

11-755— Spring 2021

Large Scale Multimedia Processing



Lecture 1/6

Multimedia capture and storage

Rita Singh

Carnegie Mellon University

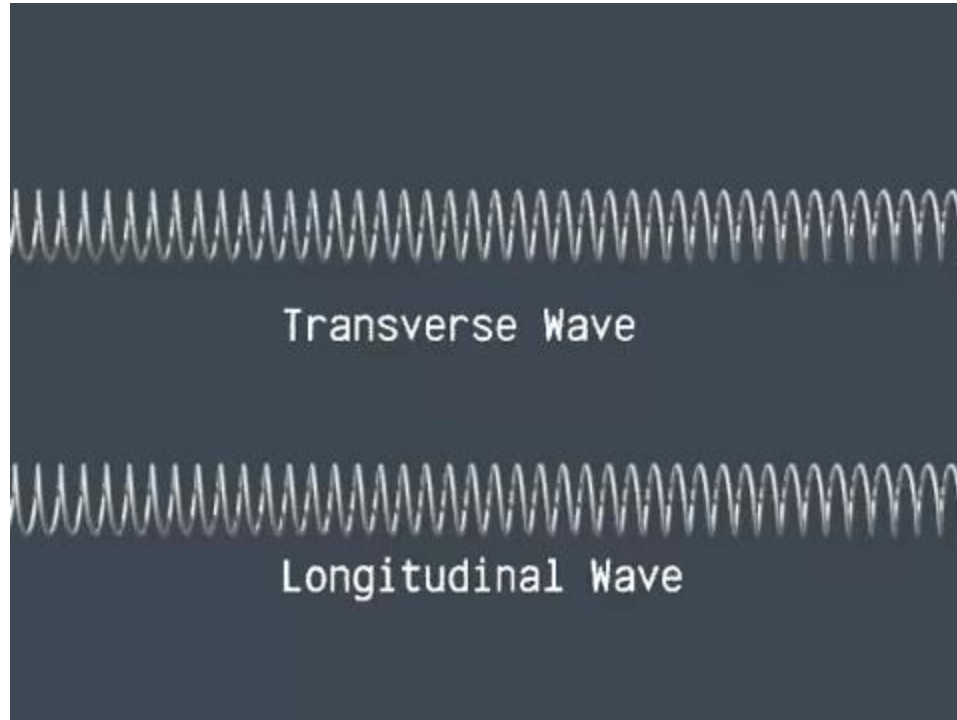
In this lecture

- Digital multimedia: Capture and storage
 - Audio
 - Images
 - Video
 - Text
- Digital multimedia: Processing
 - Audio processing
 - Some generic ML based processing techniques

Sound waves

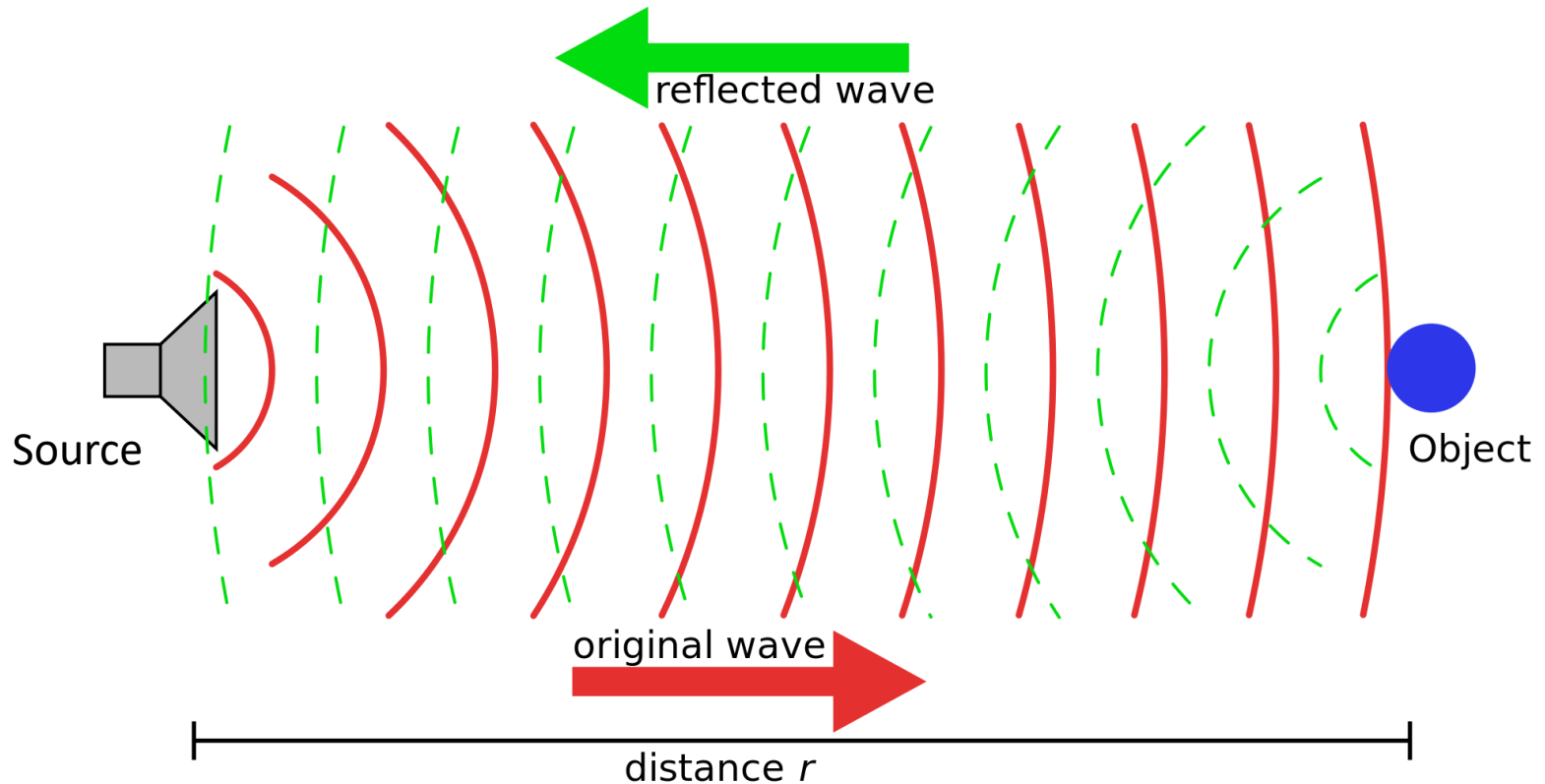
Pattern of displacement
propagates

Particles do not “travel”



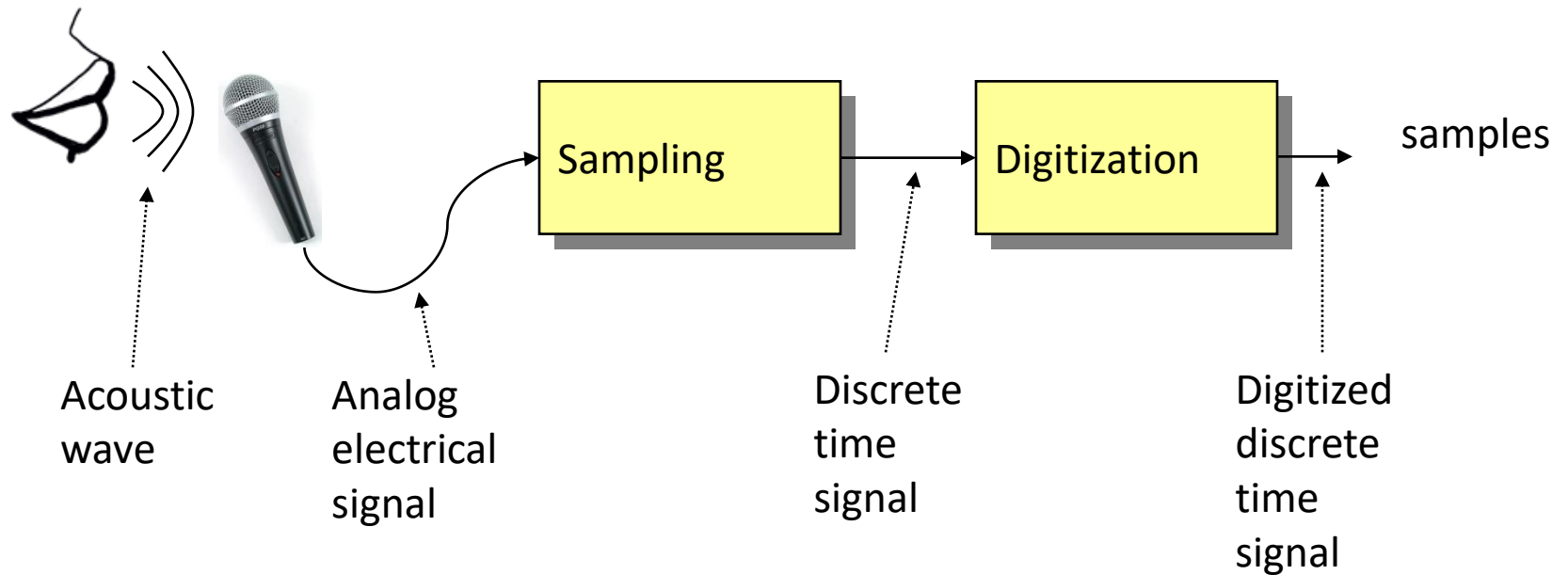
Sound is a longitudinal wave

What is sound?



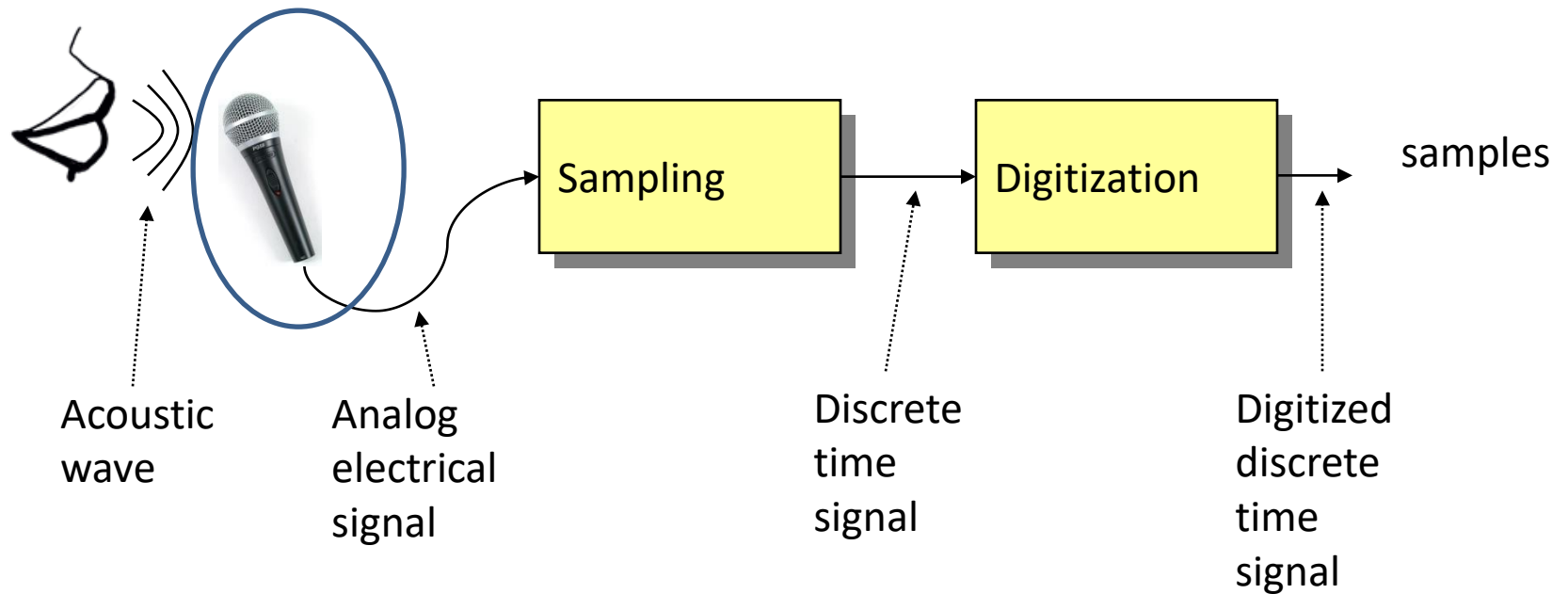
- All sounds are actually pressure waves
 - Sound exerts physical force

Audio Signal Capture



- The goal of signal capture is to convert pressure waves to a series of numbers that the computer can process
- Three stages
 - Transduction by a physical device
 - Sampling
 - Digitization

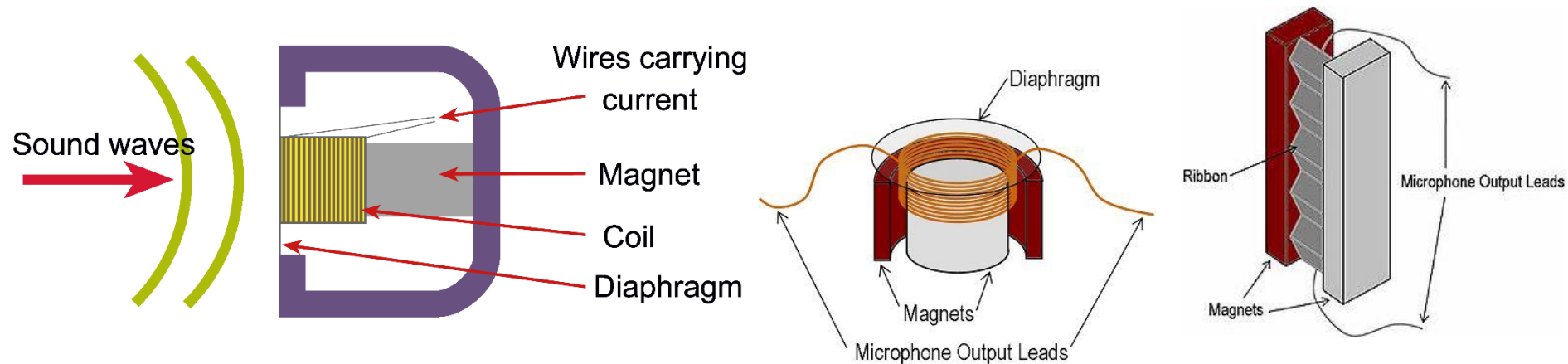
First Stage: Microphones



Microphone:

- A **transduction** mechanism that converts physical movement of a membrane into an electrical signal
 - The membrane is moved by the pressure wave
 - Electrical signal is generated based on inductance, capacitance or resistance
 - The signal is then sampled and digitized

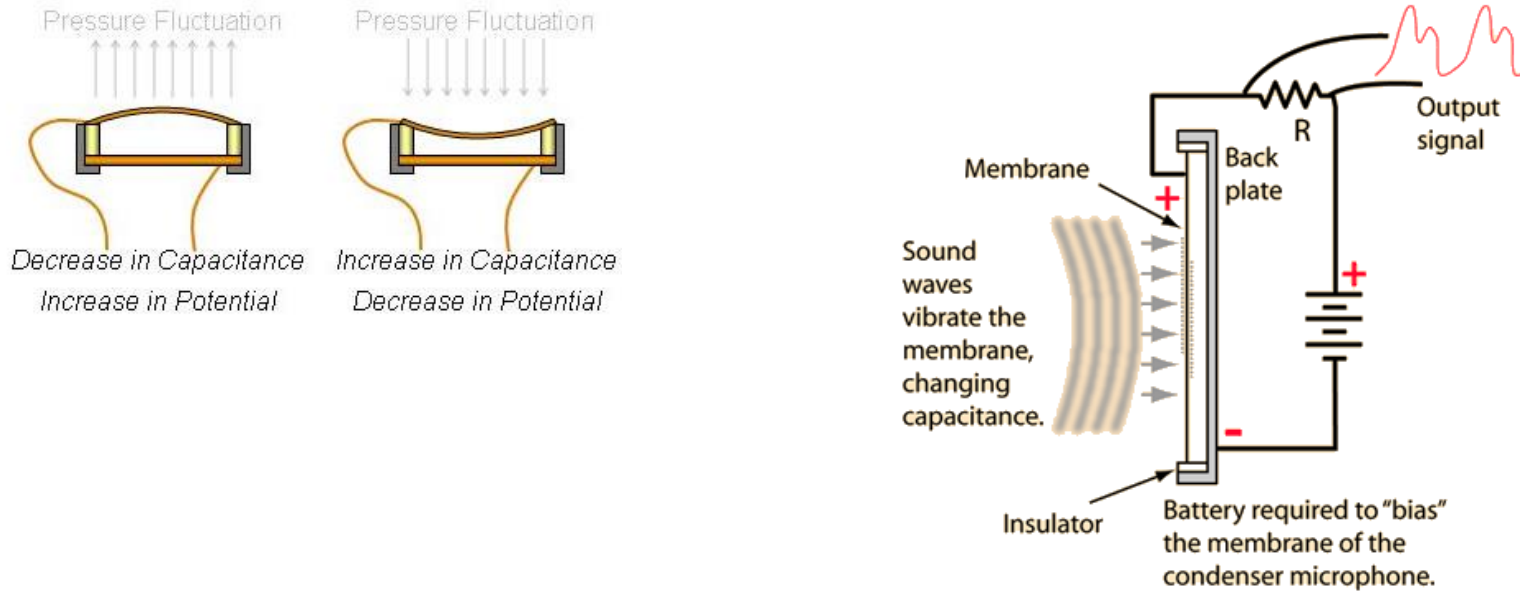
Dynamic Mics: Inductive



Based on inductance

- Time-varying voltage levels created by moving an element in an electromagnetic field
 - **Moving coil mic:** Diaphragm connected to a coil, suspended between magnets
 - Movement of diaphragms moves the coil and generates a voltage
 - **Ribbon mic:** A metal ribbon suspended between magnets
- Does not require phantom power
- Frequency response generally not flat
 - Also, small mics may be noise prone

