitch, velocity, and instrument of each note to be played. This has led to impressive results like producing Bach chorals, polyphonic music with multiple instruments, as well as minute long musical pieces. But symbolic generators have limitations – they cannot capture human voices or many of the more subtle timbres, dynamics, and expressivity that are essential to music.*<sup>10</sup>

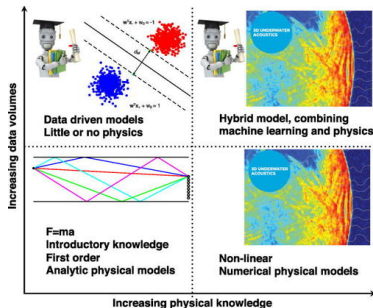
---

<sup>9</sup>HOT N FUN BASS (ver 2). URL:

<https://tabs.ultimate-guitar.com/tab/n-e-r-d-/hot-n-fun-bass-1003978>.

<sup>10</sup>Jukebox. URL: <https://openai.com/blog/jukebox/>.

# Interlude: machine learning in acoustics



**Figure:** Analytic physical models (lower left) give basic insights about physical systems. More sophisticated models, reliant on computational methods (lower right), can model more complex phenomena. Whereas physical models are reliant on rules, which are updated by physical evidence (data), ML is purely data-driven (upper left). By augmenting ML methods with physical models to obtain hybrid models (upper right), a synergy of the strengths of physical intuition and data-driven insights can be obtained<sup>11</sup>

<sup>11</sup>Bianco et al. "Machine learning in acoustics: theory and applications". In: *Acoustical Society of America* (2019).



# Sound synthesis – wavetables and sinusoidal oscillators

Wavetable synthesis is perhaps the oldest technique for creating sounds with computers. It involves **the storage of a single period of a periodic waveform** in a circular buffer. By varying the “speed” with which a read pointer is advanced through the buffer, one can achieve output waveforms of different frequencies<sup>12</sup>.

Sound synthesis algorithms are typically described and implemented using primitive signal processing “building blocks” called **unit generators**. Create complex synthesis by combining multiple oscillators (see next page)<sup>13, 14, 15</sup>.

---

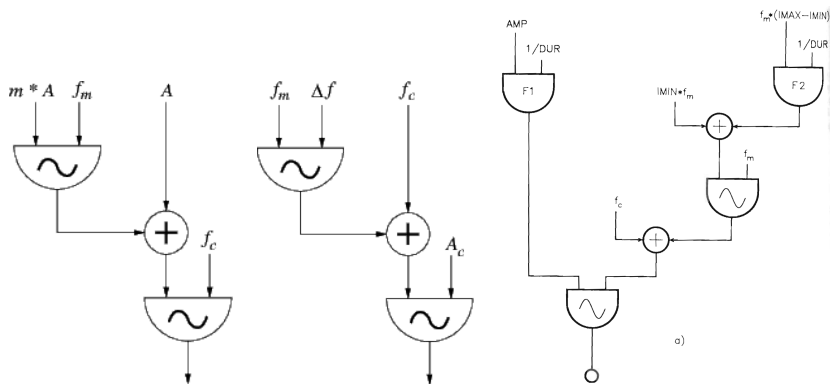
<sup>12</sup>Dr. Gary Scavone. *Wavetable Synthesis*. URL:  
<https://www.music.mcgill.ca/~gary/307/week4/wavetables.html>.

<sup>13</sup>Dr. Gary Scavone. “Classical” amplitude modulation. URL:  
<https://www.music.mcgill.ca/~gary/307/week9/node2.html>.

<sup>14</sup>Dr. Gary Scavone. *Basic frequency modulation*. URL:  
<https://www.music.mcgill.ca/~gary/307/week9/node2.html>.

<sup>15</sup>J. M. Chowning. “The Synthesis of Complex Audio Spectra by Means of Frequency Modulation”. In: *Journal of the Audio Engineering Society* (July 1973).

