

PAT 463/563 (Fall 2025)

Music & AI

Lecture 5: Audio Processing Fundamentals

Instructor: Hao-Wen Dong

Four Representative Music Representations



Symbolic music representations

Text-based

```
Program_change_0,  
Note_on_60, Time_shift_2, Note_off_60,  
Note_on_60, Time_shift_2, Note_off_60,  
Note_on_76, Time_shift_2, Note_off_67,  
Note_on_67, Time_shift_2, Note_off_67,  
...
```

MIDI

Image-based



Piano roll



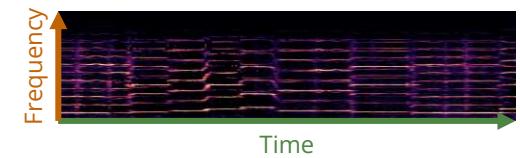
Audio-domain music representations

Time series-based



Waveform

Image-based



Spectrogram

Today's topic!

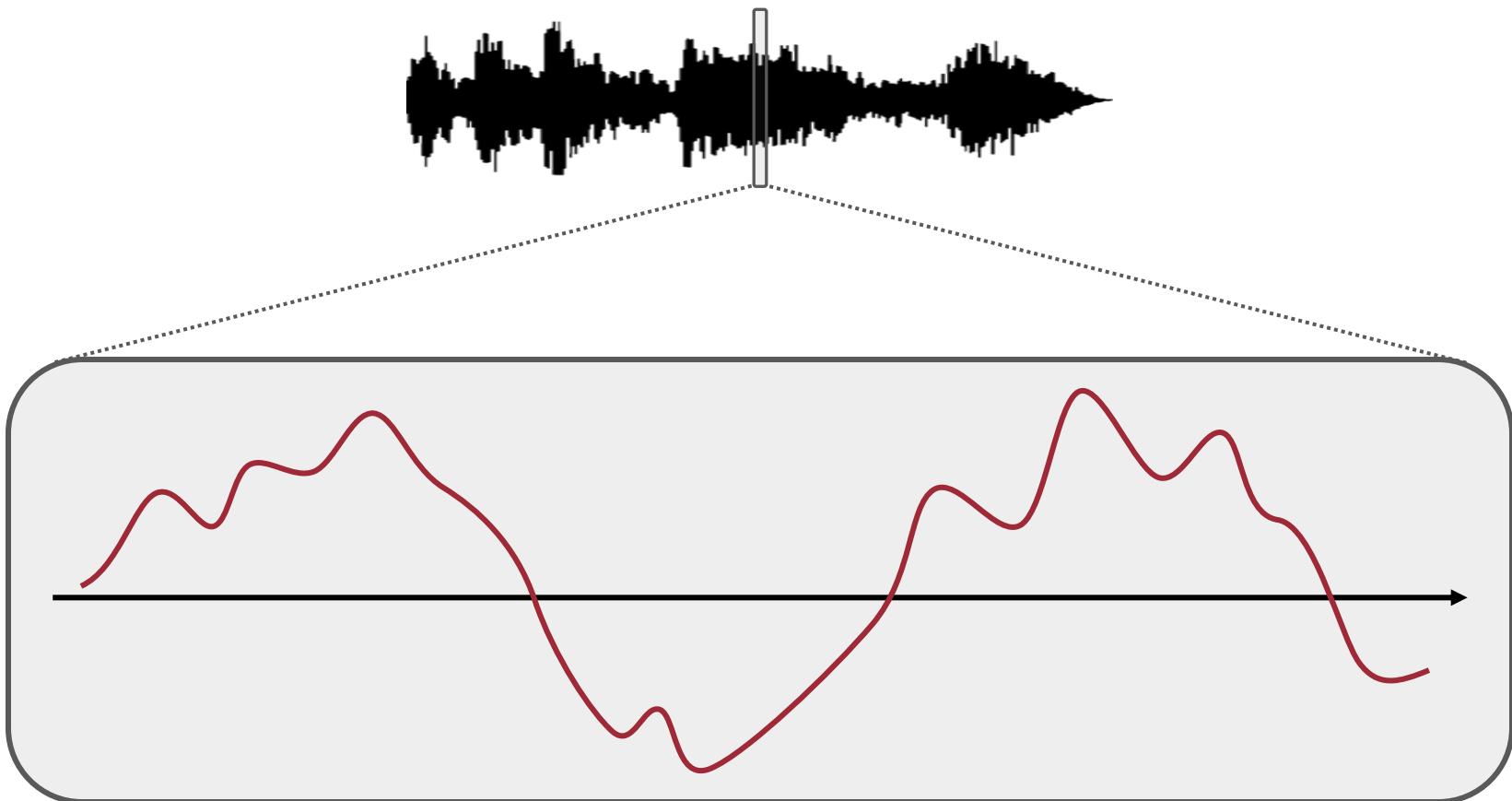
Digital Audio

Digital Audio

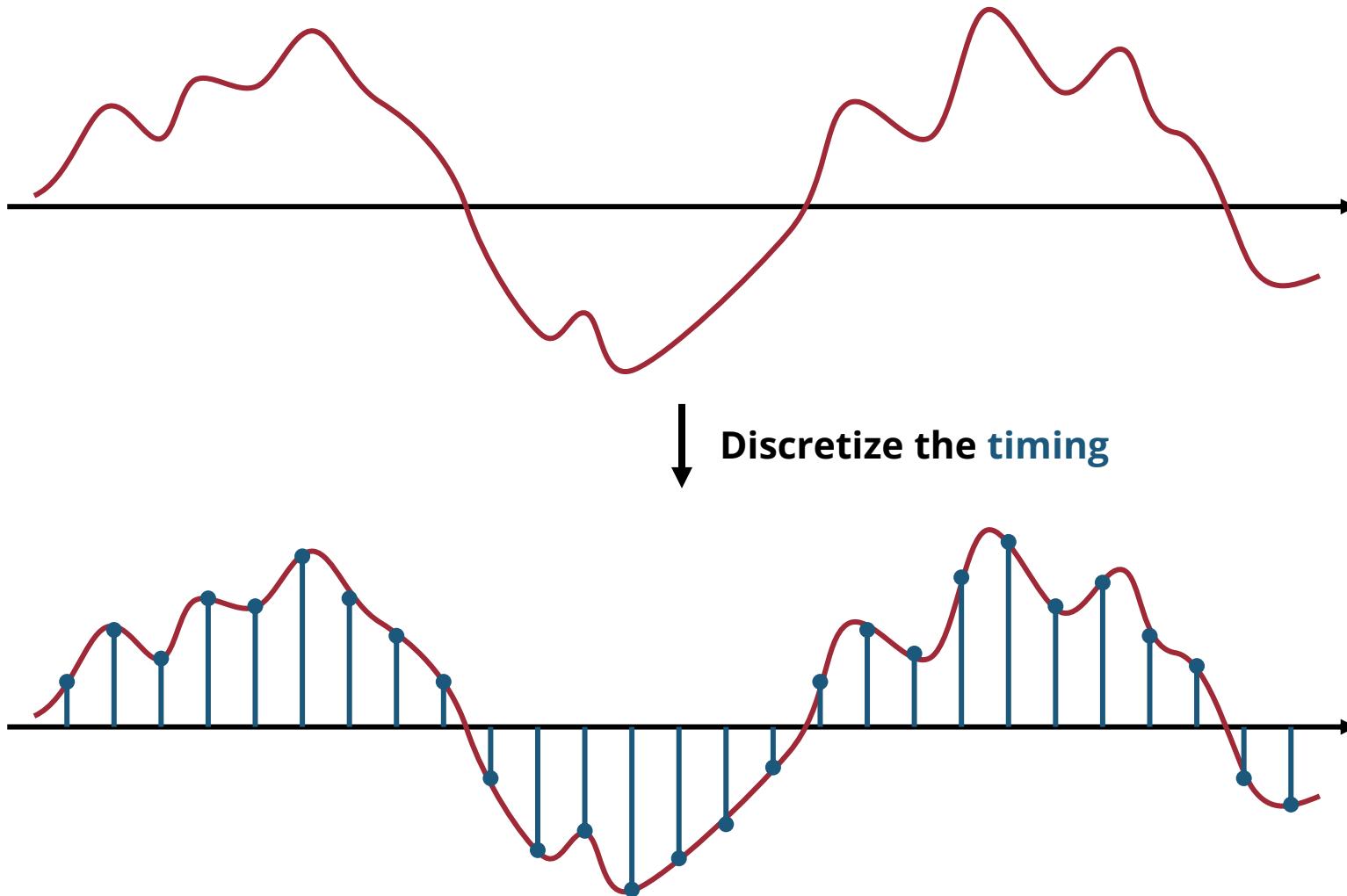


(Source: van den Oord et al., 2016)