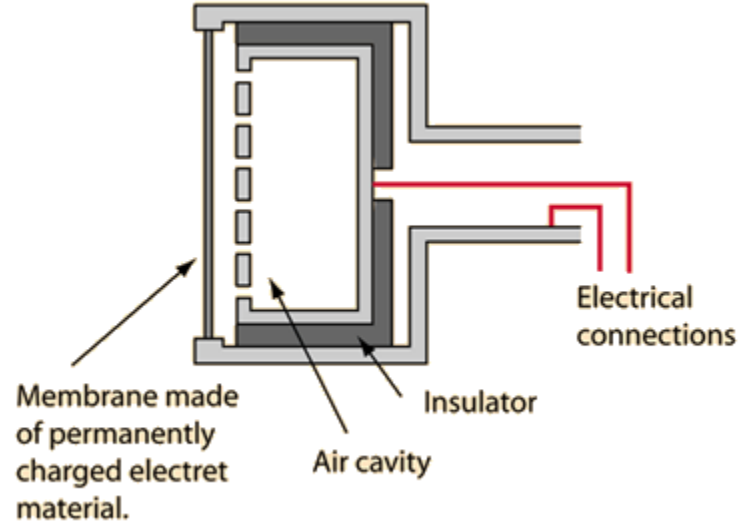
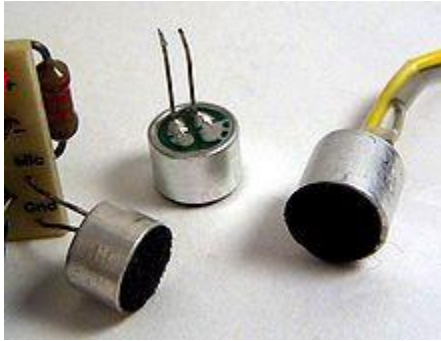
Condenser Microphones: Capacitance based



Based on *capacitance*

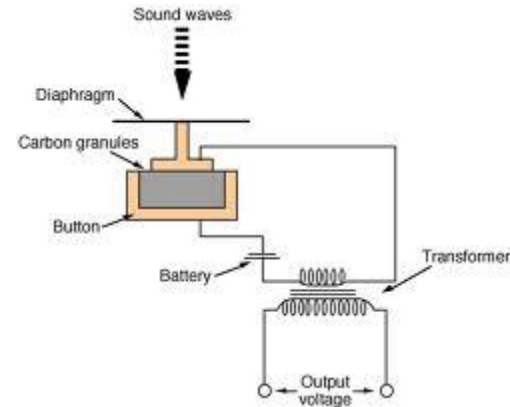
- Diaphragm and back plate maintained at different charge levels
- Pressure waves change the thickness of intervening dielectric (air)
 - Changing the capacitance of the capsule
 - Thereby changing the potential difference across the plates
- The changing potential difference is measured across the large resistor

Electret Microphones: Capacitance based



- Electret microphones are Condenser mics with fixed charge level
 - Either the diaphragm or the backplate carry a fixed charge
 - No external voltage source needed to maintain charge
 - But do require power for preamp (usually DC source)
 - Used to be low quality, no longer
- Most computer microphones are electrets
 - **Electret:** a dielectric material with a permanently embedded static electric dipole moment that will not decay for hundreds of years

Carbon Button Microphones: Resistance based



Based on *resistance*

- Carbon granules between a diaphragm and a backplate
- Motion of the diaphragm changes the resistance of the granule layer
- This changes the current in the circuit
 - Typically transformed to a higher voltage as shown
- Typical in older telephone handsets
 - Cellphones and recent handsets use electrets

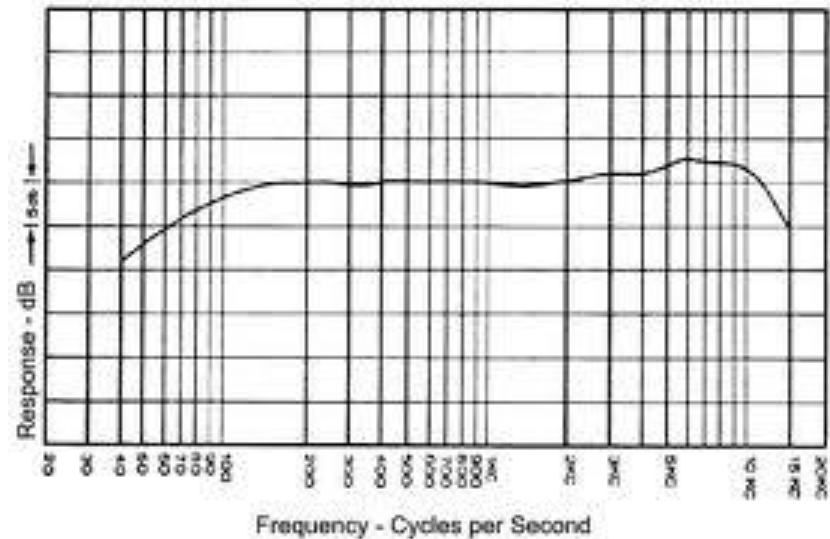
The First Microphone

- Graham Bell's "liquid" microphone
 - Metal cup filled with water, mixed with Sulphuric acid
 - Needle suspended in water (stuck to a diaphragm)
 - Sound waves move the needle up and down
 - Changing the resistance of the cup
 - A fixed voltage induces a time-varying current
- Sound can be recorded using unconventional devices!
- Bell's design was actually "copied" from Elisha Gray
 - Apparently from a drawing Bell saw in the patent office
 - In later incarnations he discarded the liquid microphone



© 1993 Smithsonian Institution

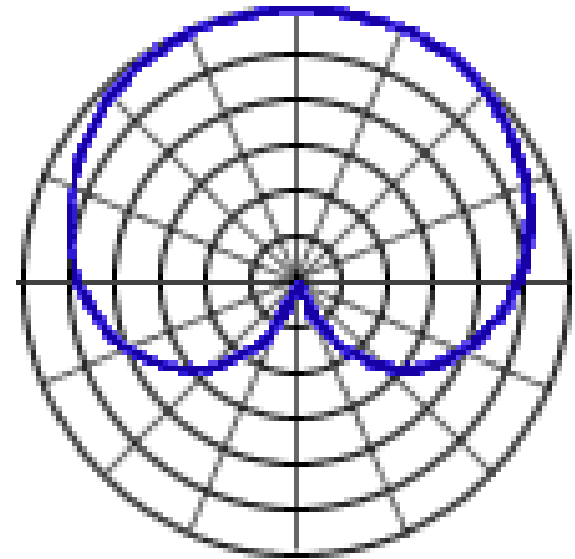
Choice of Microphones



- Characteristics to look for:
 - **Frequency response**
 - Must be flat in the frequencies of interest
 - **Directivity (or Directionality)**
 - More on this on next slide
 - **Noise level**
 - Usually stated as dB(A) SPL (Indicates the noise floor of the mic)
 - Lower is better
 - Good microphones: 20dB SPL, Ultra very good mics: 0dB
- Frequency response: Condenser microphones are usually used for high-quality recordings
 - Even cheap electret microphones nowadays can be surprisingly good

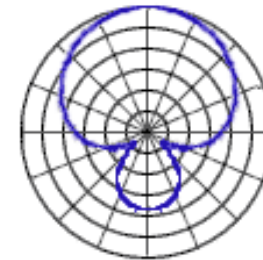
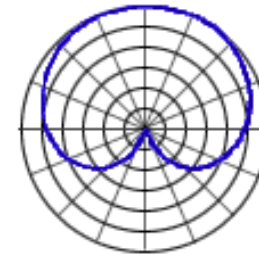
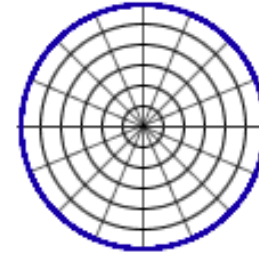
The “Directionality” of a Microphone

- The response of a microphone depends on the direction of the sound source
- Also depends on the construction of the microphone, as well as the shape of its chassis
- The response is typically represented as a polar plot
 - The distance of the curve from the “center” of the plot in any direction shows the “gain” of the mic to sounds from that direction



Typical Directivity Patterns

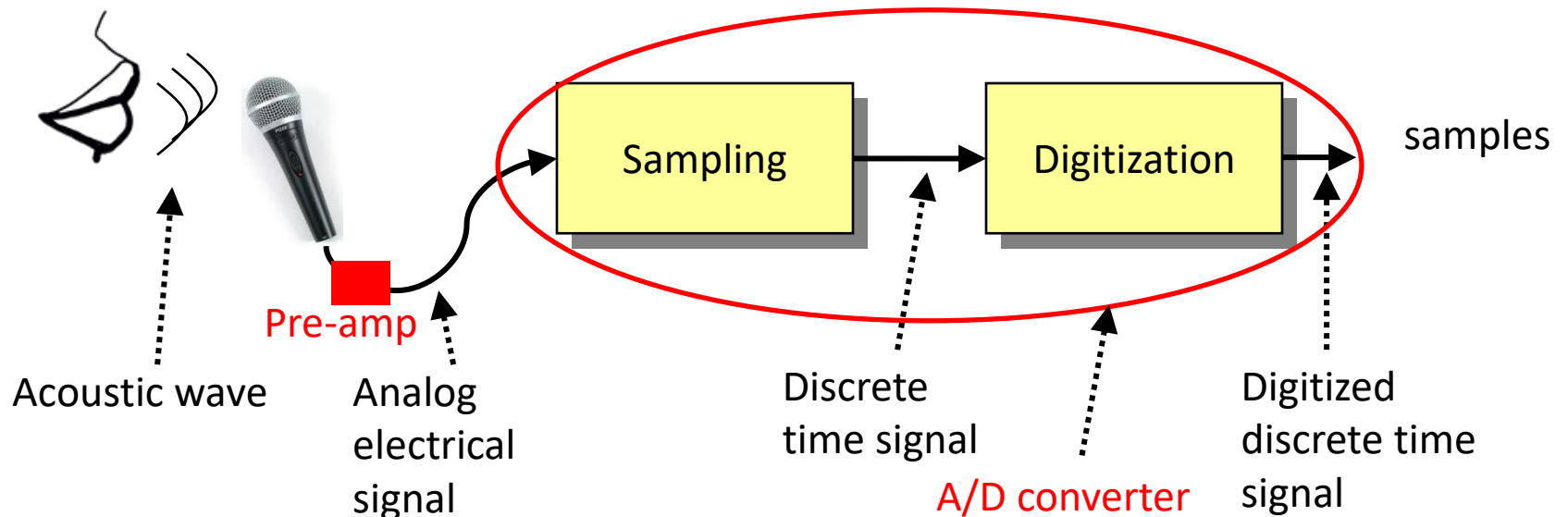
- Omnidirectional
 - Picks up sound uniformly from all directions
- Cardioid
 - Picks up sound from a relatively wide angle of directions from the front and some sound from the side
- Hyper cardioid
 - Picks up sound from a relatively narrow angle to the front, but also some noise from behind



Directional Patterns

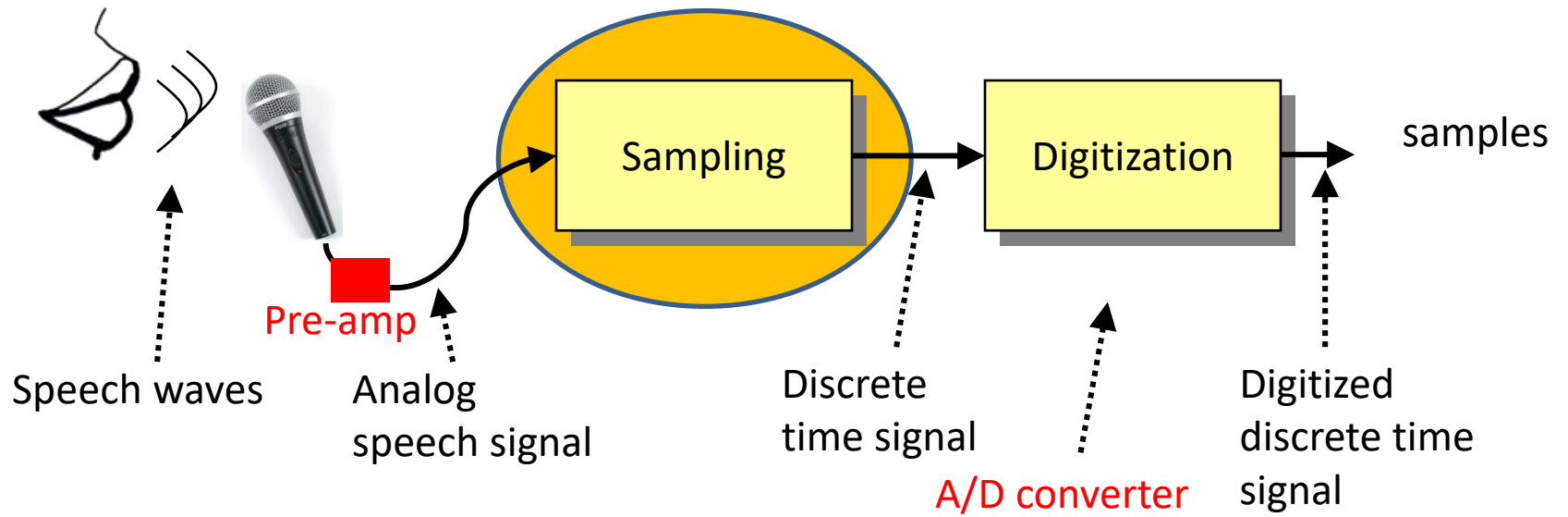
- Omnidirectional microphones are useful when we wish to pick up sounds from all directions
 - For speech applications, however, this will result in picking up a lot of noise from the surroundings
- Where the speaker location can be **somewhat localized**, Cardioid microphones are more precise
 - Pick up less non-speech noise
 - Still susceptible to noise
- Hypercardioids are a better choice when the speaker location can be **well localized**
 - Risk of picking up noise from behind
- In general, utilizing microphone directionality is the best way of minimizing extraneous noise from recordings

