

ECE160

Multimedia

Lecture 6: Spring 2011

Basics of Digital Audio

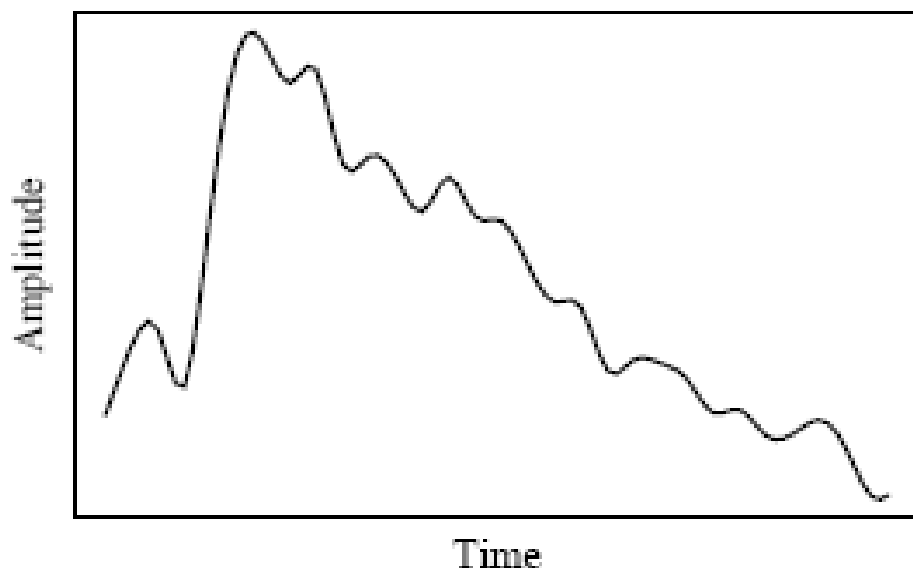
Digitization of Sound

What is Sound?

- Sound is a wave phenomenon like light, but is macroscopic and involves molecules of air being compressed and expanded under the action of some physical device.
 - (a) A speaker in an audio system moves back and forth and produces a *longitudinal* pressure wave that we perceive as sound.
 - (b) Since sound is a pressure wave, it takes on continuous analog values, as opposed to digitized ones.
- (c) Even though pressure waves are longitudinal, they still have ordinary wave properties and behaviors, i.e. reflection (bouncing), refraction (change of direction on entering a medium with a different density) and diffraction (bending around an obstacle).
 - (d) If we wish to use a digital version of sound waves we must form digitized representations of analog audio information.

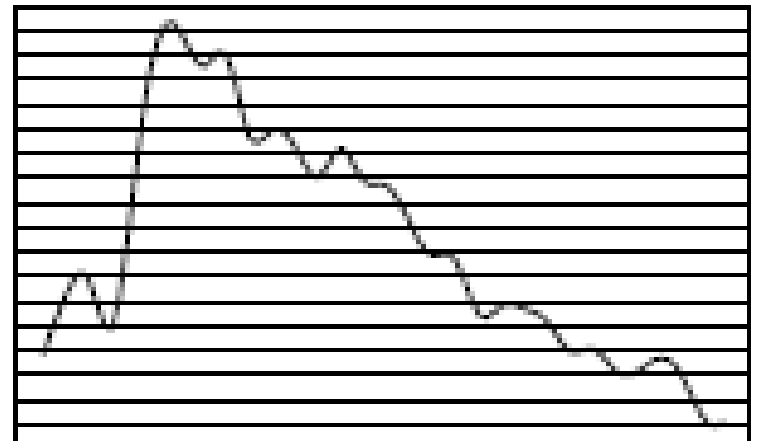
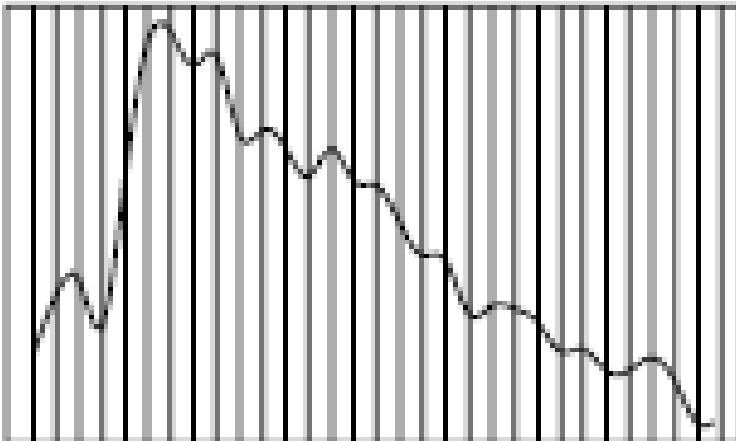
Digitization of Sound

- **Digitization** means conversion to a stream of numbers, and preferably these numbers should be integers for efficiency.
- The figure shows the 1-dimensional nature of sound: **amplitude** values depend on a 1D variable, time. (Images depend instead on a 2D set of variables, x and y).



Digitization of Sound

- The sound be made digital in both time and amplitude. To digitize, the signal is **sampled** in each dimension: in time, and in amplitude.
 - (a) Sampling means measuring the quantity in one dimension, usually at evenly-spaced intervals in the other dimension.
 - (b) The first kind of sampling, using measurements only at evenly spaced time intervals, is simply called, *sampling*. The rate at which it is performed is called the *sampling rate or frequency*.
 - (c) For audio, typical sampling rates are from 8 kHz (8,000 samples per second) to 48 kHz. This rate is determined by Nyquist theorem.
 - (d) Sampling in the amplitude dimension is called **quantization**.