# Sound synthesis – sinusoidal oscillators



**Figure:** Sinusoidal oscillators in AM synthesis, FM synthesis, Chowning clarinet model (citations on previous slide)

# Neural audio synthesis

*The question of generating musical signals has been extensively studied over the past decades. Most of the previous approaches on this topic were **defined in the spectral domain as a complex set of relatively simple manipulation (subtractive, additive or modulation synthesis)**. However, several recent breakthroughs in audio waveform generative models based on neural networks (Oord, 2016 **WaveNet**) have obtained astounding results in speech synthesis both qualitatively and quantitatively (Mehri, 2016 **SampleRNN**). These systems rely on learning the structure of audio waveforms directly from a given set of audio files, without computing a spectral transform and in an unsupervised manner.*<sup>16</sup>

---

<sup>16</sup>Crestel et al. "Generating Orchestral Music by Conditioning SampleRNN". In: *Timbre 2018: Timbre is a Many-Splendored Thing* (July 2018).

# WaveNet

Concatenative TTS: database of short speech fragments recombined to form complete utterances

Parameteric TTS: information to generate data is stored as model parameters, contents and characteristics of speech are controlled via model inputs. Parametric models generate audio through vocoders – sound less natural than concatenative

WaveNet changes this paradigm by directly modelling the raw waveform of the audio signal, one sample at a time. As well as yielding more natural-sounding speech, using raw waveforms means that WaveNet can model any kind of audio, including music.<sup>17, 18</sup>

---

<sup>17</sup> *WaveNet: A generative model for raw audio* | DeepMind. Sept. 2016. URL: <https://deepmind.com/blog/article/wavenet-generative-model-raw-audio>.

<sup>18</sup> van den Oord et al. *WaveNet: A generative model for raw audio*. 2016.

## (Artificial) Neural Networks

A kind of learning model which automatically learns non-linear functions from input to output

Biologically inspired metaphor:

- Network of computational units called neurons
- Each neuron takes scalar inputs, and produces a scalar output, very much like a logistic regression model

$$\text{Neuron}(\vec{x}) = g(a_1x_1 + a_2x_2 + \dots + a_nx_n + b)$$

As a whole, the network can theoretically compute any computable function, given enough neurons. (These notions can be formalized.)

Figure: Artificial neural network primer<sup>19</sup>

---

<sup>19</sup>Prof. Jackie Cheung. *COMP 550 NLP Lecture slides, Fall 2020, Lecture 10.* 2020.

# WaveNet – details

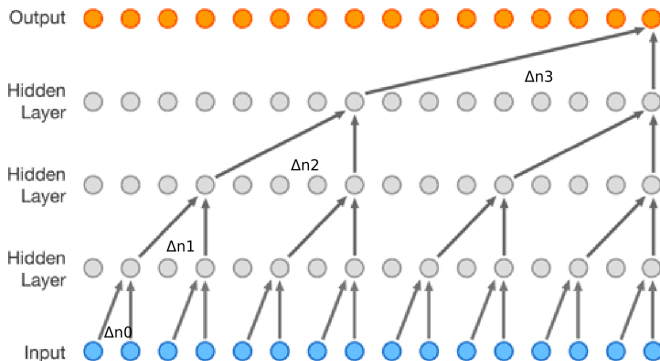
- Causality –  $y[n]$  can only depend on  $\{x[n], x[n-1], \dots\}$
- Autoregressive – audio samples are generated by estimating likely next values based on past values
- Convolutions – this is the same convolution from DSP<sup>20</sup>
- Dilated convolutions – widen the time scale of learning by increasing space between samples
- $\mu$ -law quantization: quantization mapping to critical bands, a common speech technique<sup>21</sup>

---

<sup>20</sup>Mathieu and Henaff. *Fast Training of Convolutional Networks through FFTs*. Mar. 2014.

<sup>21</sup>“Pulse Code Modulation (PCM) of voice frequencies”. In: *ITU-T. Recommendation G. 711*. (1988).

# WaveNet diagram



**Figure:** How wavenet computes samples of output  $y[n]$  from historical samples  $\{x[n], x[n-1], \dots\}$

# Inference in a FF Neural Network

Perform computations forwards through the graph:

$$\begin{aligned}\mathbf{h}^1 &= g^1(\mathbf{x}\mathbf{W}^1 + \mathbf{b}^1) \\ \mathbf{h}^2 &= g^2(\mathbf{h}^1\mathbf{W}^2 + \mathbf{b}^2) \\ \mathbf{y} &= \mathbf{h}^2\mathbf{W}^3\end{aligned}$$

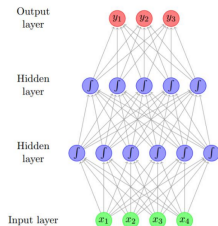


Figure 2: Feed-forward neural network with two hidden layers.