Sampling and Digitization

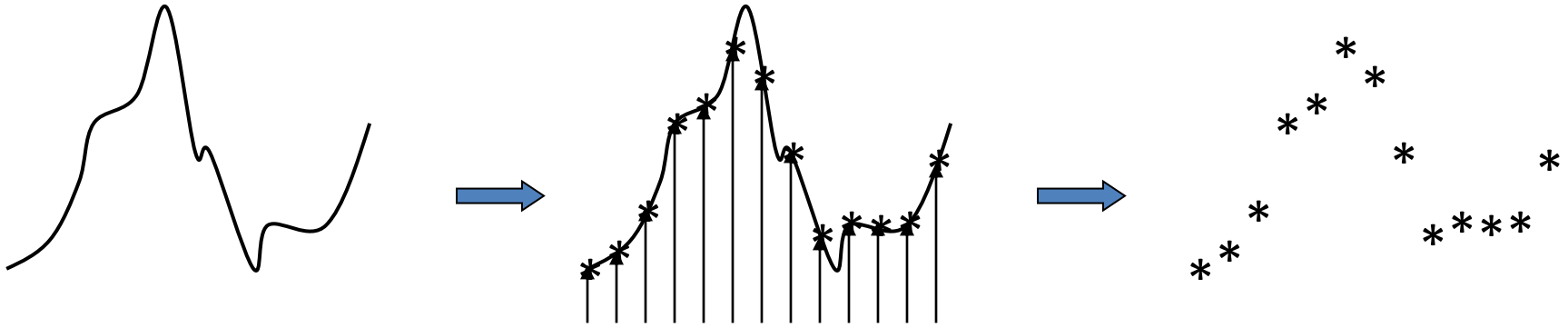


- The signal from the microphone goes through a *pre-amp*
 - The gain of which can be adjusted
 - The output of the pre-amp is a continuous-time electrical signal
 - Usually a voltage signal
- The signal is digitized by an *analog to digital converter* (A/D converter)
 1. **SAMPLING:** The signal value is sensed at regular, periodic intervals of time
 2. **DIGITIZATION/QUANTIZATION:** The value at each instant is assigned to one of a number of fixed values

Sampling



Sampling

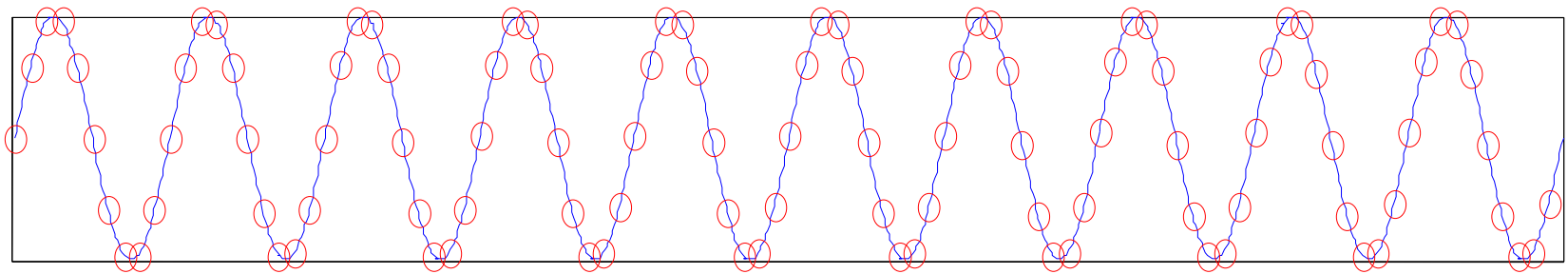


- Sampling
 - Signal values are captured at discrete intervals of time
- *Sampling frequency (or sampling **rate**)* is the number of samples taken per second
 - Usually measured in hertz (Hz) : 1 Hz = 1 sample per second

Sampling rate

Sampling requirement

- Sampling rate should be such that **all frequencies** in the analog signal are accurately represented in the digital signal
 - When the analog signal is recreated from the digital signal, it should be **exactly** like the original signal
 - Sample values are converted to levels of an electrical signal
 - The electrical signal moves a diaphragm back and forth to produce a pressure wave, sensed as sound
 - Played out on a speaker

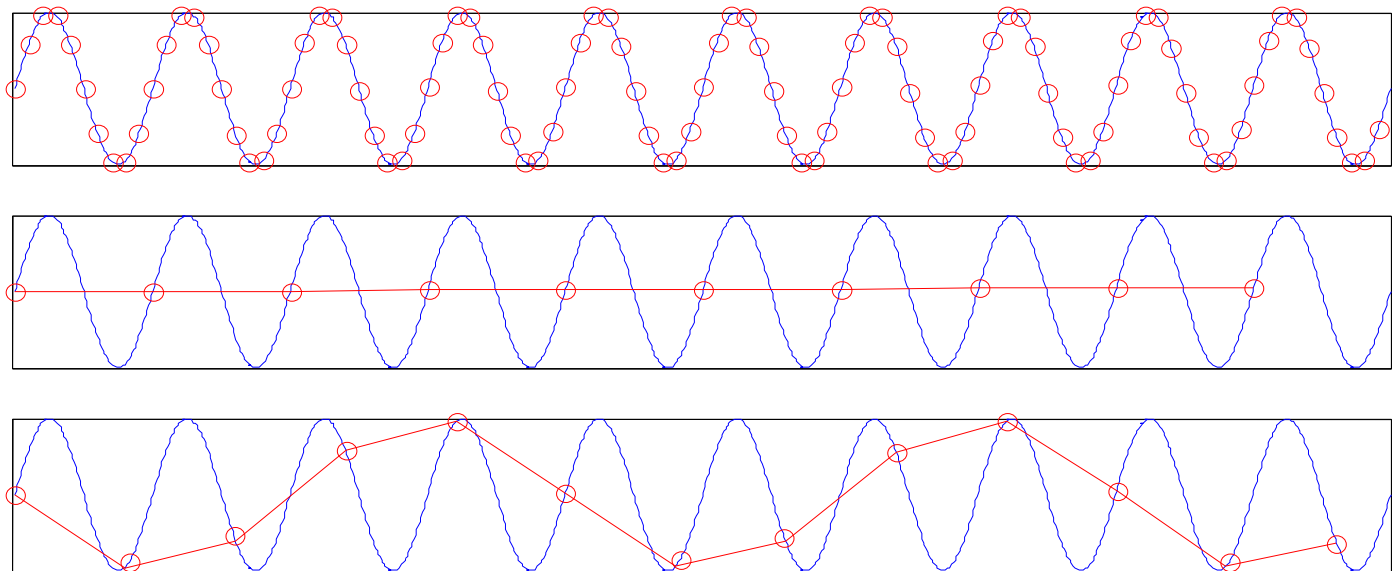


time →

Aliasing

- Low sampling rates result in *aliasing*
 - High frequencies appear as low frequencies

Aliasing:
Visual
Examples
Signal is
sampled
at less
than $2F$
in each
case

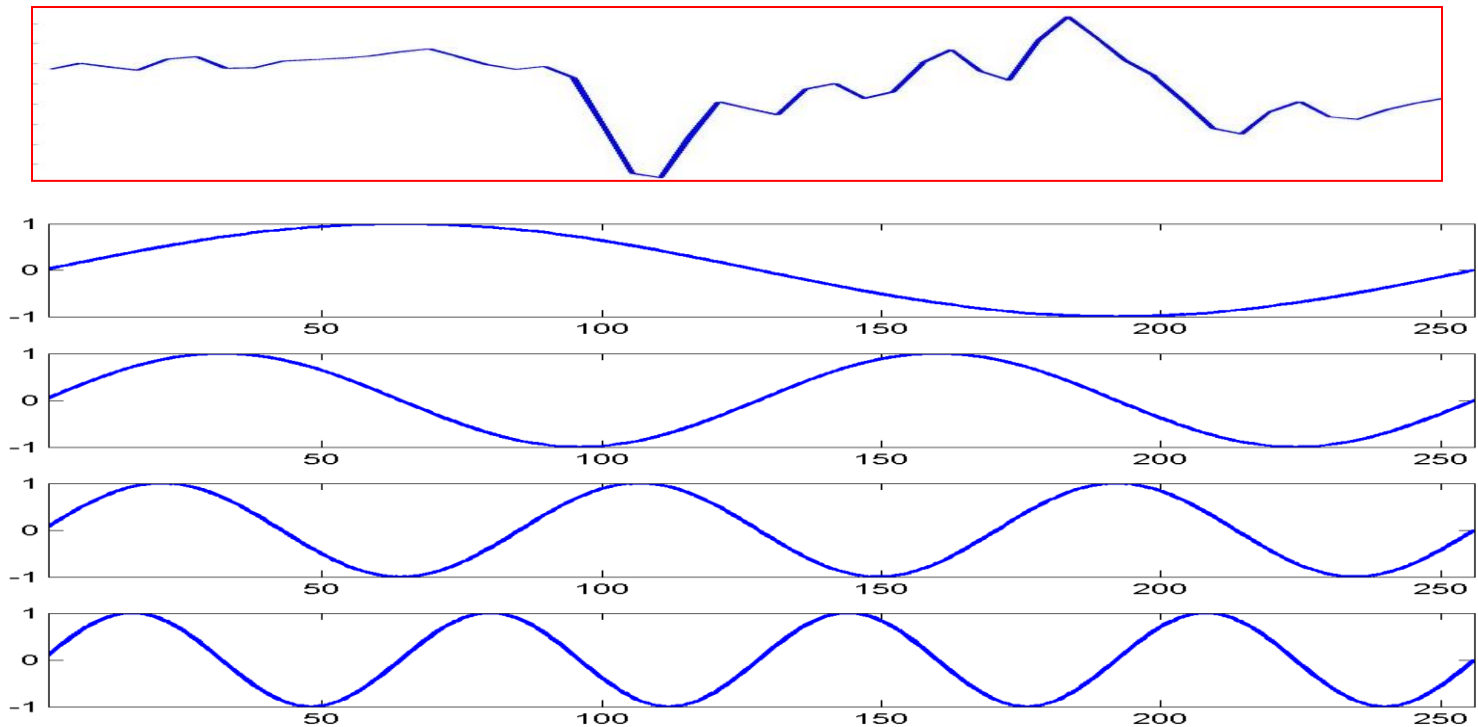


Aliasing

- The **Nyquist rate** is the minimum rate at which a finite bandwidth signal must be sampled to accurately retain or represent all component frequencies
 - **Nyquist theorem:** To represent any frequency F accurately, we must sample at the rate of at least $2F$ samples per second
 - **Nyquist frequency** (or folding frequency), named after Harry Nyquist, is one-half of the sampling rate (samples per second).
 - For a signal sampled at the rate of $2F$, the Nyquist frequency is F .
- Low sampling rates result in *aliasing*
 - Frequency $(F + f)$ will appear as $(F - f)$
 - e.g. 8000 sampling, Nyquist frequency is 4k, 5k Hz will appear as 3 k Hz

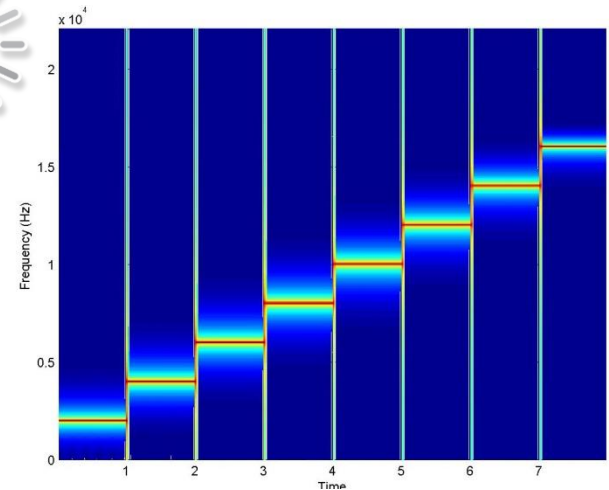
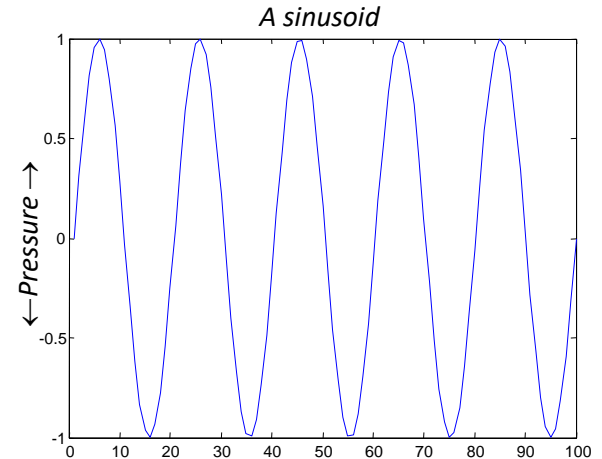
Sampling rate

- A sound wave is usually a superposition of waves of many frequencies
- **Fourier Theorem:** Any signal can be represented as a sum of sinusoids



How many samples a second

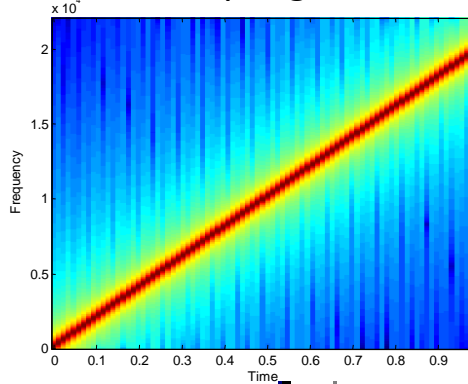
- Sounds can be modeled as the sum of many sinusoids of different frequencies
 - We sense the individual frequencies
 - Each hair cell in our inner ear is tuned to specific frequency
- Human hearing range
 - We can hear frequencies up to 16000Hz
 - Frequency components above 16000Hz can be heard by children and some young adults
 - Nearly nobody can hear over 20000Hz
 - This has design implications in the capture of sound



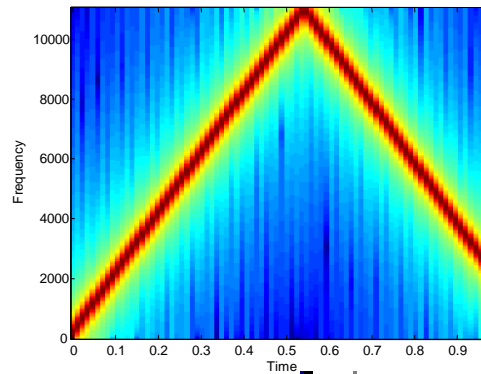
Aliasing examples

Sinusoid sweeping from 0Hz to 20kHz (a Chirp signal)

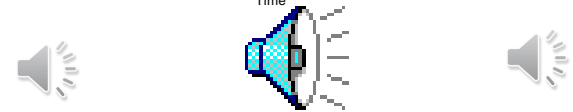
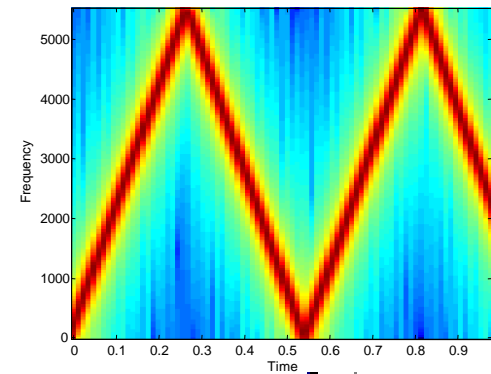
44kHz Sampling Rate, is ok



22kHz SR, aliasing!

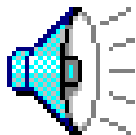


11kHz SR, double aliasing!

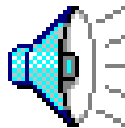


On real sounds

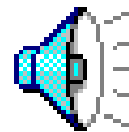
at 44kHz



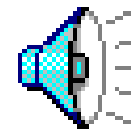
at 22kHz



at 11kHz



at 3kHz

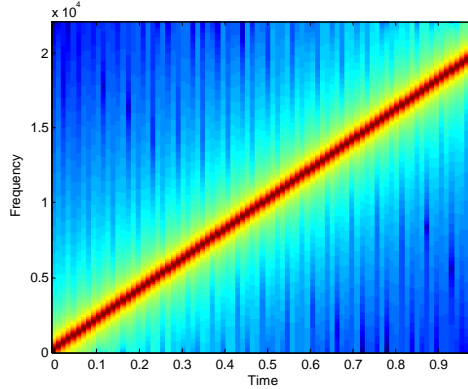


Anti-Aliasing

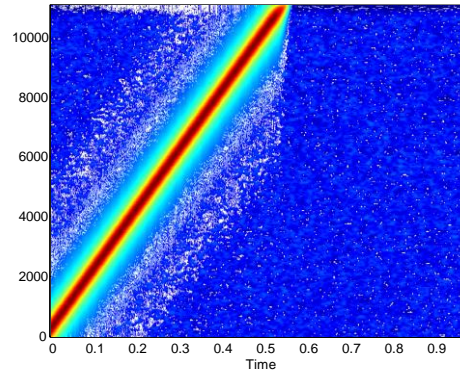
- To mitigate the effects of aliasing, we use anti-aliasing filters
 - Anti-aliasing is done to remove the frequencies above the Nyquist frequency from the signal

Anti Aliasing examples: Sweep 0-20kHz

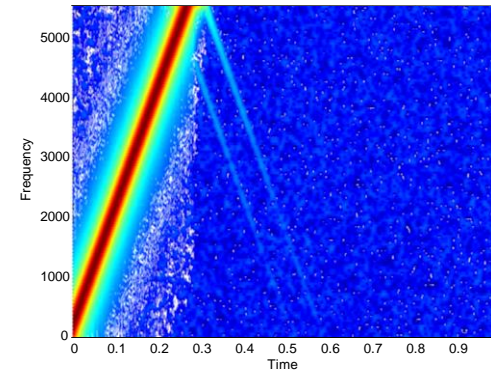
44kHz SR, is ok



22kHz SR, antialiased!