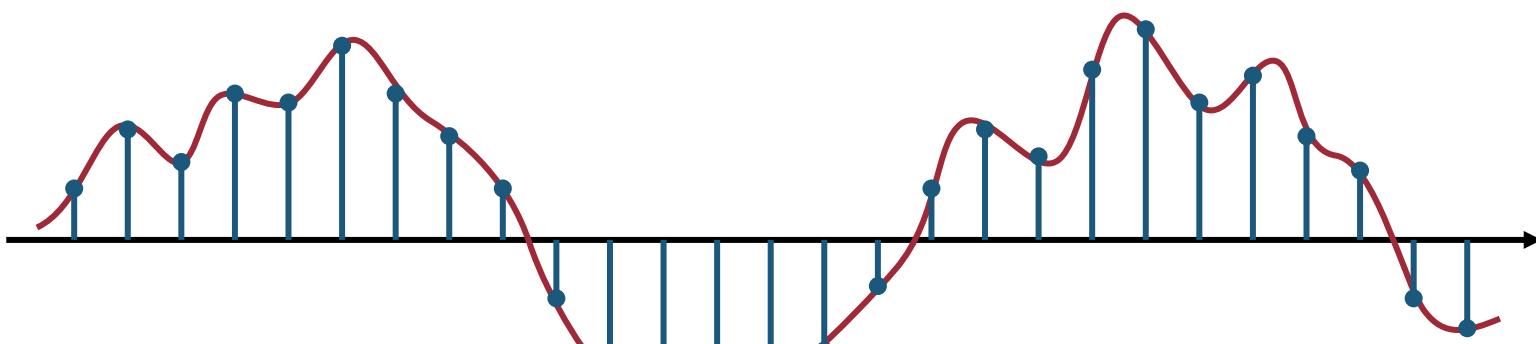
| Waveform



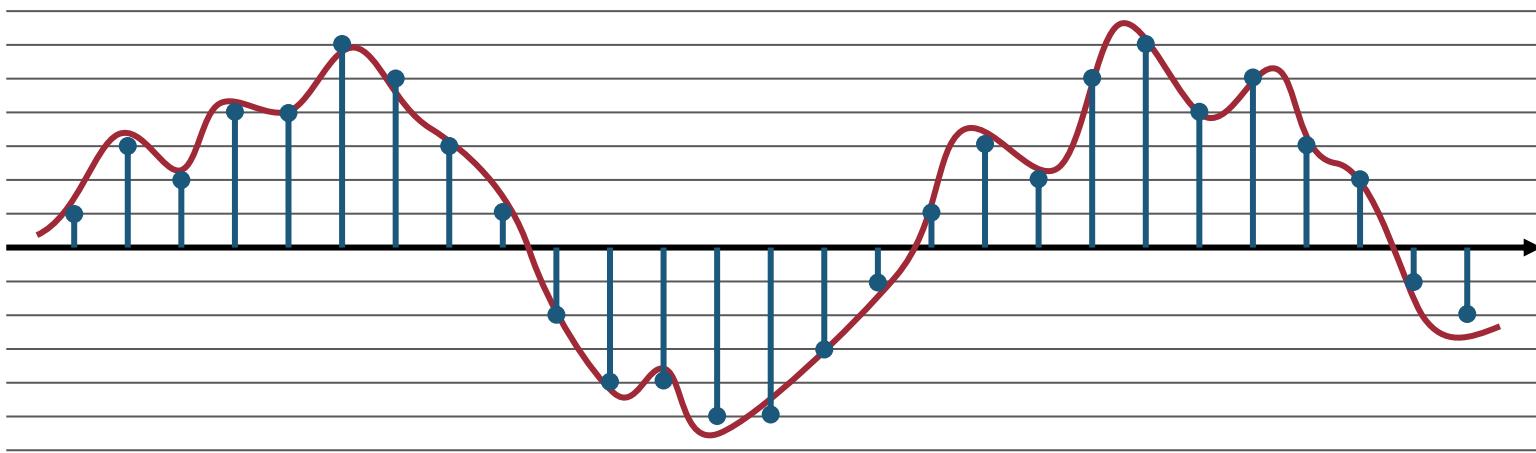
Digitalizing Audio: Timing



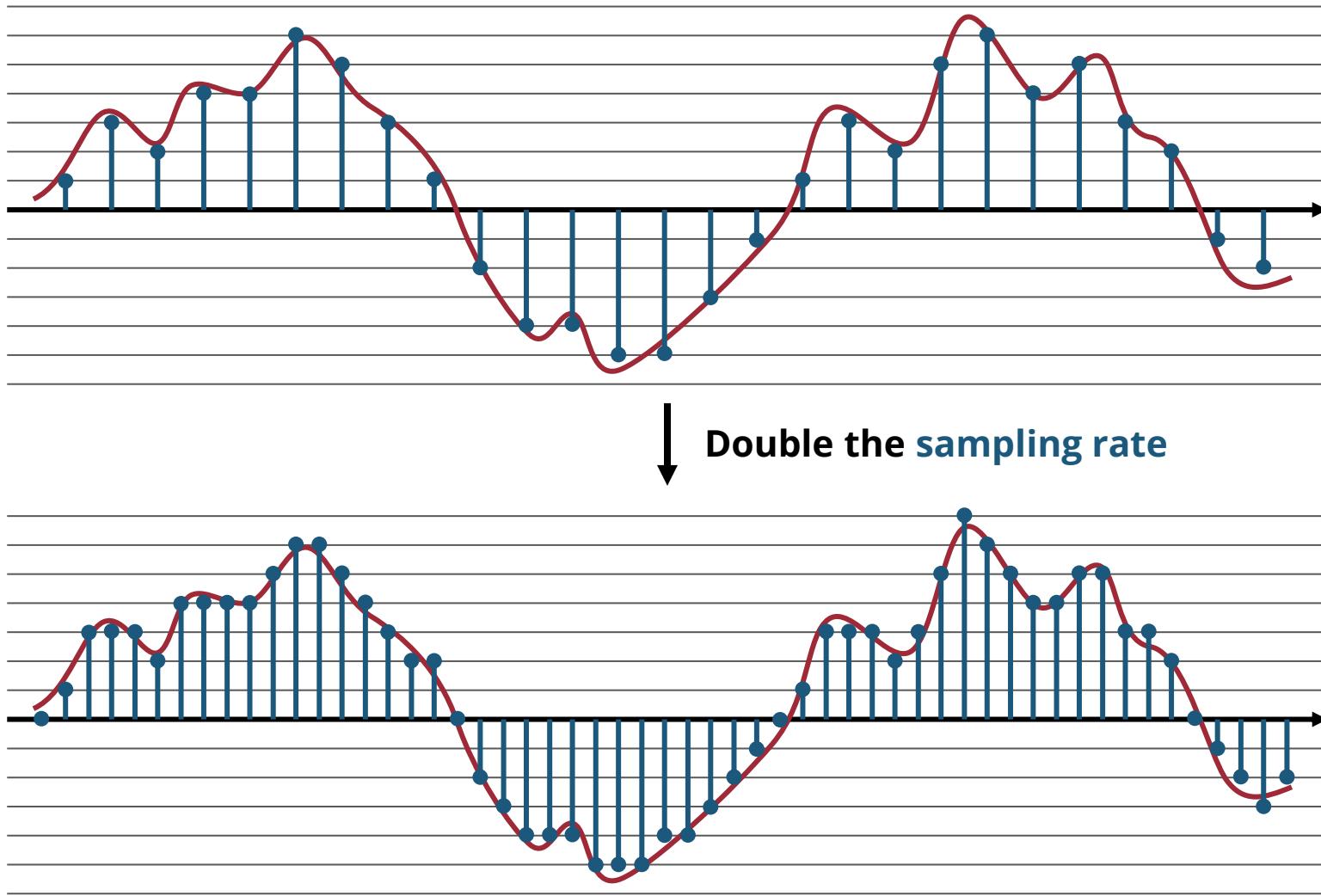
Digitalizing Audio: Amplitude



↓ Discretize the amplitude