Note that we are now representing each layer as a vector; combining all of the weights in a layer across the units into a weight matrix

**Figure:** Example of how the previous diagram maps to some equations<sup>22</sup>

<sup>22</sup>Prof. Jackie Cheung. *COMP 550 NLP Lecture slides, Fall 2020, Lecture 10.* 2020.



# WaveNet examples

All examples are from WaveNet blog post<sup>23</sup>

- Unconditioned speech (i.e. let the machine do whatever it wants):  
CLICK TO PLAY
- Conditioned speech, English (i.e. train machine to learn a specific phrase\*): CLICK TO PLAY
- Conditioned speech, Mandarin: CLICK TO PLAY
- Unconditional music\*\*: CLICK TO PLAY
- Unconditional music: CLICK TO PLAY

\*: it's unclear to me how – the paper has precious few details

\*\* : they provide no examples of conditioned (i.e. structured) music

---

<sup>23</sup> *WaveNet: A generative model for raw audio* | DeepMind. Sept. 2016. URL: <https://deepmind.com/blog/article/wavenet-generative-model-raw-audio>.

# SampleRNN

Use RNNs to learn long-term patterns:

*raw audio signals are challenging to model because they contain structure at very different scales: correlations exist between neighboring samples as well as between ones thousands of samples apart. SampleRNN helps to address this challenge by using a hierarchy of modules, each operating at a different temporal resolution.*<sup>24</sup>

Contrast with WaveNet dilations:

*In order to deal with long-range temporal dependencies needed for raw audio generation, we develop new architectures based on dilated causal convolutions*<sup>25</sup>

---

<sup>24</sup> [Mehri et al. SampleRNN: An unconditional end-to-end neural audio generation model. 2017.](#)

<sup>25</sup> [WaveNet: A generative model for raw audio | DeepMind. Sept. 2016. URL: <https://deepmind.com/blog/article/wavenet-generative-model-raw-audio>.](#)

# SampleRNN network

*Recurrent neural networks, also known as RNNs, are a class of neural networks that allow previous outputs to be used as inputs while having hidden states.*<sup>26</sup>

RNNs are commonly used to model sequential data which can be formulated as:

$$h_t = \mathcal{H}(h_{t-1}, x_{i=t}) \quad (2)$$

$$p(x_{i+1}|x_1, \dots, x_i) = \text{Softmax}(\text{MLP}(h_t)) \quad (3)$$

with  $\mathcal{H}$  being one of the known memory cells, Gated Recurrent Units (GRUs) (Chung et al., 2014), Long Short Term Memory Units (LSTMs) (Hochreiter & Schmidhuber, 1997), or their deep variations (Section 3). However, raw audio signals are challenging to model because they contain structure at very different scales: correlations exist between neighboring samples as well as between ones thousands of samples apart.

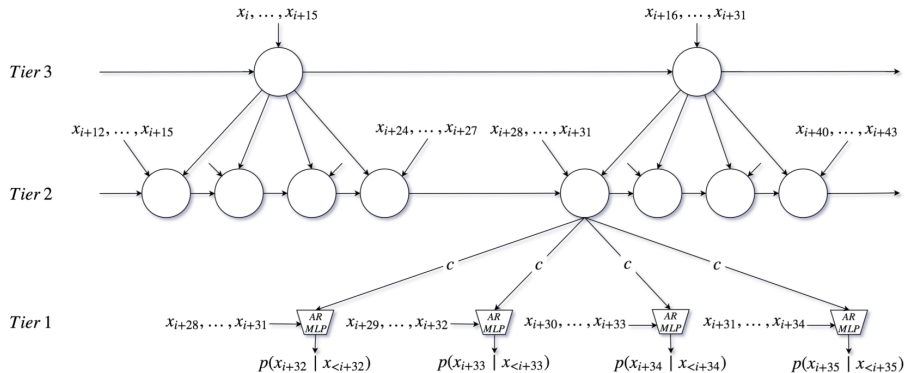
SampleRNN helps to address this challenge by using a hierarchy of modules, each operating at a different temporal resolution. The lowest module processes individual samples, and each higher module operates on an increasingly longer timescale and a lower temporal resolution. Each module conditions the module below it, with the lowest module outputting sample-level predictions. The entire hierarchy is trained jointly end-to-end by backpropagation.