11kHz SR, antialiased!



Prior to 22kHz
downsampling

Prior to 11kHz
downsampling

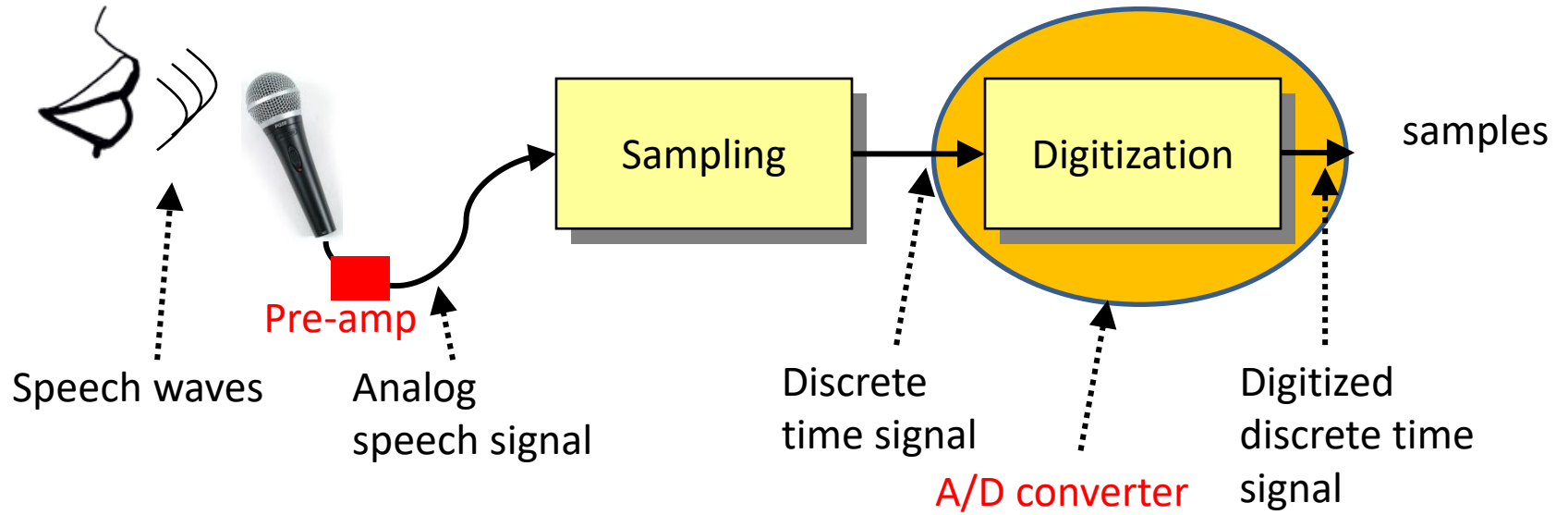
Anti-aliasing filter

Anti-aliasing filter

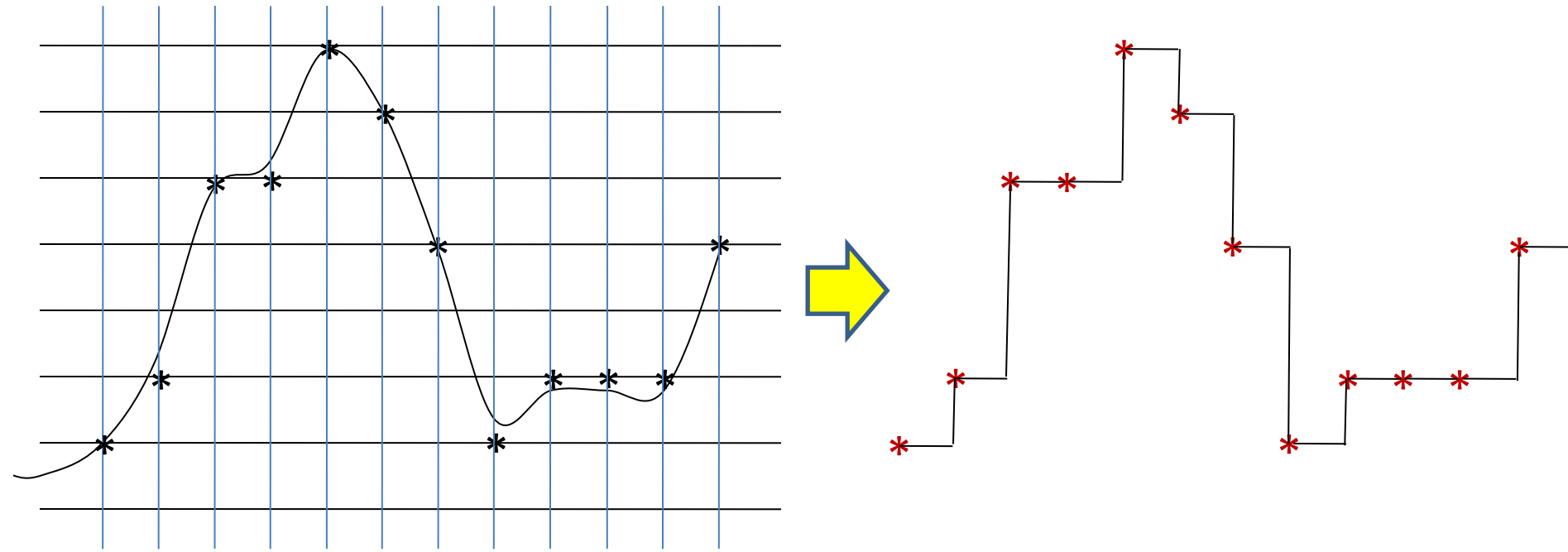
Typical Sampling Frequencies

- Audio hardware in practical systems typically supports several standard rates
 - CD recording employs 44.1 kHz per channel – high enough to represent most signals you can hear most faithfully
 - Telephone data is narrowband and has frequencies only up to 4 kHz
 - Good microphones provide a wideband speech signal. 16kHz sampling can represent audio frequencies up to 8 kHz. This is considered sufficient for capturing all speech content
 - Used for speech recognition, where possible
 - For speech 8kHz, 11.5kHz and 16kHz and 44.1 kHz are all used
 - Sampling hardware will usually employ appropriate anti-aliasing
 - Not always – be careful

Digitization

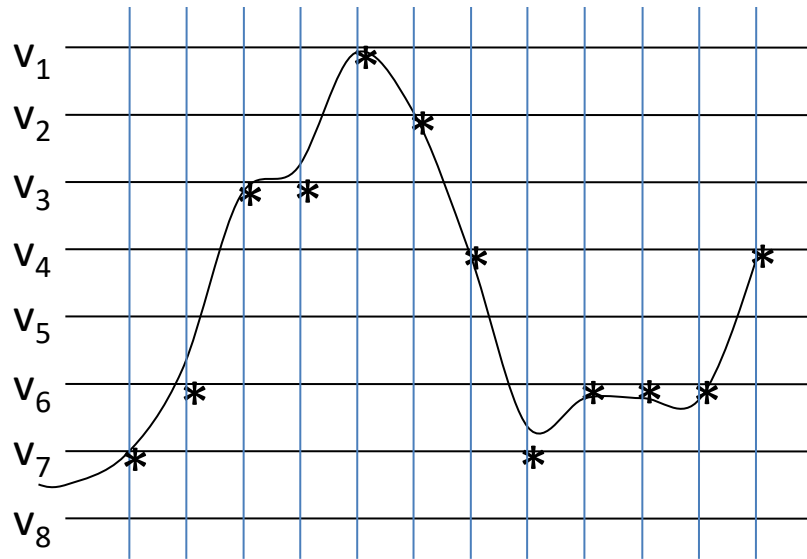


Sample Resolution: the need to quantize



- Samples cannot take just any value
 - Each sample is represented by a limited number of bits
 - Sample values are usually integers, restricted to a finite set of values
 - Typically 2^N integer values, for $N = 8, 16, 32$
 - Signed 16-bit integer is $-32,768$ (-1×2^{15}) through $32,767$ ($2^{15} - 1$);
 - Unsigned 16-bit integer: can take values 0 through $65,535$ ($2^{16} - 1$)

Process of Quantization



| index | Value |
|-------|-------|
| 0 | v_1 |
| 1 | v_2 |
| 2 | v_3 |
| 3 | v_4 |
| 4 | v_5 |
| 5 | v_6 |
| 6 | v_7 |
| 7 | v_8 |

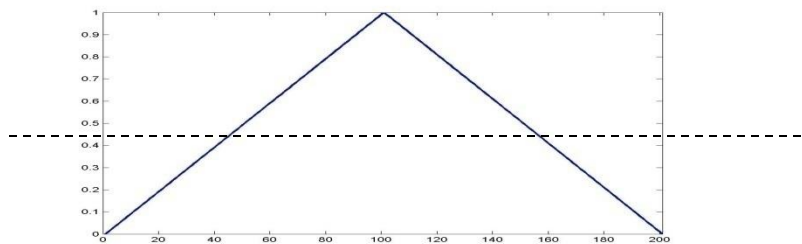
- Map the sample value read from the signal to the closest integer value
 - Or use an integer index mapped to some close value
 - Such mappings are specified in a table of values
 - When recreating the signal, The actual signal value is read off from the table
 - When indices are the same as values
 - **Only the indices are stored!**

Sample resolution

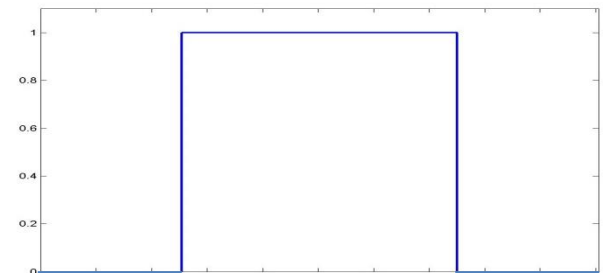
- Example of **1-bit Uniform** quantization table
 - The range of signal values is binned into equal intervals

| Signal Value | Bit sequence | Mapped to |
|---------------|--------------|--------------------|
| $S > 2.5v$ | 1 | $1 * \text{const}$ |
| $S \leq 2.5v$ | 0 | 0 |

Only 2 levels, using a threshold of 2.5 volts
(this varies with hardware)



Original Signal

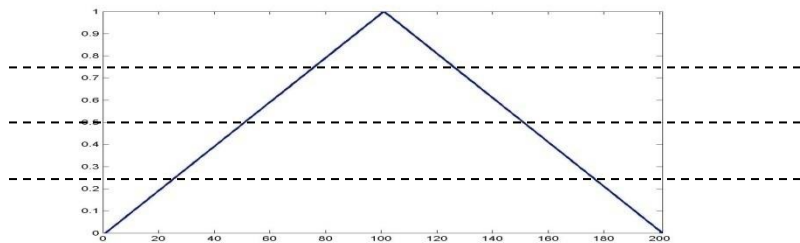


Quantized approximation

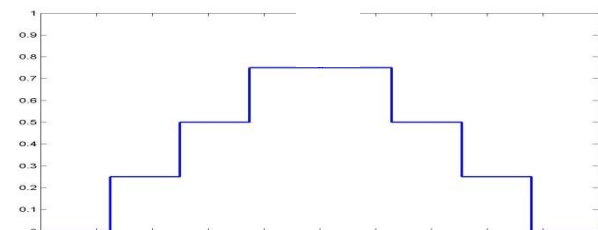
Mapping signals into bits

- Example of **2-bit Uniform** quantization table:
 - 4 equally-spaced levels, (voltage threshold/range of 1.25 for each bin)

| Signal Value | Bit sequence | Mapped to |
|-------------------------------------|--------------|--------------------|
| $S \geq 3.75\text{v}$ | 11 | $3 * \text{const}$ |
| $3.75\text{v} > S \geq 2.5\text{v}$ | 10 | $2 * \text{const}$ |
| $2.5\text{v} > S \geq 1.25\text{v}$ | 01 | $1 * \text{const}$ |
| $1.25\text{v} > S \geq 0\text{v}$ | 00 | 0 |



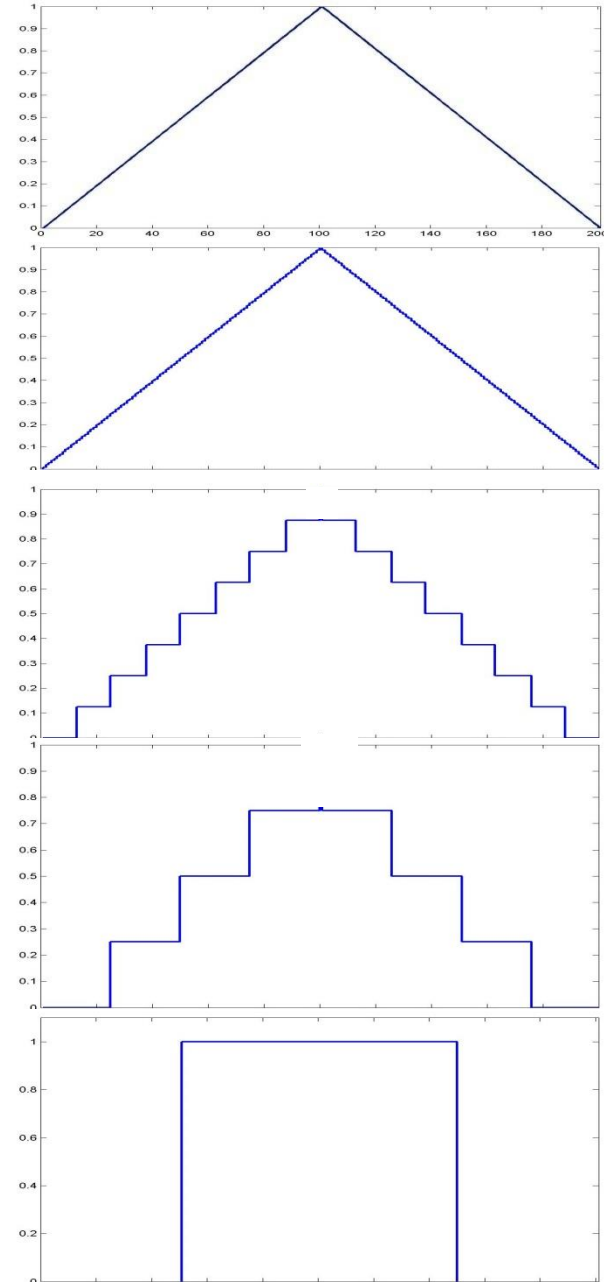
Original Signal



Quantized approximation

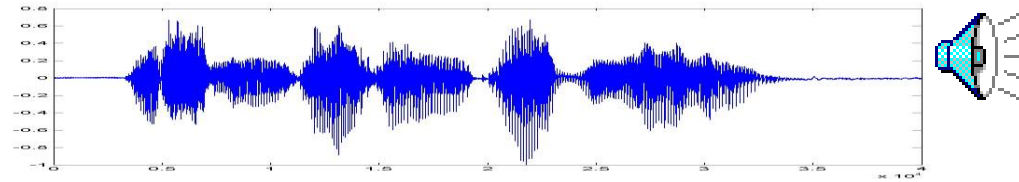
Uniform quantization, different levels

- The original signal
- 8 bit quantization
- 3 bit quantization
- 2 bit quantization
- 1 bit quantization