Digitization of Sound

- Thus to decide how to digitize audio data we need to answer the following questions:
 1. What is the sampling rate?
 2. How finely is the data to be quantized, and is quantization uniform?
 3. How is audio data formatted?
(file format)

Nyquist Theorem

- Signals can be decomposed into a sum of sinusoids. The figure shows how weighted sinusoids can build up quite a complex signal.

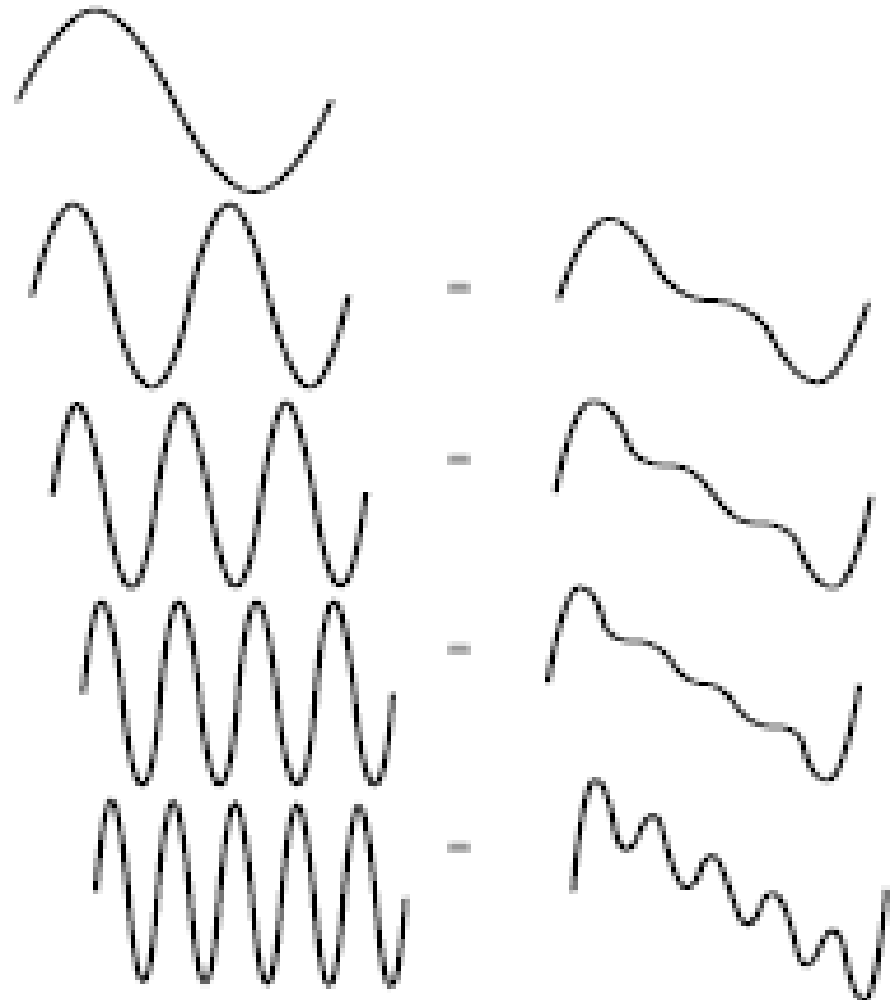
Fundamental
frequency

+ 0.5 ×
2 × fundamental

+ 0.33 ×
3 × fundamental

+ 0.25 ×
4 × fundamental

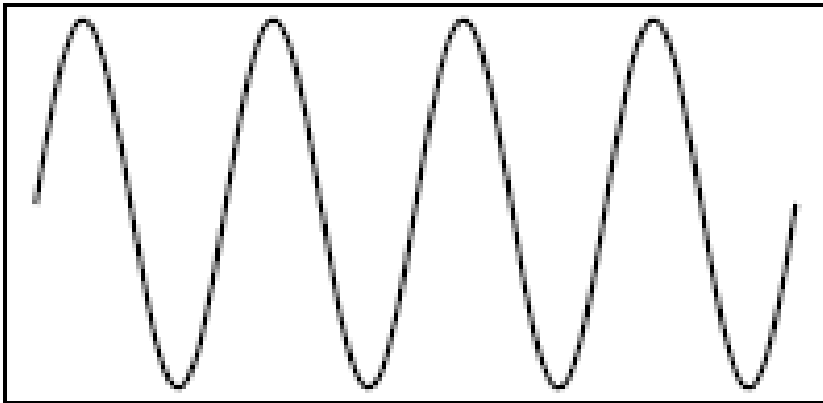
+ 0.5 ×
5 × fundamental



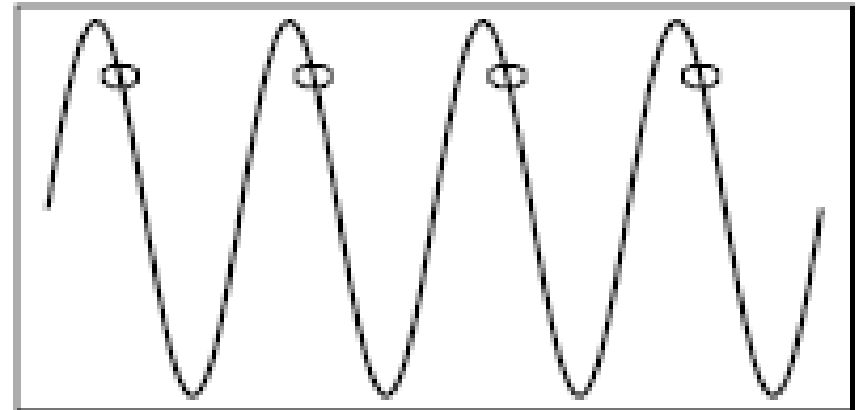
Nyquist Theorem

The Nyquist theorem states how frequently we must sample in time to be able to recover the original sound.