---

<sup>26</sup>Shervine Amidi. CS 230 - Recurrent Neural Networks Cheatsheet. URL: <https://stanford.edu/~shervine/teaching/cs-230/cheatsheet-recurrent-neural-networks>.

## SampleRNN diagram

Note how similar it is to WaveNet. WaveNet “dilates” the convolutions to go from a small time step to a large one – SampleRNN “upsamples” the vector to bump up the temporal resolution



**Figure:** How SampleRNN computes samples of output  $y[n]$  from historical samples  $\{x[n], x[n-1], \dots\}$

# SampleRNN in practise

So what does SampleRNN actually look like? Up-to-date implementation, PRiSM SampleRNN based on Tensorflow 2<sup>27</sup>,<sup>28</sup>. Your best bet at getting up and running (works for me).

- Own an NVIDIA GPU + modern Linux to set up Python dependencies
- Gather training data (wav files), run train.py and generate.py (README has example commands and description of parameters)
- How to choose parameters?

*Training neural networks is, ironically, more of an art than a science, and depends on a lot of trial and error... So I'd say just dive into it with the default settings initially, and then see what results you get.*<sup>29</sup>

---

<sup>27</sup>Dr. Christopher Melen. A Short History of Neural Synthesis. URL: <https://www.rncm.ac.uk/research/research-centres-rncm/prism/prism-blog/a-short-history-of-neural-synthesis/>.

<sup>28</sup>PRiSM SampleRNN – Neural sound synthesis with Tensorflow2. URL: <https://github.com/rncm-prism/prism-samplernn>.

<sup>29</sup>Best parameters for a layperson - Issue #5. URL: <https://github.com/rncm-prism/prism-samplernn/issues/5>.

# SampleRNN examples

Examples are from independent blog post<sup>30\*</sup>

- Unconditioned piano music: [CLICK TO PLAY](#)
- Unconditioned music, Dawn of Midi: [CLICK TO PLAY](#)
- Unconditioned music, Animals as Leaders<sup>\*\*</sup>: [CLICK TO PLAY](#)
- Unconditioned music, Animals as Leaders: [CLICK TO PLAY](#)
- Unconditioned music, Animals as Leaders: [CLICK TO PLAY](#)

\*: Karl Hiner mentions that he never got results as good as WaveNet 's public releases. Purports that they do lots of specialized training

\*\* : from my own experiments of running PRiSM SampleRNN<sup>31</sup>

---

<sup>30</sup> [Karl Hiner](#). *Generating Music with WaveNet and SampleRNN*. URL: [https://karlhiner.com/music\\_generation/wavenet\\_and\\_samplernn/](https://karlhiner.com/music_generation/wavenet_and_samplernn/).

<sup>31</sup> *PRiSM SampleRNN – Neural sound synthesis with Tensorflow2*. URL: <https://github.com/rncm-prism/prism-samplernn>.

# Traits of WaveNet and SampleRNN

Referring to WaveNet generated speech quality:

*Unconditional generation from this model manifests as “babbling” due to the lack of longer term structure<sup>32</sup>*

Referring to WaveNet and SampleRNN listening experiments trained on classical and techno music:

*we first compare the related audio file sets of both experiments. After listening to the music samples, we conclude that none of the four sets contain audio that even slightly resembles music. The files generally sound noisy and random.<sup>33</sup>*

---

<sup>32</sup>Engel et al. *Neural Audio Synthesis of Musical Notes with WaveNet Autoencoders.*

<sup>33</sup>Sinit Tafla. *AN ANALYSES OF THE CURRENT CAPABILITIES OF NEURAL NETWORKS TO PRODUCE MUSIC-RESEMBLING AUDIO.* Dec. 2018.

## Unstructured can be good

The dadabots<sup>34</sup> use the incoherence of SampleRNN to their advantage. Who are they? “We make raw audio neural networks that can imitate bands”. Example: Relentless Doppleganger

*[...] we want the output to overfit short timescale patterns (timbres, instruments, singers, percussion) and underfit long timescale patterns (rhythms, riffs, sections, transitions, compositions) so that it sounds like a recording of the original musicians playing new musical compositions in their style.*<sup>35</sup>

---

<sup>34</sup> dadabots.com. URL: <https://dadabots.com/>.