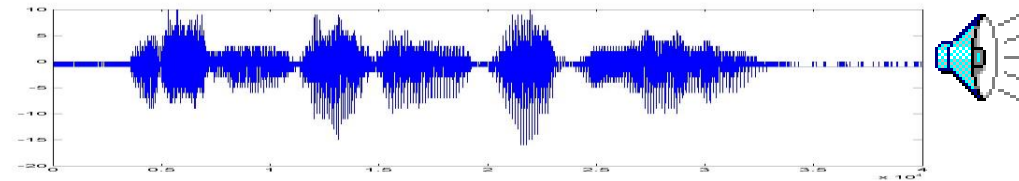


Tom Sullivan Says his Name

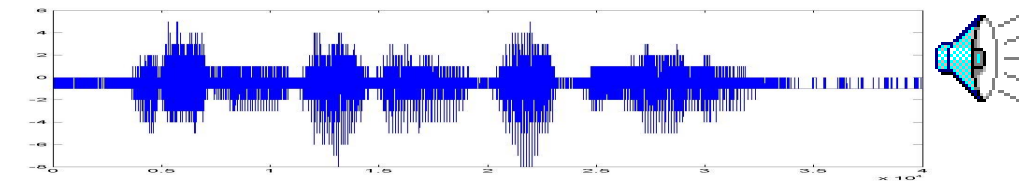
- 16 bit sampling



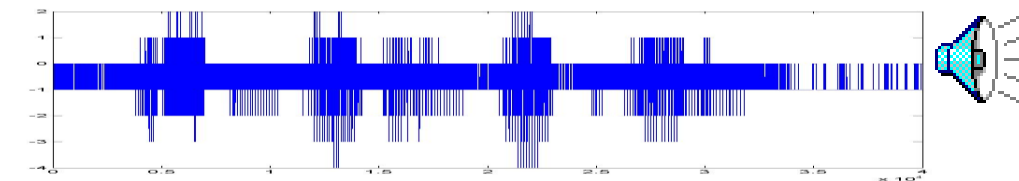
- 5 bit sampling



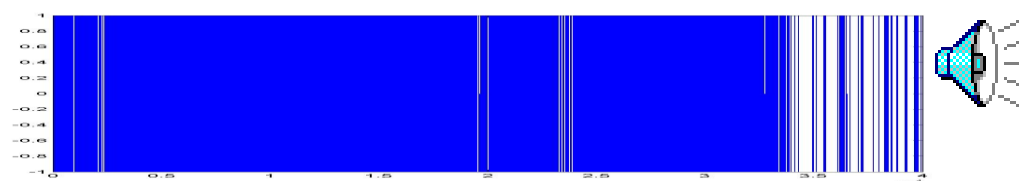
- 4 bit sampling



- 3 bit sampling

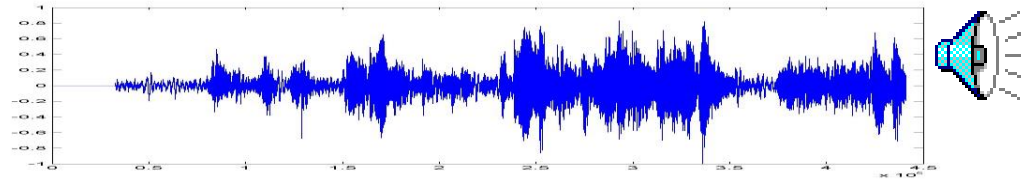


- 1 bit sampling

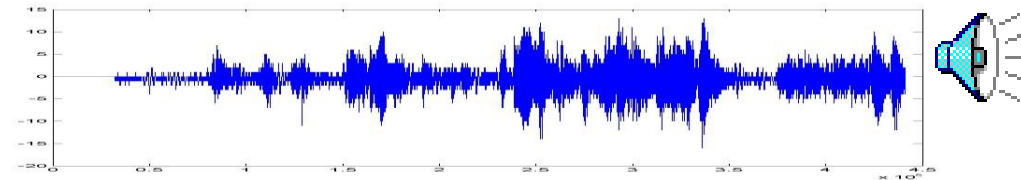


A Schubert Piece

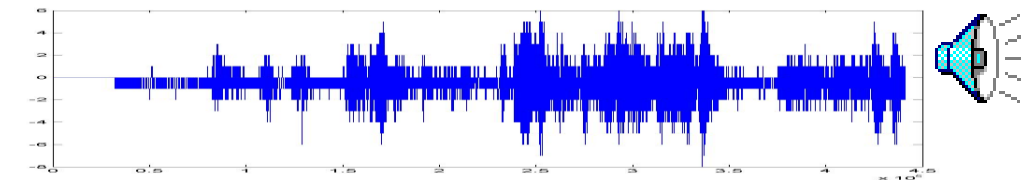
- 16 bit sampling



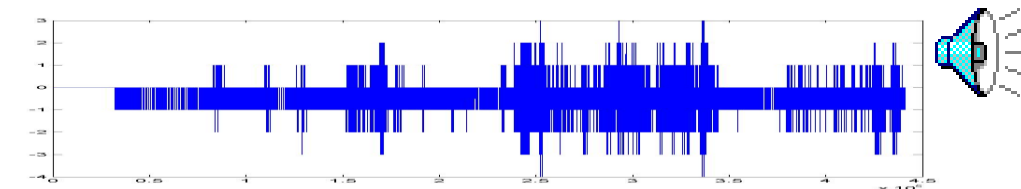
- 5 bit sampling



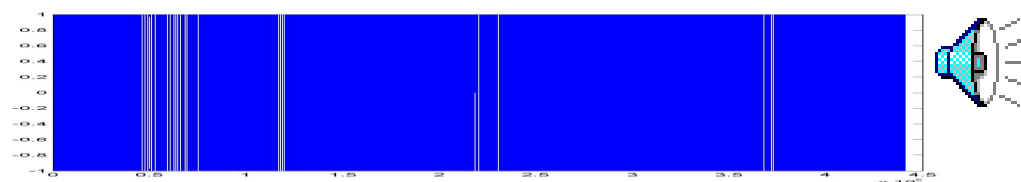
- 4 bit sampling



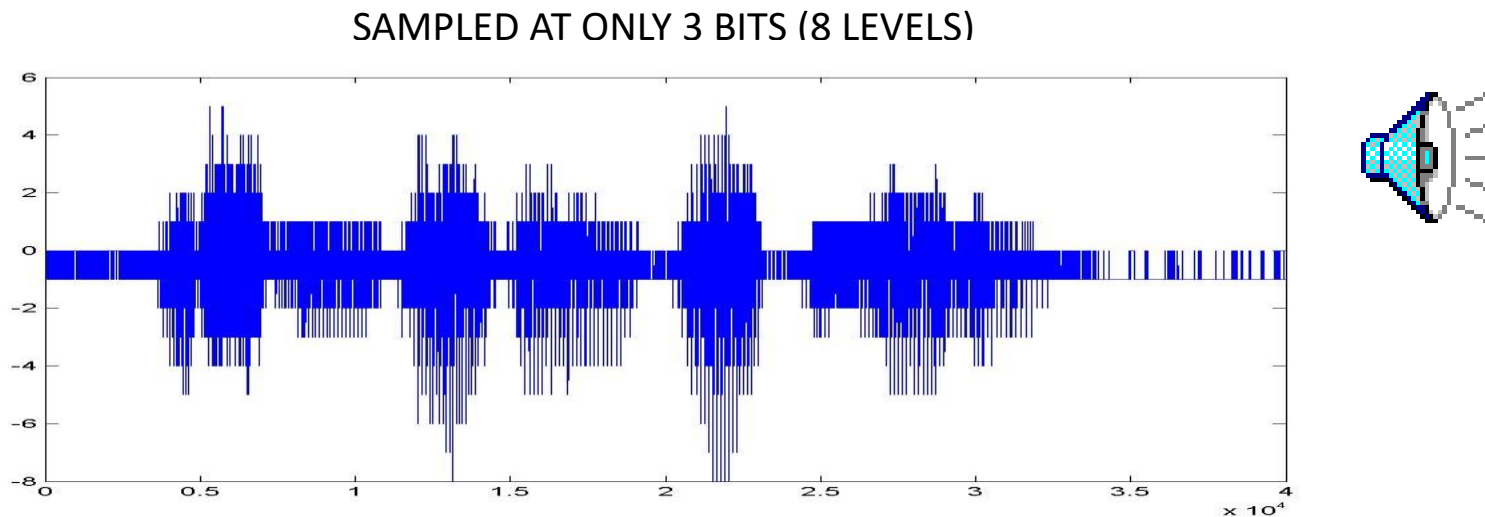
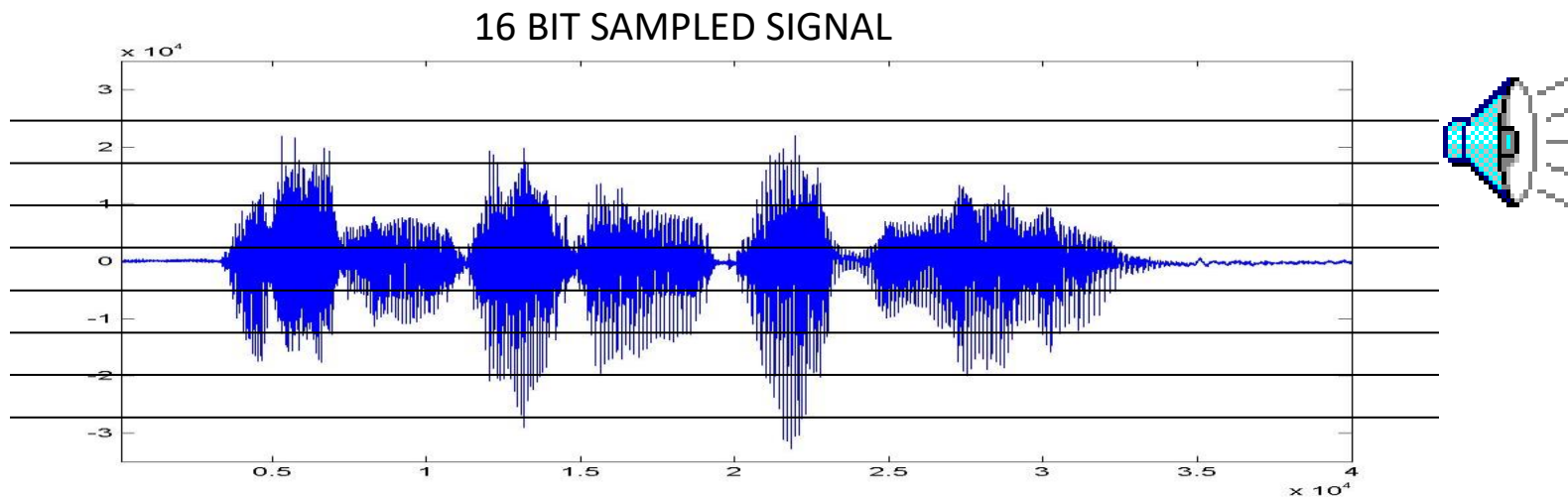
- 3 bit sampling



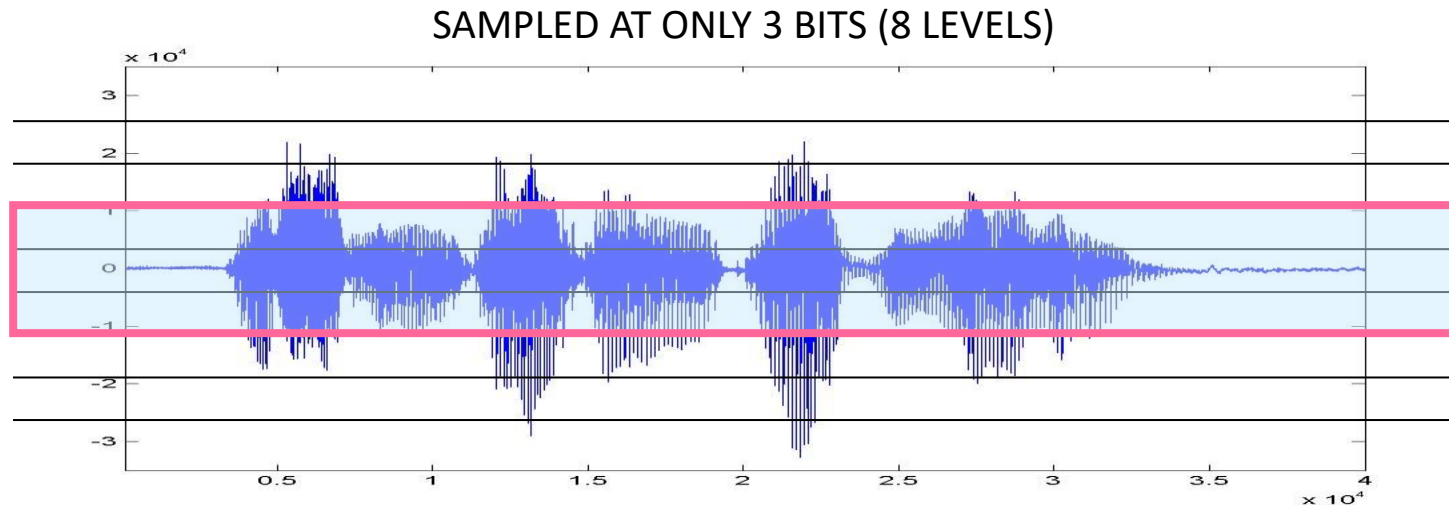
- 1 bit sampling



Improving on Uniform quantization

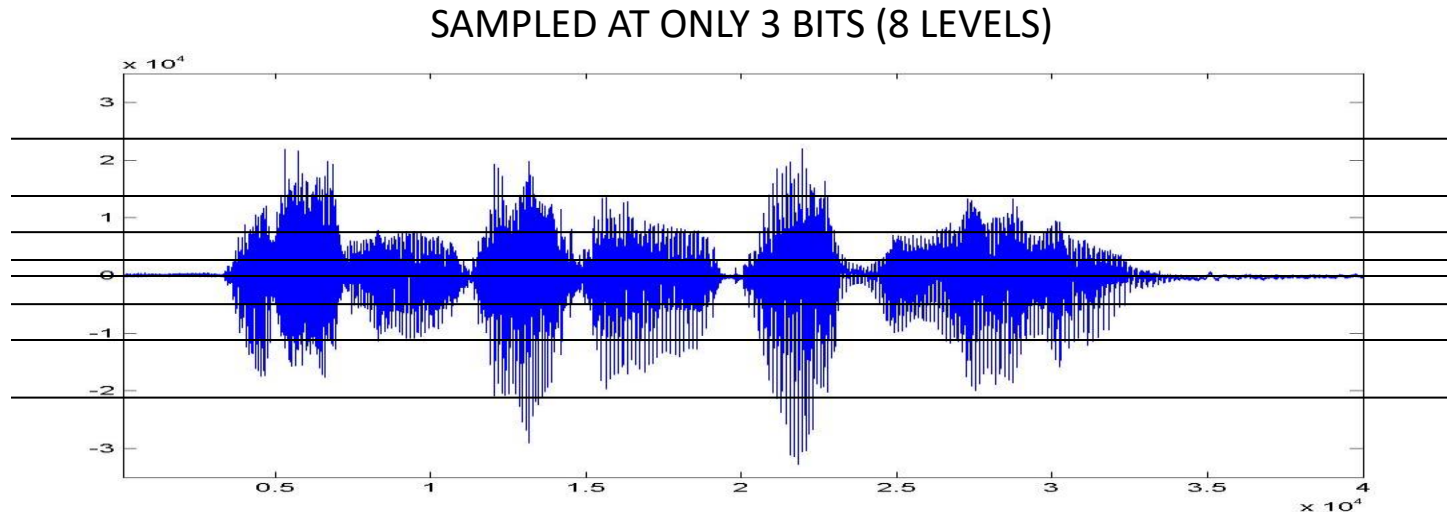


Improving on Uniform quantization



- There is a lot more action in the central region than outside.
- Assigning only four levels to the busy central region and four entire levels to the sparse outer region is inefficient