- (a) The figure shows a single sinusoid: it is a single, pure, frequency (only electronic instruments can create such sounds).
- (b) If sampling rate just equals the actual frequency, the figure shows that a false signal is detected: a constant, with zero frequency.



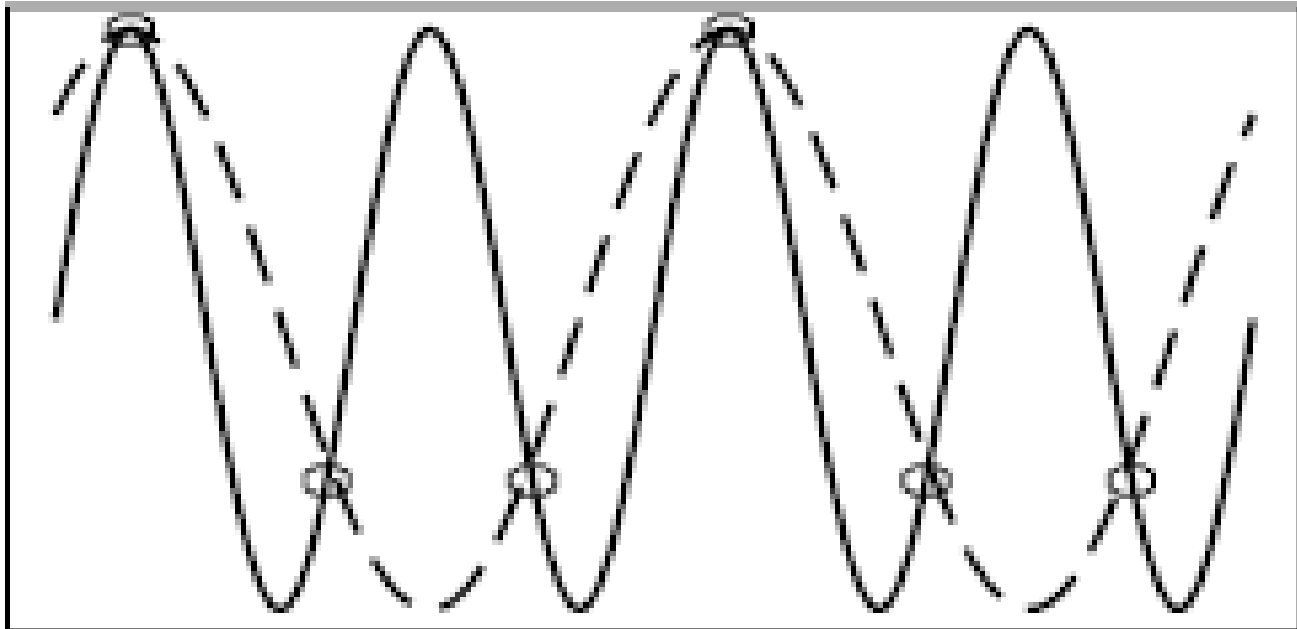
(a)



(b)

Nyquist Theorem

- (c) If we sample at 1.5 times the actual frequency, the figure shows that we obtain an incorrect (**alias**) frequency that is lower than the correct frequency - it is half the correct frequency.
- (d) For correct sampling we must use a sampling rate equal to at least *twice the maximum frequency* content in the signal. This rate is called the **Nyquist rate**.



Nyquist Theorem

- **Nyquist Theorem:** If a signal is **band-limited**, i.e., there is a lower limit f_1 and an upper limit f_2 of frequency components in the signal, then the sampling rate should be at least $2(f_2 - f_1)$.
- **Nyquist frequency:** half of the Nyquist rate.
 - Since it would be impossible to recover frequencies higher than Nyquist frequency in any event, most systems have an **antialiasing filter** that restricts the frequency content in the input to the sampler to a range at or below Nyquist frequency.
- The relationship among the Sampling Frequency, True Frequency, and the Alias Frequency is

$$f_{alias} = f_{sampling} - f_{true}, \quad \text{for} \quad f_{true} < f_{sampling} < 2 \times f_{true}$$

Pitch

- Whereas **frequency** is an absolute measure, **pitch** is generally relative - a perceptual subjective quality of sound.
 - (a) Pitch and frequency are linked by setting the note A above middle C to exactly 440 Hz.
 - (b) An **octave** above that note takes us to another A note. An octave corresponds to *doubling the frequency*. Thus with the middle "A" on a piano ("A4" or "A440") set to 440 Hz, the next "A" up is at 880 Hz, or one octave above.
 - (c) **Harmonics**: any series of musical tones whose frequencies are integral multiples of the frequency of a fundamental tone.
 - (d) If we allow non-integer multiples of the base frequency, we allow non-"A" notes and have a more complex resulting sound.

Signal to Noise Ratio (SNR)

- The ratio of the power of the correct signal and the noise is called the *signal to noise ratio* (**SNR**) - a measure of the quality of the signal.
- The SNR is usually measured in decibels (**dB**), where 1 dB is a tenth of a **bel**. The SNR value, in units of dB, is defined in terms of base-10 logarithms of squared voltages:

$$SNR = 10 \log_{10} \frac{V_{signal}^2}{V_{noise}^2} = 20 \log_{10} \frac{V_{signal}}{V_{noise}}$$

Signal to Noise Ratio (SNR)

- a) The power in a signal is proportional to the square of the voltage. For example, if the signal voltage V_{signal} is 10 times the noise, then the SNR is $20 \log_{10}(10)=20\text{dB}$.
- b) In terms of power, if the power from ten violins is ten times that from one violin playing, then the ratio of power is 10dB, or 1B.

Sound Levels

The usual levels of sound we hear around us are described in terms of decibels, as a ratio to the quietest sound we are capable of hearing. The table shows approximate levels for these sounds.