<sup>35</sup> CJ Carr and Zack Zukowski. “Generating Albums with SampleRNN to Imitate Metal, Rock, and Punk Bands”. In: (Nov. 2018).



## Towards structure

Jukebox<sup>36,37</sup> addresses the long input problem with an autoencoder that compresses raw audio to a lower-dimensional space by discarding perceptually irrelevant information. Then trains a model to generate audio in this compressed space, and upsample back to the raw audio space.

Tacotron 2<sup>38,39</sup> uses an 80-dimensional audio spectrogram with frames computed every 12.5 milliseconds to capture not only pronunciation of words, but also various subtleties of human speech, including volume, speed and intonation. These features are converted to a waveform using a WaveNet-like architecture.

---

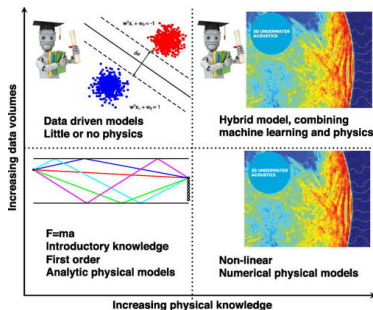
<sup>36</sup>Jukebox. URL: <https://openai.com/blog/jukebox/>.

<sup>37</sup>Dhariwal et al. Jukebox: A Generative Model for Music. Apr. 2020.

<sup>38</sup>Tacotron 2: Generating Human-Like Speech from Text. URL: <https://ai.googleblog.com/2017/12/tacotron-2-generating-human-like-speech.html>.

<sup>39</sup>Shen et al. Natural TTS Synthesis by Conditioning WaveNet on Mel Spectrogram Predictions. Feb. 2018.

# Recall



**Figure:** By augmenting ML methods with physical models to obtain hybrid models (upper right), a synergy of the strengths of physical intuition and data-driven insights can be obtained<sup>40</sup>

<sup>40</sup>Bianco et al. "Machine learning in acoustics: theory and applications". In: *Acoustical Society of America* (2019).

## Middle ground

**Traditional synthesis:** structured building blocks based on frequency and/or symbolic representation. Fiddly parameters, imperfect recreations

**Fully neural audio synthesis:** realistic timbre, dynamics, delays, “human” traits, create natural sounds. Unstructured black box

**Differentiable DSP**<sup>41, 42</sup>

*a collection of linear filters and sinusoidal oscillators can create the sound of a realistic violin if the frequencies and responses are tuned in just the right way. However, it is difficult to dynamically control all of these parameters by hand, which is why synthesizers with simple controls often sound unnatural and “synthetic”.*

*With DDSP, we use a neural network to convert a user’s input into complex DSP controls that can produce more realistic signals.*

---

<sup>41</sup>DDSP: Differentiable Digital Signal Processing. URL: <https://magenta.tensorflow.org/ddsp>.

<sup>42</sup>Engel et al. DDSP: Differentiable Digital Signal Processing. Aug. 2020.

# Differentiability

*Derivatives, mostly in the form of gradients and Hessians, are ubiquitous in machine learning. Automatic differentiation (AD), also called algorithmic differentiation or simply “autodiff”, is a family of techniques similar to but more general than backpropagation for efficiently and accurately evaluating derivatives of numeric functions expressed as computer programs.*<sup>43</sup>

The gradient, or slope, of the loss function. How does machine learning know that one set of random parameters is better than another? How does it “learn” to improve? By moving in a direction that reduces the slope of error, i.e. gradient descent.

---

<sup>43</sup>Baydin et al. “Automatic Differentiation in Machine Learning”. In: *JMLR* (2018).

# Conclusions

- Modeling audio in the waveform domain is now feasible with the computational power of modern machines and neural networks
- WaveNet and SampleRNN were the trendsetters, and there is a lot of derivative work since then (WaveRNN, Jukebox, etc.)
- Advantages include learning realistic timbre, dynamics, intonations, etc. directly from the waveform. No reconstruction problems
- Disadvantages include unstructured/unreliable outputs, black box (difficult to understand) computational model