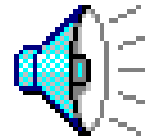
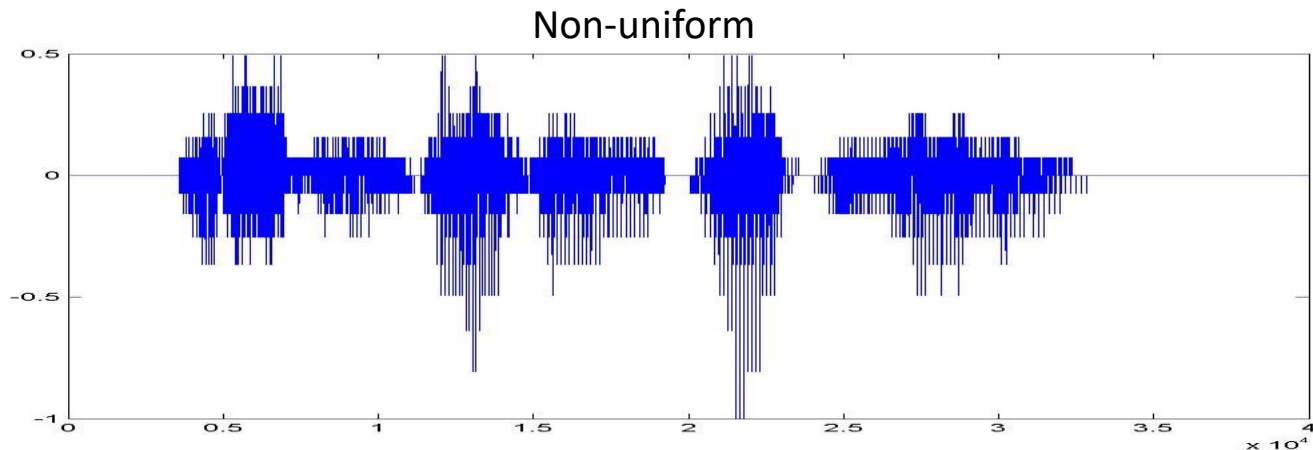
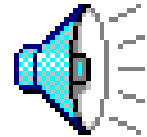
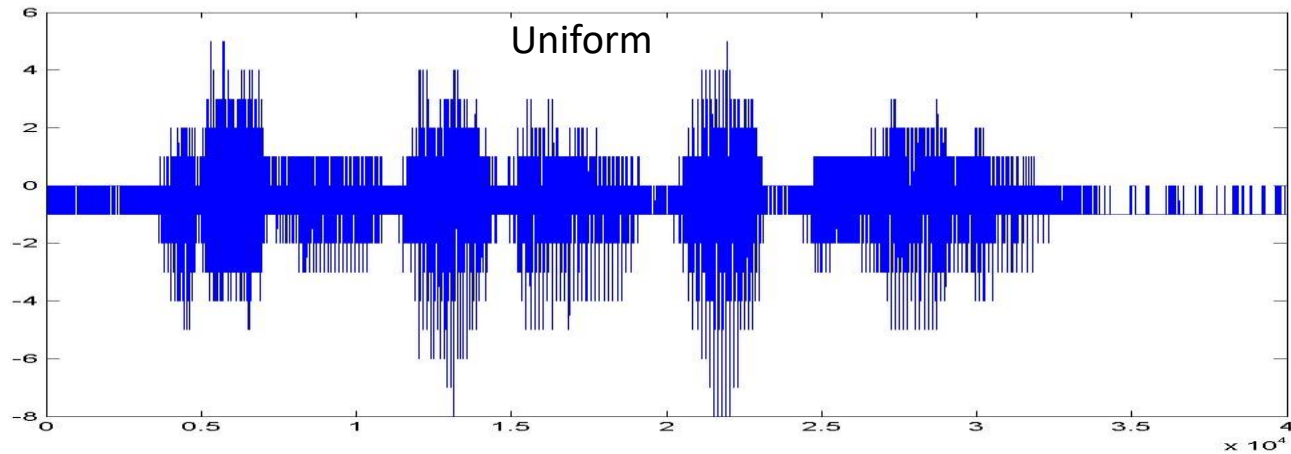
Improving on Uniform quantization



- Assigning more levels to the central region and less to the outer region can give better fidelity
 - for the same overall number of levels

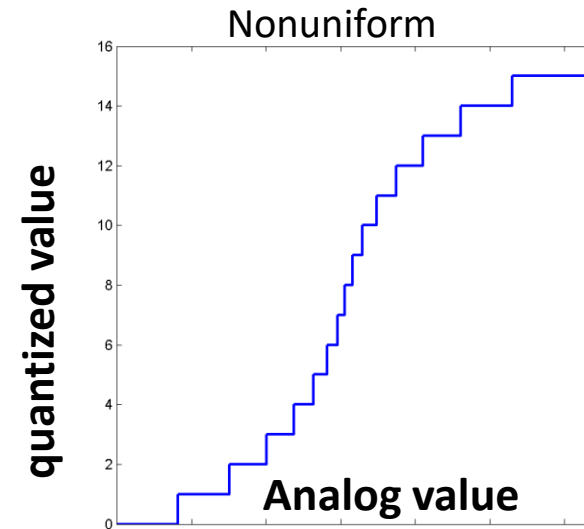
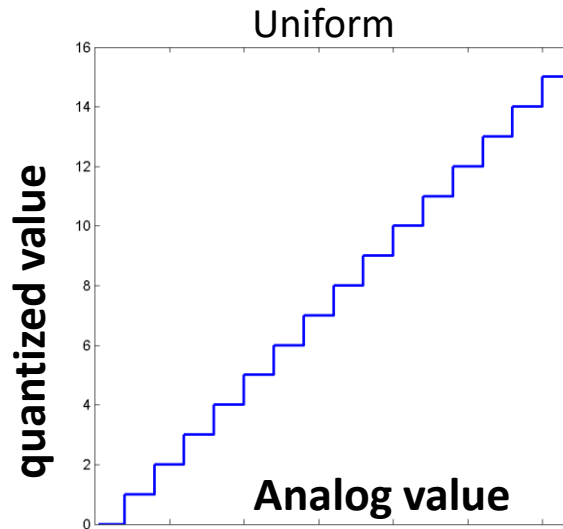
Uniform vs. Non-uniform quantization

SAMPLED AT ONLY 3 BITS (8 LEVELS)



- Assigning more levels to the central region and less to the outer region can give better fidelity for the same storage

Uniform vs. Non-uniform quantization

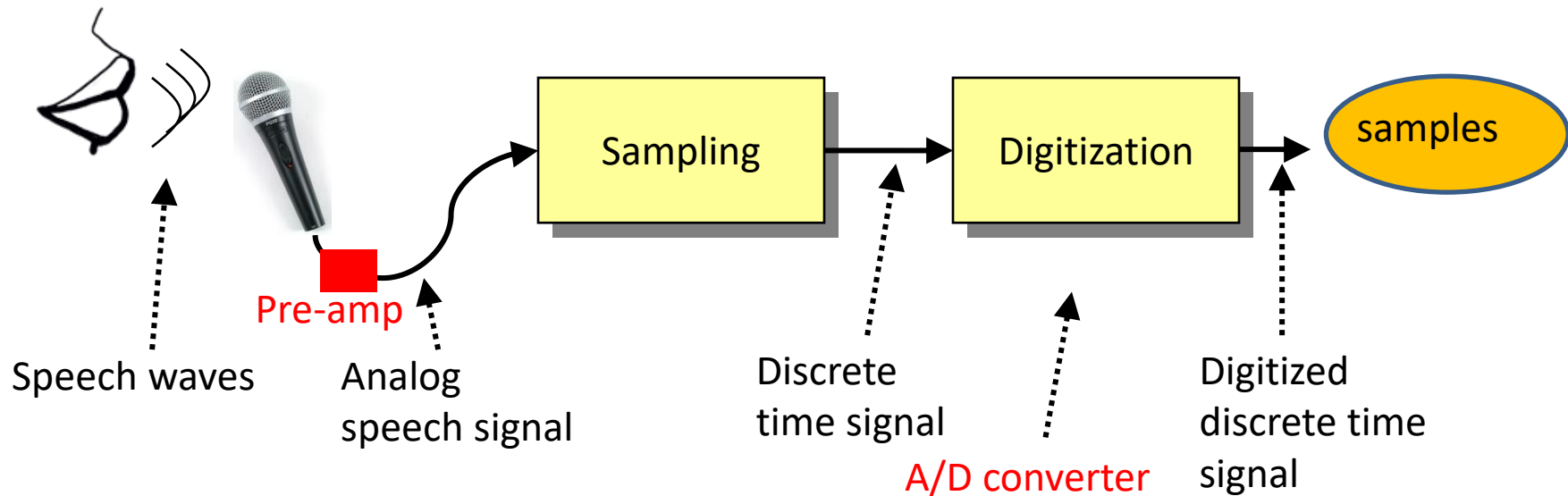


- Uniform quantization maps uniform widths of the analog signal to units steps of the quantized signal
- In non-uniform sampling the step sizes are smaller near 0 and wider farther away
 - In practical devices one can get the same perceptual effect with 8 bits of non-linear quantization as 12 bits of linear quantization
 - Most effective functions follow a logarithmic law:
 - Mu-Law: $Y = C \cdot \log(1 + \mu X/C)/(1+\mu)$
 - A-Law: $Y = C \cdot (1 + \log(a \cdot X)/C)/(1+a)$

For playback or processing....

- The numbers returned by an ADC (analog-to-digital converter) or read from a file are *indices* into the table
 - The indices must be converted back to actual values
 - For uniformly sampled data the indices are proportional to the value and can be used directly
 - Called “Linear PCM” or “PCM” encoding. For sound, usually 16-bit sample values are sufficient in the PCM signal are sufficient
 - For non-uniform sampling, table lookup (or function inversion) is required. Often done automatically by most audio libraries when reading audio

The samples are generated...



- How are they stored?
 - The sequence of numbers generated is a “channel” of audio
 - May not necessarily produce just one sequence of numbers for a given sampling rate
 - May be multiple sequences multiplexed with each other (2= stereo signal, 2 or more = multichannel signal)
 - Depends on the hardware and how we set up the recording

Storing Audio/Speech Files

- The data are typically written in binary, but the files typically have headers that can be read as ascii text.
 - Headers store critical information such as byte order, no. of samples, coding type, bits per sample (“bit-depth”), sampling rate etc
 - **There are many storage formats in use. 3 categories of audio formats:**
 - **Uncompressed**
 - **Lossy Compressed**
 - **Lossless Compressed**
- Audio files must be converted from stored format to linear PCM format for further processing
 - Audio I/O library routines will usually process the headers to obtain the necessary information and perform the appropriate conversion. We can also write the I/O ourselves

Audio file formats and codecs

UNCOMPRESSED FORMATS

- Digitized signal stored without further processing, take up a lot of storage space
 - Especially if it a stereo or multichannel signal

PCM (“raw” data, *.raw)

- PCM (Pulse-Code Modulation): Analog sound waves sampled at regular intervals (pulses) and stored. This term is interchangeably used with LPCM (Linear Pulse-Code Modulation) in which samples are taken at linear intervals
- Commonly used in CDs and DVDs

NIST (*.sph)

- container
- 1024 byte ascii header, followed by PCM format

SUN (*.au, *.snd)

- Legacy

Audio file formats and codecs

WAV (*.wav)

- **Microsoft PCM and ADPCM:** “windows” audio.
 - Standard for this was developed by Microsoft and IBM in 1991.
 - In PCM, data for .WAV files is stored using linear samples 8 bits per sample
 - ADPCM (Adaptive Delta (or Differential) Pulse Code Modulation)
 - Stores deltas between samples at 4 bits per sample
 - ADPCM takes up half the disk space as PCM. Suitable for longer files
 - Makes more assumptions that allow data reduction. Because of these assumptions, low frequencies are properly reproduced, but high frequencies tend to get distorted. The distortion is easily audible in 11 kHz ADPCM files
- **WAV (*.wav)** (When Waveform Audio V)
 - Windows **container** for audio formats
 - WAV files can contain both compressed and uncompressed audio, but mostly always contain uncompressed audio in PCM format

Audio file formats and codecs

AIFF (*.aiff)

- AIFF (Audio Interchange File Format)
 - Developed by Apple for Mac systems in 1988. Windows systems can also handle them.
 - A **container**
 - Can contain multiple kinds of audio, e.g. AIFF-C (contains compressed audio)
 - Most AIFF files contain uncompressed audio in PCM format.

Audio file formats and codecs

LOSSY COMPRESSION

- Some information is lost, mostly imperceptible to the human ear
 - If compressed too much or too often, audio can have perceptual artifacts

MP3 (*.mp3)

- MP3 (MPEG-1 Audio Layer 3, MPEG-2 Audio Layer 3, NOT THE SAME AS MPEG-3)
 - MPEG-1 and MPEG-2 are multimedia containers
 - MPEG-3 : A container. Was designed to handle HDTV signals at 1080p in the range of 20 to 40 megabits per second
- Introduced in 1993. Popular for music. Almost all devices can handle it today
- Based on perceptual coding. Removes frequencies that are outside of hearing capabilities of people (decided based on psychoacoustic principles)
 - Reduce the quality of sounds that are not easy to hear, compress the important parts of the audio as efficiently as possible

Audio file formats and codecs

AAC (*.aac)

- AAC (Advanced Audio Coding).
- Developed in 1997 as the successor to MP3. The compression algorithm is more advanced. Codec is more efficient. At the same bitrate, AAC has better sound quality than MP3
- Standard audio compression method used by YouTube, Android, iOS, iTunes, PlayStations etc.

OGG (Vorbis) (*.ogg)

- A **multimedia container** that can hold all kinds of compression formats, but is most commonly used to hold Vorbis files. Vorbis was first released in 2000. Opensource. Better than most other lossy compression formats (smaller file size for equivalent audio quality). Good for mid to high quality (8kHz-48.0kHz, 16+ bit, polyphonic) audio and music at fixed and variable bitrates from 16 to 128 kbps/channel)
- **Vorbis** is a free and open-source software project headed by the Xiph.Org Foundation. The project produced Ogg.

WMA (*.wma)

- WMA (Windows Media Audio). Released in 1999. Proprietary to Microsoft
- Developed to address flaws in the MP3 compression.
- Compression algorithm quite similar to AAC and OGG. Better quality than MP3
- Not many devices/platforms support it since it is proprietary

Audio file formats and codecs

LOSSLESS COMPRESSION

- Compresses without any loss from original source.
- Not as efficient as lossy compression. Equivalent files can be 2x to 5x larger
- Do not confuse lossless compression with high-resolution audio (which is most likely a scam)

FLAC (*.flac or *.fla)

- FLAC (Free Lossless Audio Codec). Developed in 2001.
- Can compress up to 60% without losing a single bit of information. Open source.
- Alternative to MP3 for CD audio.
- Full quality of raw uncompressed audio in half the file size

Audio file formats and codecs

ALAC (*.alac)

- ALAC (Apple Lossless Audio Codec). Developed in 2004. Proprietary initially, became opensource in 2011
- less efficient than FLAC. But Apple users are forced to use it because iTunes and iOS support ALAC and do not support FLAC at all.

WMA (*.wma)

- Microsoft created the format to avoid the licensing issues associated with the MP3 format
- WMA (Windows Media Audio). Lossless alternative to lossy WMA This version is called “WMA Lossless”. Used the same extension as the lossy version.
- Worse than FLAC and ALAC in compression efficiency. Proprietary. Supported by Windows and Mac o/s.