Threshold of hearing	0
Rustle of leaves	10
Very quiet room	20
Average room	40
Conversation	60
Busy street	70
Loud radio	80
Train through station	90
Riveter	100
Threshold of discomfort	120
Threshold of pain	140
Damage to ear drum	160

Signal to Quantization Noise Ratio (SQNR)

- Aside from any noise present in the original analog signal, there is also an additional error that results from quantization.
 - (a) If voltages are 0 to 1 but we have only 8 bits to store values, then we force **all** continuous values into only 256 different values.
 - (b) This introduces a roundoff error. It is not really “noise”. Nevertheless it is called **quantization noise** (or quantization error).
- The quality of the quantization is characterized by the Signal to Quantization Noise Ratio (**SQNR**).
 - (a) **Quantization noise**: the difference between the actual value of the analog signal, for the particular sampling time, and the nearest quantization interval value.
 - (b) At most, this error can be as much as half of the interval.

Signal to Quantization Noise Ratio (SQNR)

(c) For a quantization accuracy of N bits per sample, the SQNR can be simply expressed:

$$\begin{aligned} SQNR &= 20 \log_{10} \frac{V_{\text{signal}}}{V_{\text{quan_noise}}} = 20 \log_{10} \frac{2^{N-1}}{\frac{1}{2}} \\ &= 20 \times N \times \log 2 = 6.02N \text{ (dB)} \end{aligned}$$

Notes:

(a) We map the maximum signal to $2^{N-1} - 1$ ($\sim 2^{N-1}$) and the most negative signal to -2^{N-1} .

(b) The equation is the *Peak* signal-to-noise ratio, PSQNR: peak signal and peak noise.

Signal to Quantization Noise Ratio (SQNR)

- (c) The *dynamic range* is the ratio of maximum to minimum absolute values of the signal: V_{max}/V_{min} . The max abs. value V_{max} gets mapped to $2^{N-1}-1$; the min abs. value V_{min} gets mapped to 1. V_{min} is the smallest positive voltage that is not masked by noise. The most negative signal, $-V_{max}$, is mapped to -2^{N-1} .
- (d) The quantization interval is $V=(2V_{max})/2N$, since there are $2N$ intervals. The whole range V_{max} down to $(V_{max}-V/2)$ is mapped to $2N-1 - 1$.
- (e) The maximum noise, in terms of actual voltages, is half the quantization interval: $V/2=V_{max}/2N$.
- (f) **6.02N is the worst case.** If the input signal is sinusoidal, the quantization error is statistically independent, and its magnitude is uniformly distributed between 0 and half of the interval, then it can be shown that the expression for the SQNR becomes:

$$SQNR = 6.02N + 1.76(dB)$$

Linear and Non-linear Quantization

- **Linear format:** samples are typically stored as uniformly quantized values.
- **Non-uniform quantization:** set up more finely-spaced levels where humans hear with the most acuity.
 - Weber's Law stated formally says that equally perceived differences have values proportional to absolute levels:

$$\Delta \text{Response} \propto \Delta \text{stimulus} / \text{Stimulus}$$

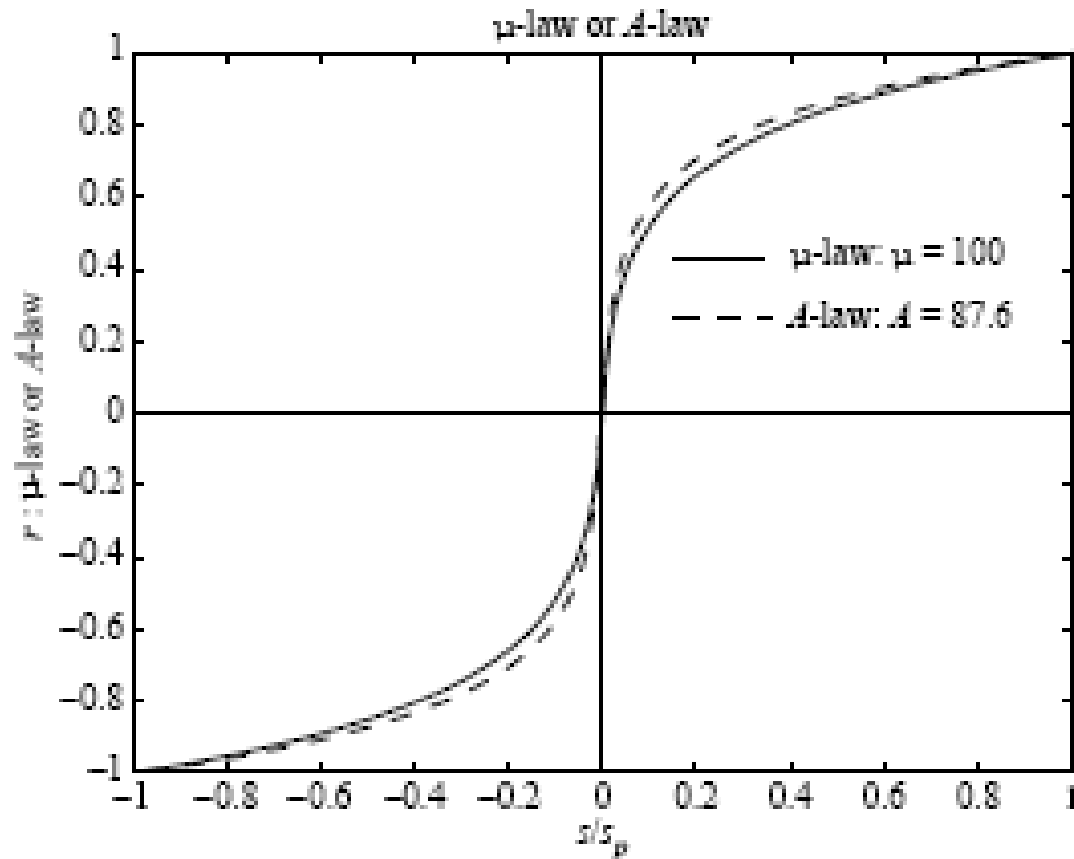
- Inserting a constant of proportionality k , we have a differential equation that states:

$$dr = k(1/s) ds$$

with response r and stimulus s .