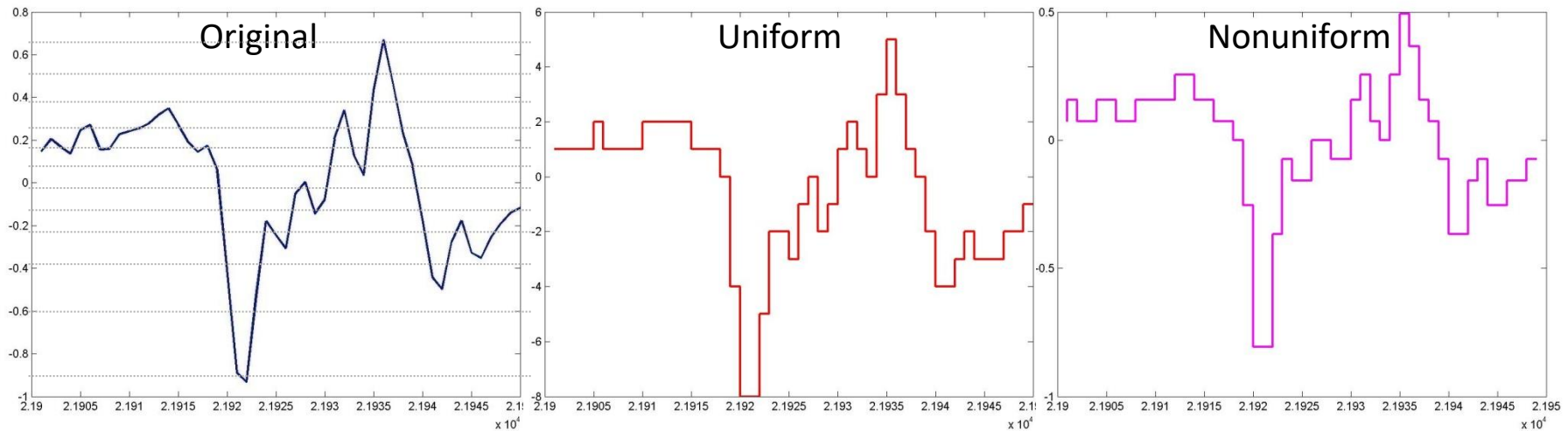
Editing audio

- Should be done on raw/uncompressed format

Signal Quality

- Must be mindful of signal quality when processing audio files of an application
- The **quality** of the digitized signal depends critically on many factors:
 - **The electronics performing sampling and digitization**
 - **Microphone quality**
 - **Poor quality electronics** can severely degrade signal quality
 - *E.g.* Disk or memory bus activity can inject noise into the analog circuitry
 - **Anti-aliasing (proper/improper)**
 - Not using an anti-aliasing filter is a cause for many problems
 - **Quantization levels (sufficient/insufficient)**
 - Minimally 16 bit PCM or 8 bit Mu / A-law is needed
 - **Proper setting of the recording level**
 - Too low a level underutilizes available signal range, increasing susceptibility to noise
 - Too high a level can cause *clipping*
 - **Ambient noise in recording environment**
 - **Recording channel: bandwidth and distortion**
 - **Echoes, reverberation, other sound-altering effects of the recording environment**

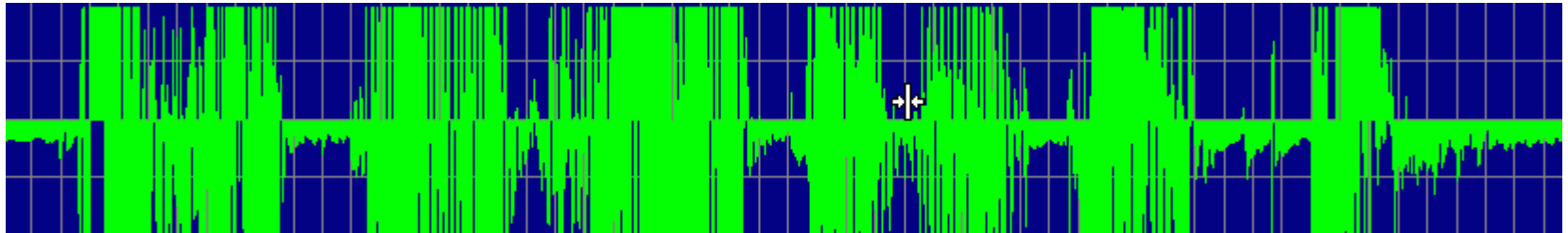
Proper recording levels: clipping



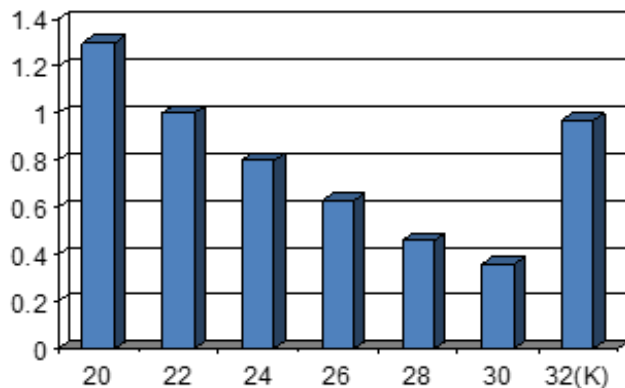
- At each sampling instant, the value of the waveform is rounded off to the nearest quantization level
- Samples outside the quantization range are given either the highest or lowest quantized values
 - E.g. max value for 8 bit quantization is 255. If sample = 300, it is set to 255
- If the recording level is too high, most values go outside the range, and are set to the highest value we can represent. The waveform then has a tabletop shape.
 - **This is called clipping**

Clipping

- Clipping and non-linear distortion are the most common problems in audio recordings
 - While recording, clipping can be fixed by reducing signal gain (but AGC is not good)

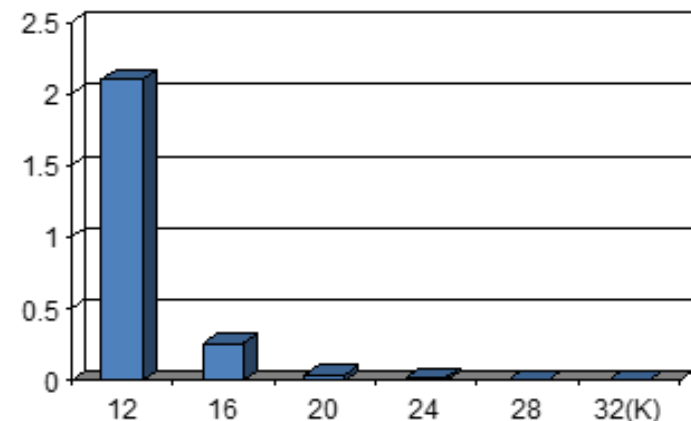


Clipped signal histogram



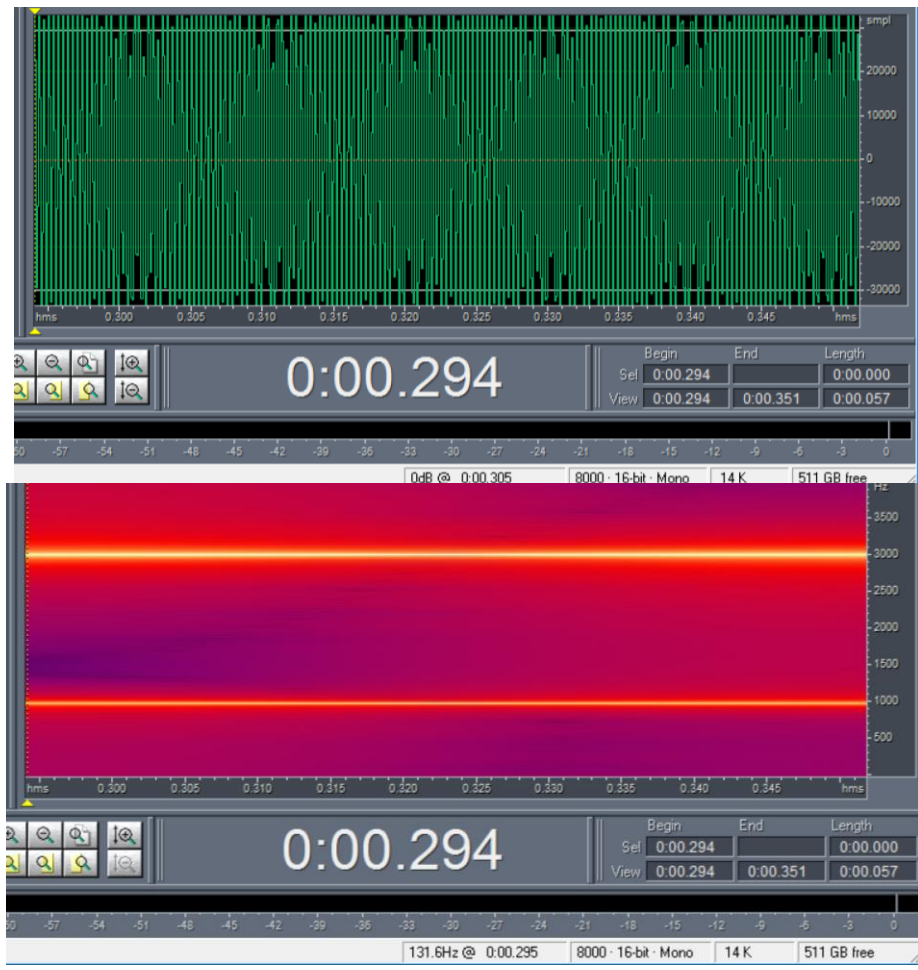
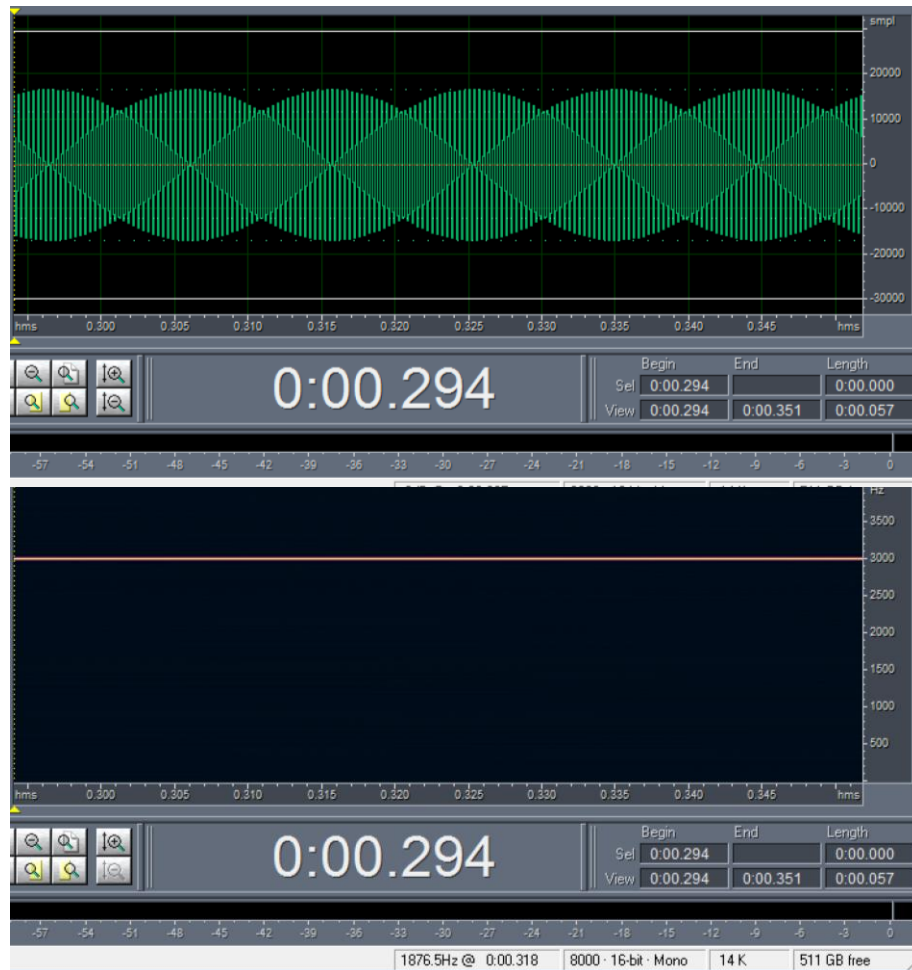
Absolute sample value

Normal signal histogram



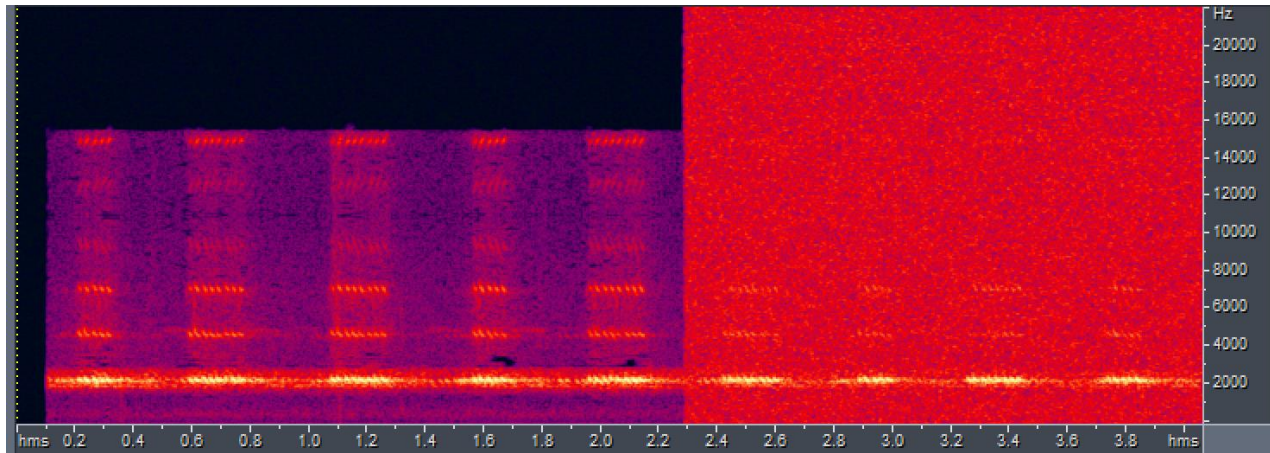
Absolute sample value

What does clipping do?



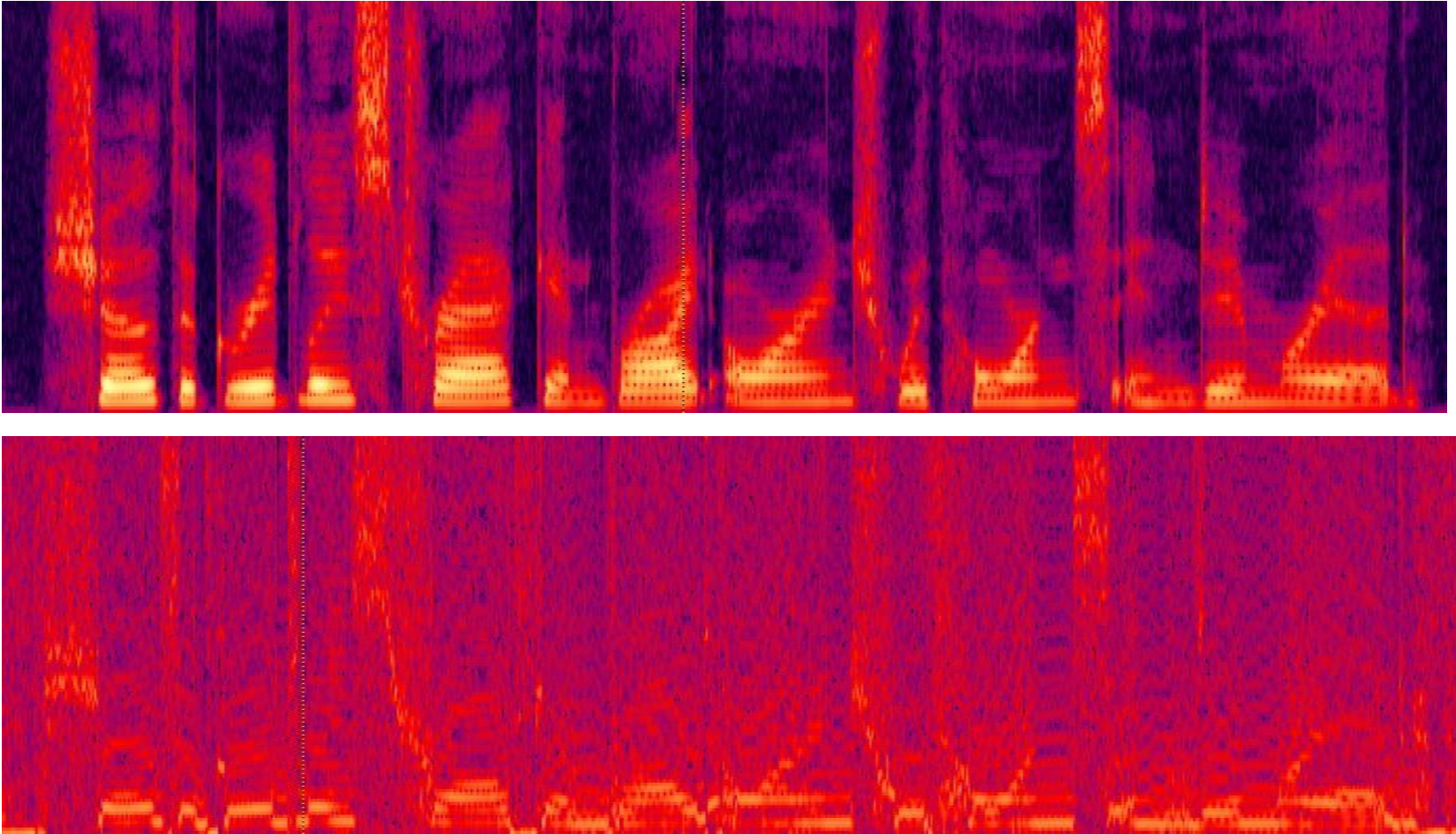
- It causes aliasing

Background Noise



- Signals in the background can corrupt the recording
 - Diffuse background noise, e.g. in an automobile
 - Localized sounds, e.g. air conditioner
 - Background talkers, music..
 - Denoising algorithms are used for cancelling noise, but they leave artifacts behind