Non-linear Quantization

- The figure shows these curves. The parameter μ is set to $\mu = 100$ or $\mu = 255$; the parameter A for the A -law encoder is set to $A = 87.6$.



Audio Filtering

- Prior to sampling and AD conversion, the audio signal is also usually *filtered* to remove unwanted frequencies. The frequencies kept depend on the application:
 - (a) For speech, typically from 50Hz to 10kHz is retained, and other frequencies are blocked by a **band-pass filter** that screens out lower and higher frequencies.
 - (b) An audio music signal will typically contain from about 20Hz up to 20kHz.
 - (c) At the DA converter end, high frequencies may reappear in the output - because of sampling and then quantization - smooth input signal is replaced by a series of step functions containing all possible frequencies.
 - (d) At the decoder side, after the DA converter, a **lowpass filter** is used to remove those steps.

Audio Quality vs. Data Rate

- The uncompressed data rate increases as more bits are used for quantization. Stereo: double the bandwidth. to transmit a digital audio signal.

Quality	Sample Rate (KHz)	Bits per Sample	Mono/ Stereo	Data Rate (uncompressed) (kB/sec)	Frequency Band (KHz)
Telephone	8	8	Mono	8	0.200-3.4
AM Radio	11.025	8	Mono	11.0	0.1-5.5
FM Radio	22.05	16	Stereo	88.2	0.02-11
CD	44.1	16	Stereo	176.4	0.005-20
DAT	48	16	Stereo	192.0	0.005-20
DVD Audio	192 (max)	24 (max)	6 channels	1,200.0 (max)	0-96 (max)

MIDI: Musical Instrument Digital Interface

MIDI Overview

(a) MIDI is a scripting language with “events” for the production of sounds, including values for the pitch of a note, its duration, and its volume.

(b) MIDI is a standard adopted by electronic music for controlling devices, such as synthesizers and sound cards, that produce music.

(c) MIDI standard is supported by most synthesizers, so sounds created on one synthesizer can be played and manipulated on another synthesizer and sound reasonably close.

(d) Computers have a MIDI interface incorporated into most sound cards, with both D/A and A/D converters.

MIDI Concepts

- **MIDI channels** are used to separate messages.
 - (a) There are 16 channels numbered from 0 to 15. The channel forms the last 4 bits (the least significant bits) of the message.
 - (b) Usually a channel is associated with a particular instrument: e.g., channel 1 is the piano, channel 10 is the drums, etc.
 - (c) Nevertheless, one can switch instruments midstream, if desired, and associate another instrument with any channel.

MIDI Concepts

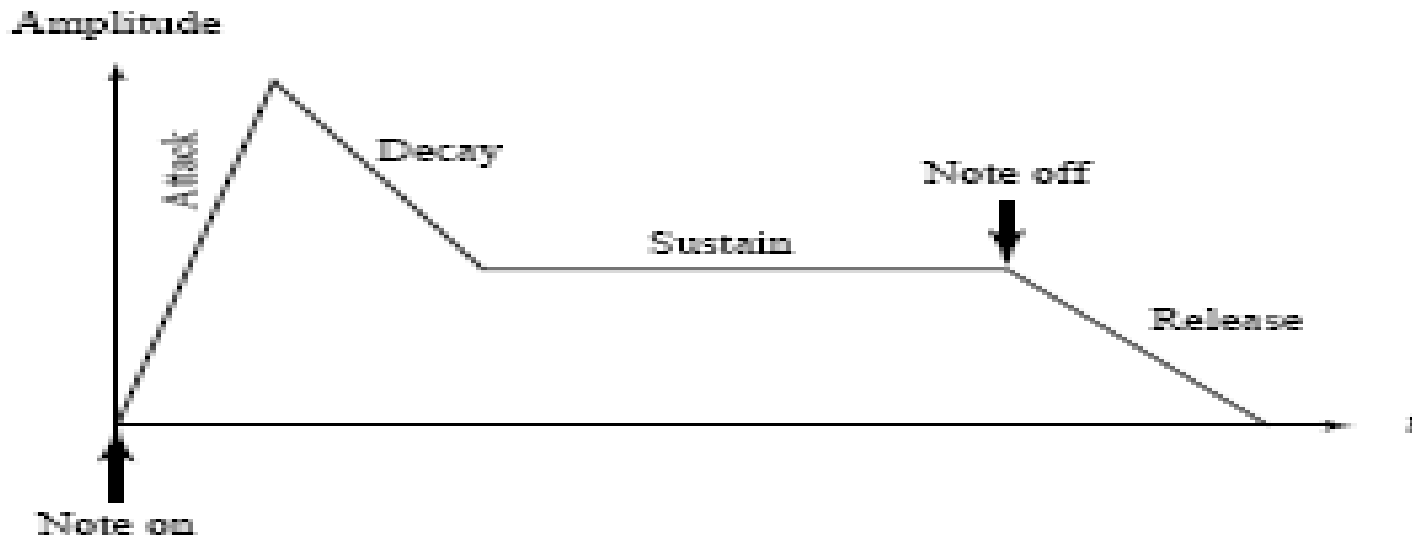
- **System messages**
 - (a) Several other types of messages, e.g. a general message for all instruments indicating a change in tuning or timing.
 - (b) If the first 4 bits are all 1s, then the message is interpreted as a **system common** message.
- The way a synthetic musical instrument responds to a MIDI message is usually by simply ignoring any **play sound** message that is not for its channel.
 - If several messages are for its channel, then the instrument responds, provided it is **multi-voice**, i.e., can play more than a single note at once.

MIDI Concepts

- It is easy to confuse the term **voice** with the term **timbre** - the latter is MIDI terminology for just what instrument that is trying to be emulated, e.g. a piano as opposed to a violin: it is the quality of the sound.
 - (a) An instrument (or sound card) that is **multi-timbral** is one that is capable of playing many different sounds at the same time, e.g., piano, brass, drums, etc.
 - (b) On the other hand, the term **voice**, while sometimes used by musicians to mean the same thing as timbre, is used in MIDI to mean every different timbre and pitch that the tone module can produce at the same time.
- Different timbres are produced digitally by using a **patch** - the set of control settings that define a particular timbre.
- Patches are often organized into databases, called **banks**.

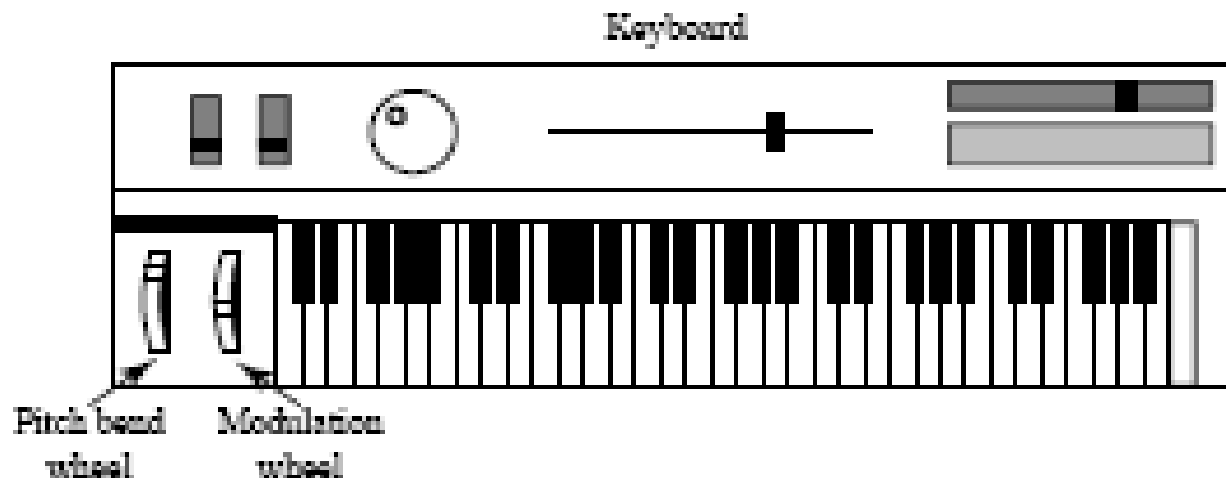
MIDI Concepts

- A MIDI device is capable of **programmability**, and can change the **envelope** describing how the amplitude of a sound changes over time.
- The figure shows a model of the response of a digital instrument to a Note On message:



Hardware Aspects of MIDI

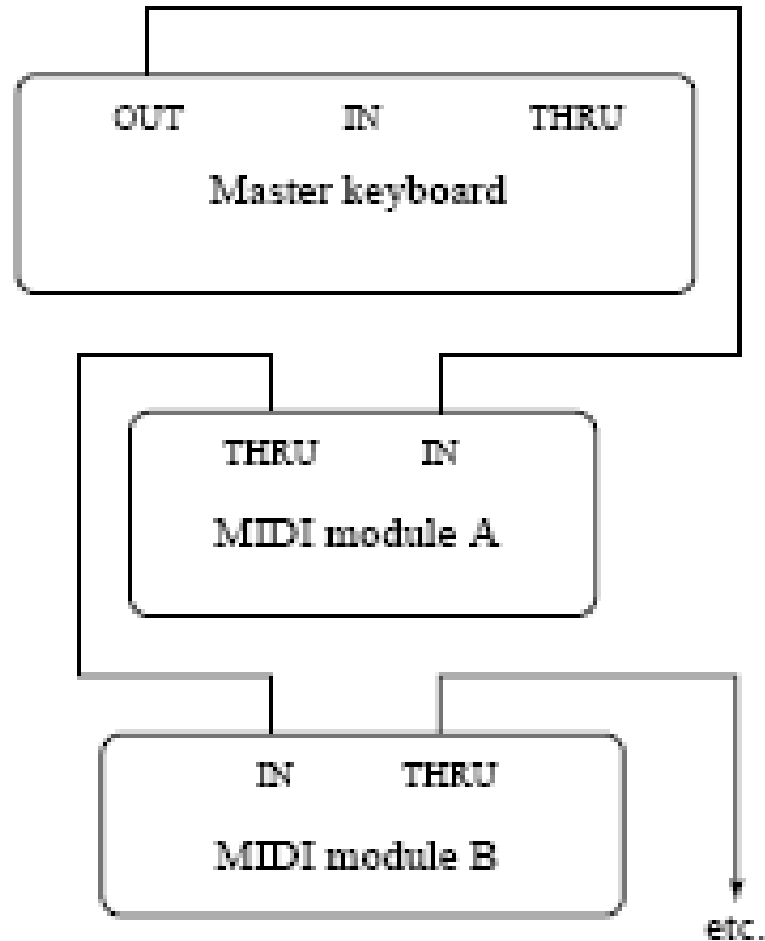
- The MIDI hardware setup consists of a 31.25 kbps serial connection. Usually, MIDI-capable units are either Input devices or Output devices, not both.
- A traditional synthesizer is shown