Background Noise



- Clean and noisy speech

Recording Channel

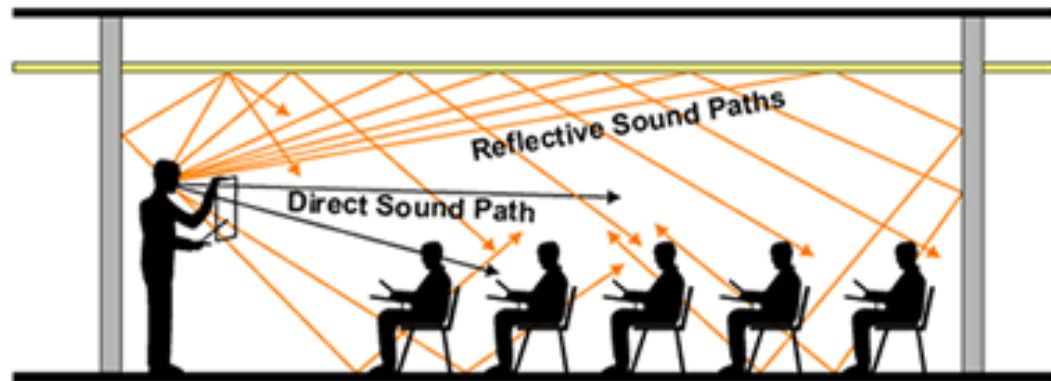
- Recording channels/media can result in signal distortions
 - Telephone channels reduce the bandwidth of the signal
 - 300-3300Hz
 - Cellphone / VOIP channels introduce *coding* distortion
 - Coding and decoding speech introduces distortions
 - They can also result in *discontinuities* from dropped packets
 - Poor microphones can result in spectral distortions
 - Computer recordings may be affected by memory and disk activity
- Distortions cannot usually be recovered from
- **Tampered signals may show anomalous patterns in distortion**

Reverberation

REVERBERATION TIME

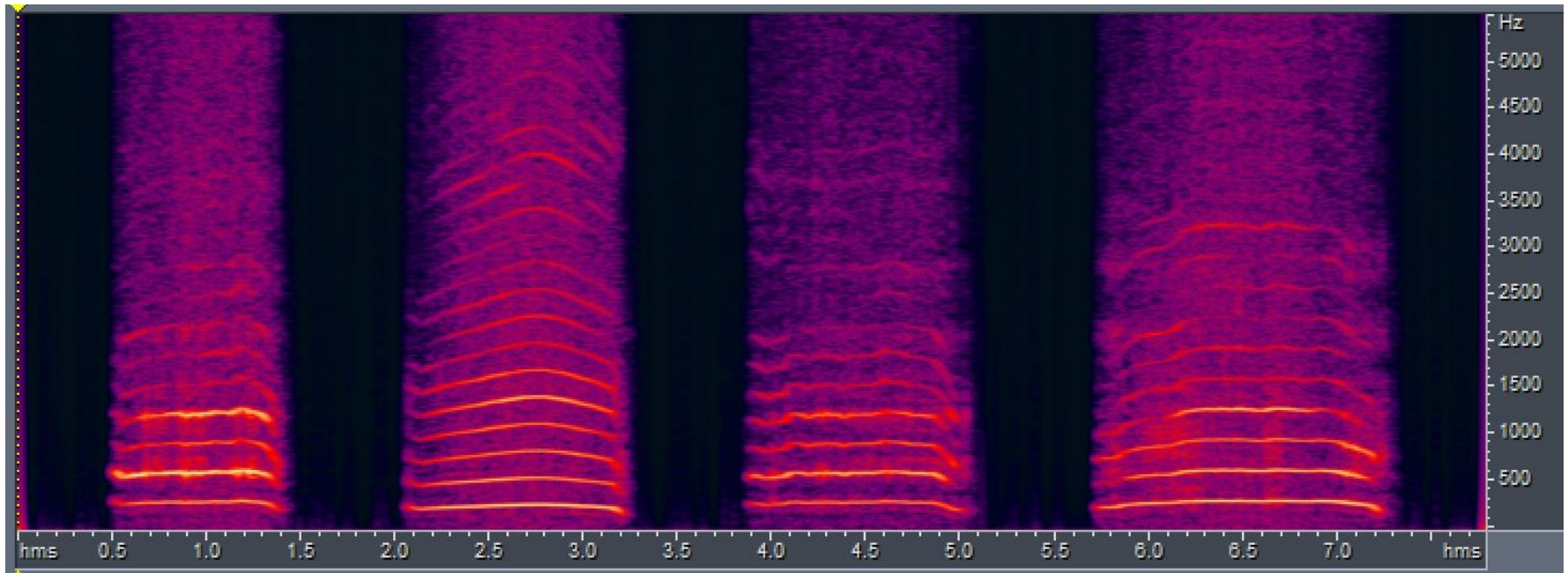
The time it takes for reflected sound to die down by 60 decibels from the cessation of the original sound signal (measured in seconds).

- Reflected sound tends to "build up" to a level louder than direct sound. Reflected sounds **MASK** direct sound.
- Late arriving reflections tend to **SMEAR** the direct sound signal.



- Each receiver "hears" the sound slightly differently

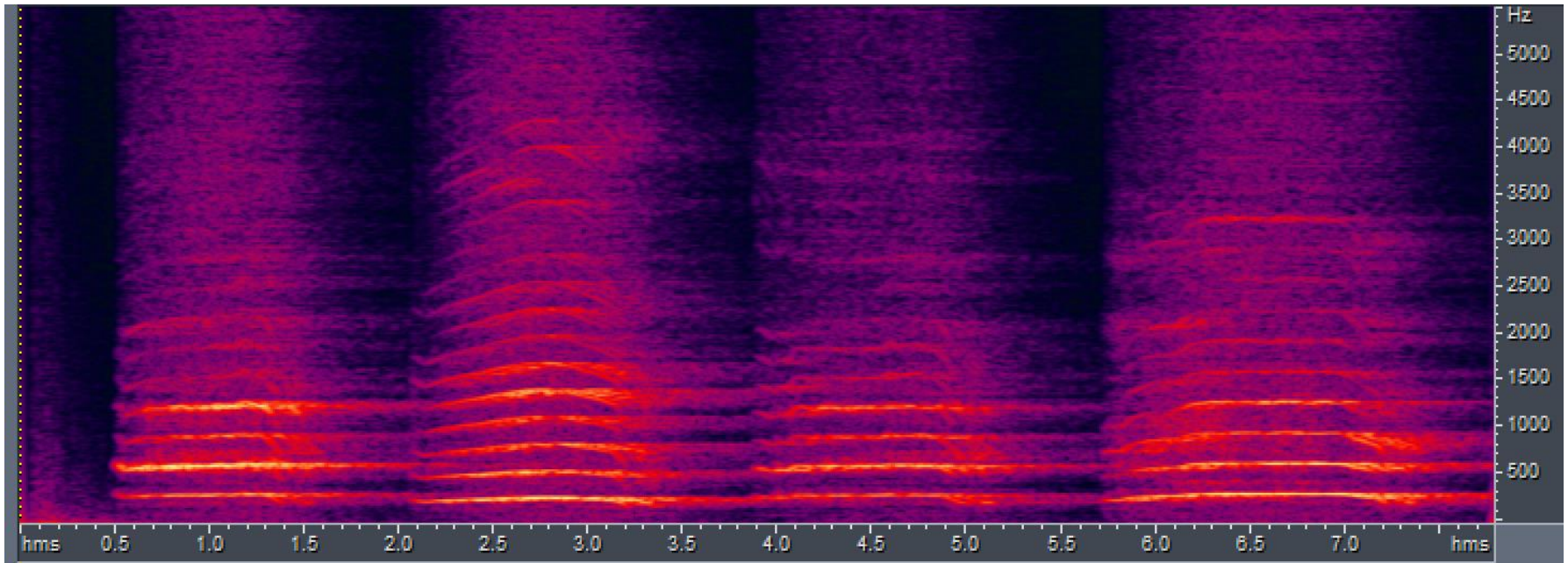
Reverberation example



- Lama



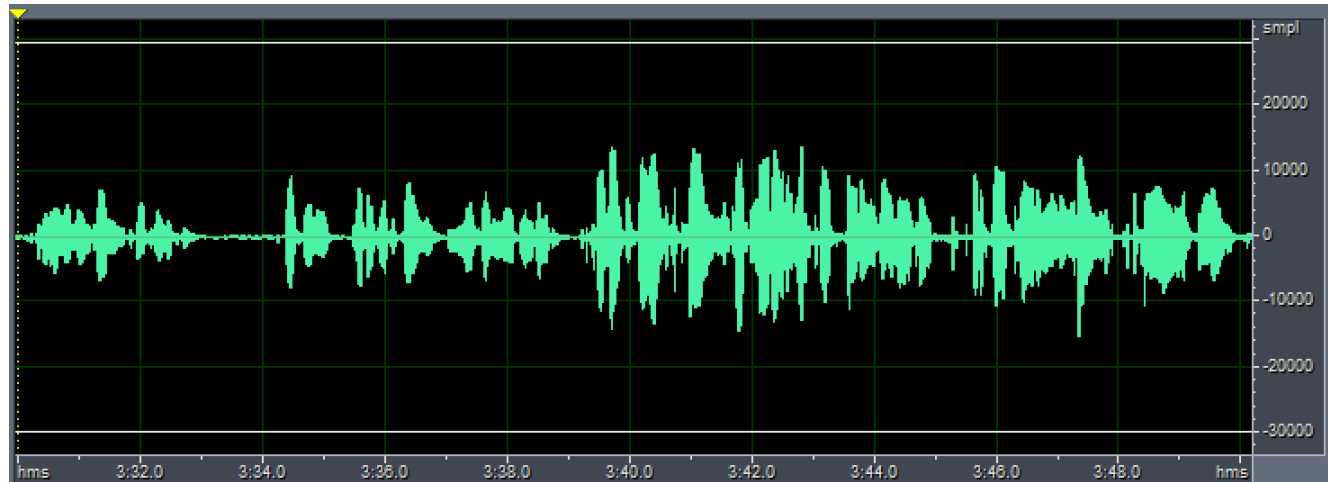
Reverberation example



- Lama, reverb in large hall

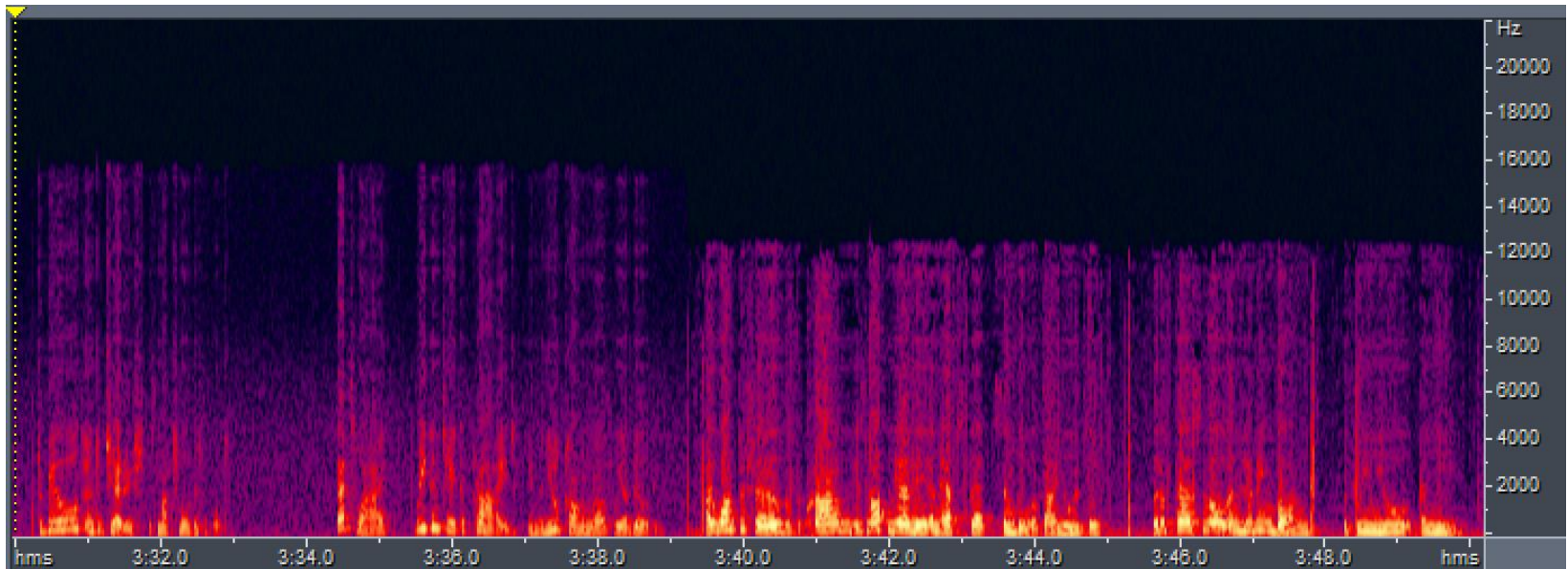


Mixing signals of different qualities: Splicing boundary effect



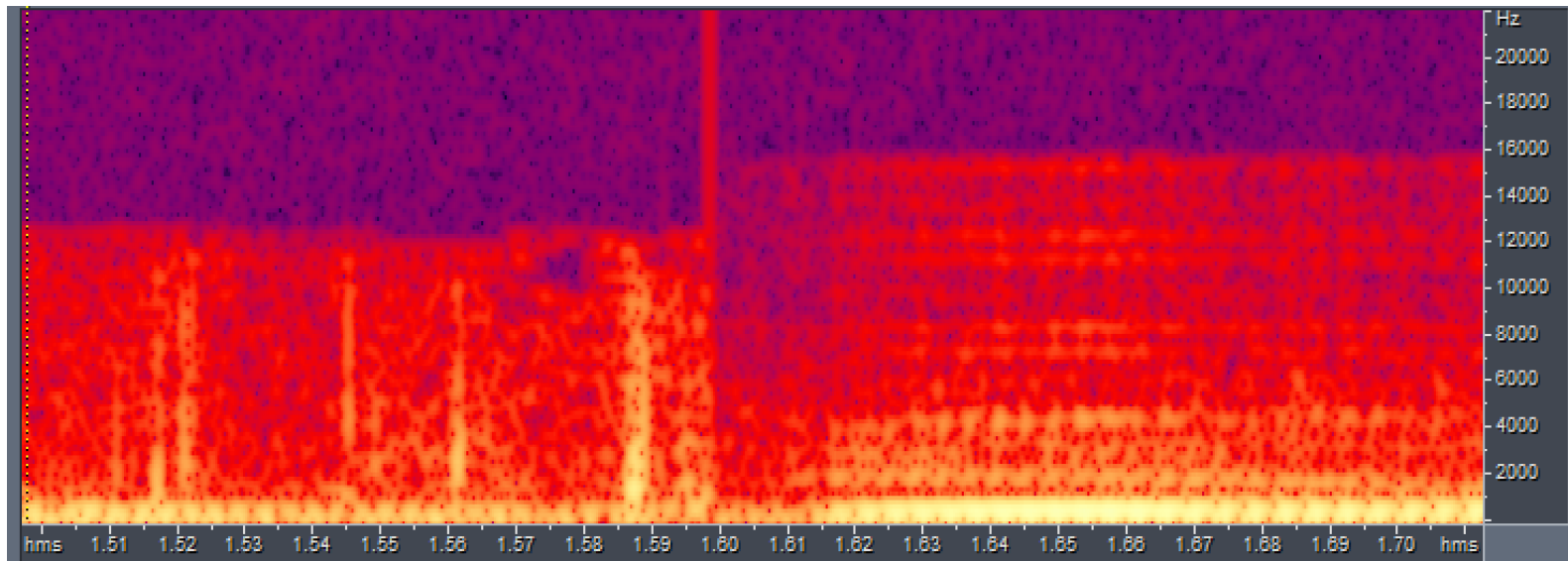
- The queen of England, 1953
- Tampered signal in time domain

Splicing boundary effect



- Snippets from two different speeches of the Queen of England in 1953

Splicing boundary effect



- The queen of England, 1953: “I want to sell the building”

In the next lecture

- Digital multimedia: Recording and devices
 - Audio
 - Images
 - Video
 - Text
- Digital multimedia: Processing
 - Audio processing
 - Some generic ML based processing techniques