Hardware Aspects of MIDI

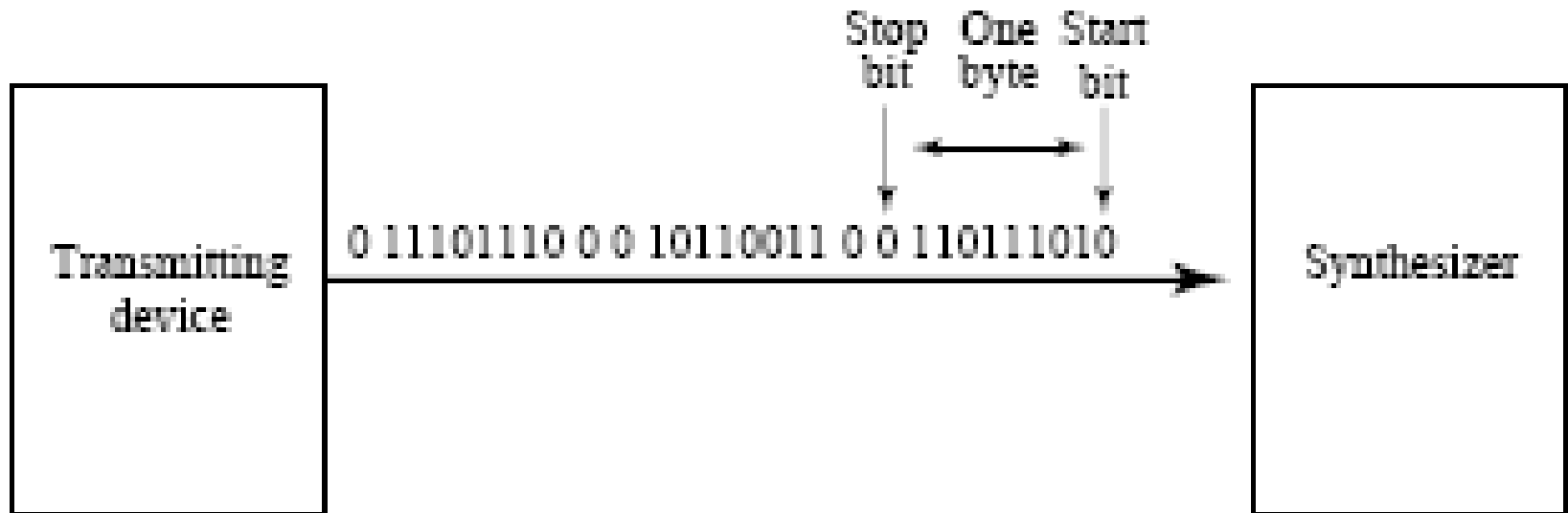
- The physical MIDI ports consist of 5-pin connectors for IN and OUT, as well as a third connector called THRU.
 - (a) MIDI communication is half-duplex.
 - (b) MIDI IN is the connector via which the device receives all MIDI data.
 - (c) MIDI OUT is the connector through which the device transmits all the MIDI data it generates itself.
 - (d) MIDI THRU is the connector by which the device echoes the data it receives from MIDI IN. Note that it is only the MIDI IN data that is echoed by MIDI THRU - all the data generated by the device itself is sent via MIDI OUT.

A typical MIDI sequencer setup



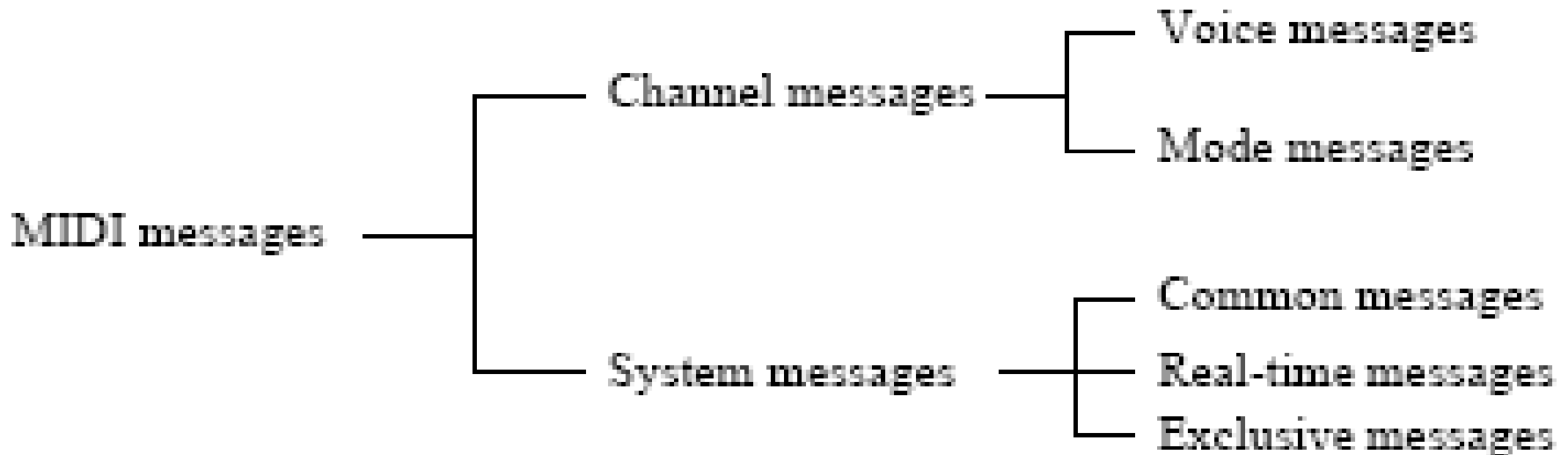
MIDI data stream

- A stream of 10-bit bytes for MIDI messages consist of {Status byte, Data Byte, Data Byte} = {Note On, Note Number, Note Velocity}.
MIDI bytes are 10-bit, with a 0 start and 0 stop bit.



Structure of MIDI Messages

- MIDI messages can be classified into two types: channel messages and system messages



MIDI Voice Messages

Voice messages:

- a) This type of channel message controls a voice, i.e., sends information specifying which note to play or to turn off, and encodes key pressure.
- b) Voice messages are also used to specify controller effects such as sustain, vibrato, tremolo, and the pitch wheel.

Voice Message	Status Byte	Data Byte1	Data Byte2
Note Off	&H8n	Key number	Note Off velocity
Note On	&H9n	Key number	Note On velocity
Poly. Key Pressure	&HAN	Key number	Amount
Control Change	&HBn	Controller num.	Controller value
Program Change	&HCn	Program number	None
Channel Pressure	&HDn	Pressure value	None
Pitch Bend	&HEn	MSB	LSB

General MIDI

- General MIDI is a scheme for standardizing the assignment of instruments to patch numbers.
 - a) A standard percussion map specifies 47 percussion sounds.
 - b) Where a “note” appears on a musical score determines what percussion instrument is being struck: a bongo drum, a cymbal.
 - c) Other requirements for General MIDI compatibility: MIDI device must support all 16 channels; a device must be multitimbral (i.e., each channel can play a different instrument/program); a device must be polyphonic (i.e., each channel is able to play many voices); and there must be a minimum of 24 dynamically allocated voices.
- **General MIDI Level2:** An extended general MIDI has recently been defined, with a standard .smf “Standard MIDI File” format defined - inclusion of extra character information, such as karaoke lyrics.

MIDI to WAV Conversion

- Some programs, such as early versions of Premiere, cannot include .mid files - instead, they insist on .wav format files.
 - a) Various shareware programs exist for approximating a reasonable conversion between MIDI and WAV formats.

Csound

- Csound is a unit generator-based, user-programmable computer music system, written by Barry Vercoe at MIT in 1984. Since then Csound has received numerous contributions from researchers and musicians from around the world.
- Csound runs on many varieties of UNIX and Linux, Microsoft DOS and Windows, all versions of the Macintosh operating system including Mac OS X, and others.
- Csound can be considered one of the most powerful musical instruments ever created.
- To make music with Csound:
 1. Write an **orchestra** (.orc file) that creates instruments and signal processors by connecting unit generators (also called opcodes, in Csound-speak) using Csound's simple programming language.
 2. Write a **score** (.sco file) that specifies a list of notes and other events to be rendered by the orchestra.
 3. Run Csound to compile the orchestra and score, run the sorted and preprocessed score through the orchestra.

pluck -- Produces a naturally decaying plucked string or drum sound.

pluck *kamp*, *kcps*, *icps*, *ifn*, *imeth* [, *iparm1*] [, *iparm2*]

- *kamp* -- the output amplitude.
- *kcps* -- the resampling frequency in cycles-per-second. An audio buffer, filled at *i*-time according to *ifn*, is sampled with periodicity *kcps* and multiplied by *kamp*. The buffer is smoothed to simulate the effect of natural decay.
- *icps* -- intended pitch value in Hz, sets up a buffer of audio samples smoothed by a chosen decay.
- *ifn* -- number of a function used to initialize the cyclic decay buffer. If *ifn* = 0, a random sequence will be used.
- *imeth* -- method of natural decay. There are six, some of which use parameter values that follow.
 1. Simple averaging. A simple smoothing process, uninfluenced by parameter values.
 2. Stretched averaging. As above, with smoothing time stretched by a factor of *iparm1* (=1).
 3. Simple drum. The range from pitch to noise is controlled by a 'roughness factor' in *iparm1* (0 to 1). Zero gives the plucked string effect, while 1 reverses the polarity of every sample (octave down, odd harmonics). The setting .5 gives an optimum snare drum.
 4. Stretched drum. Combines both roughness and stretch factors. *iparm1* is roughness (0 to 1), and *iparm2* the stretch factor (=1).
 5. Weighted averaging. As method 1, with *iparm1* weighting the current sample (the status quo) and *iparm2* weighting the previous adjacent one. *iparm1* + *iparm2* must be ≤ 1 .
 6. 1st order recursive filter, with coefs .5. Unaffected by parameter values.
- *iparm1*, *iparm2* (optional) -- parameter values for use by the smoothing algorithms (above).
- Plucked strings (1,2,5,6) are best realized by starting with a random noise source, which is rich in initial harmonics. Drum sounds (methods 3,4) work best with a flat source (wide pulse), which produces a deep noise attack and sharp decay.

streson -- A string resonator with variable fundamental frequency.

- **streson** *asig*, *kfr*, *ifdbgain*
- *asig* -- the input audio signal.
- *kfr* -- the fundamental frequency of the string.
- *streson* passes the input *asig* through a network composed of comb, low-pass and all-pass filters, similar to the one used in the Karplus-Strong algorithm, creating a string resonator effect. The fundamental frequency of the “string” is controlled by *kfr*. This opcode is used to simulate sympathetic resonances to an input signal.
- *ifdbgain* -- feedback gain, between 0 and 1, of the internal delay line. A value close to 1 creates a slower decay and a more pronounced resonance. Small values may leave the input signal unaffected. Depending on the filter frequency, typical values are $> .9$

Csound Score

Cscore is a program for generating and manipulating numeric score files. It comprises function subprograms, called by a user-written control program, invoked either as a stand alone score preprocessor, or as part of the Csound run-time system:

A **Score** (a collection of score statements) is divided into time-ordered sections by the *s statement*. Before being read by the orchestra, a score is preprocessed one section at a time. Each section is processed by 3 routines: *Carry*, *Tempo*, and *Sort*.

Carry

- Determines how long a note or sequence of notes should last.

Tempo

- Time warps a score section according to the information in a *t statement*. The tempo operation converts p2 (and, for *i statements*, p3) from original beats into real seconds, since those are the units required by the orchestra.

Sort

- This routine sorts all action-time statements into chronological order by p2 value.