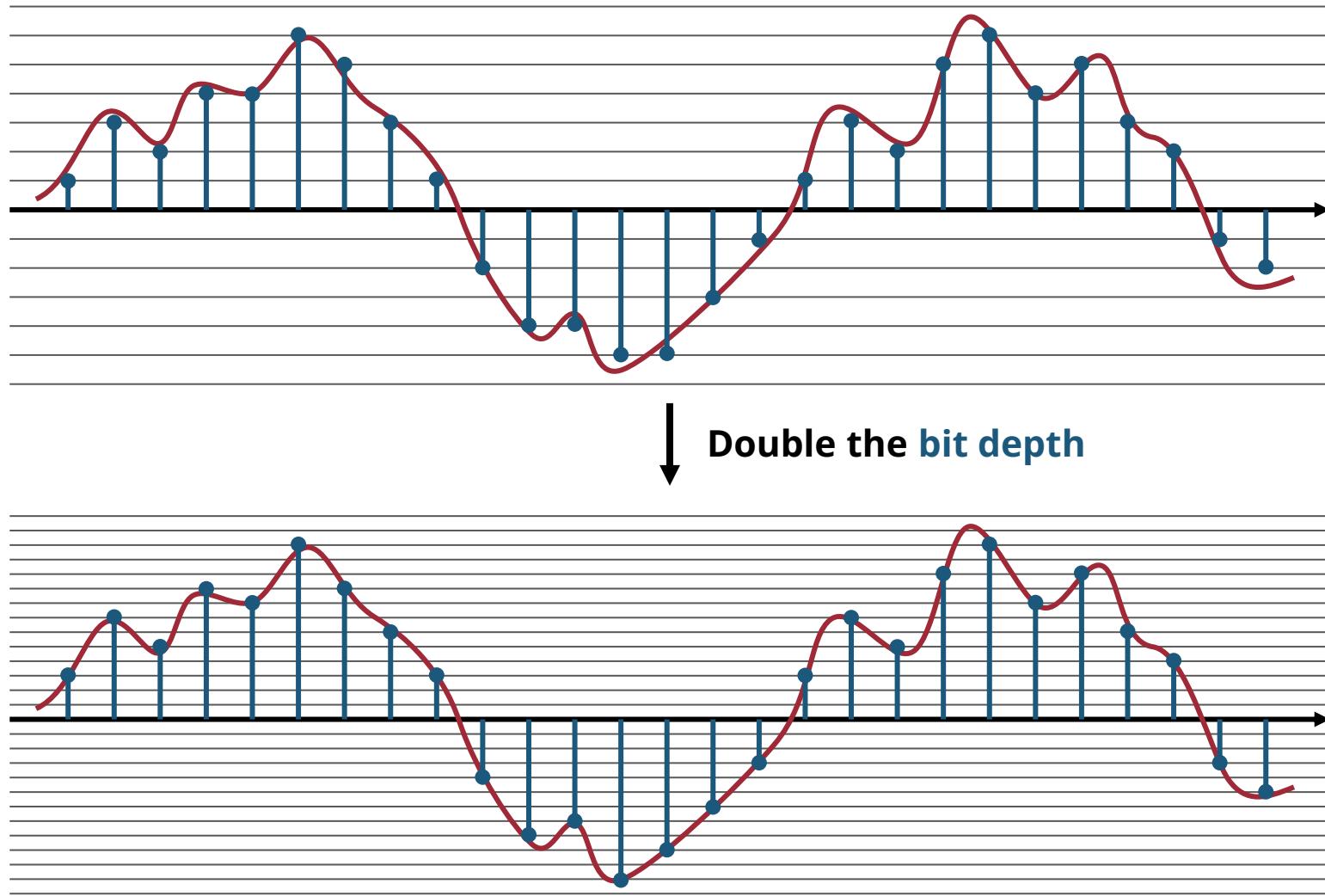
Resolution: Sampling Rate



Sampling Rate

- **Definition:** **Number of samples per second**
 - How many times the “sound pressure” is measured per second
 - The higher the sampling rate, the lower the timing distortion
- **Common sampling rates**
 - **Telephone:** 8 kHz
 - **CD:** 44.1 kHz
 - **DVD:** 48 kHz
 - **Modern audio interfaces & DAWs:** 96 kHz, 192 kHz

Resolution: Bit Depth



Bit Depth

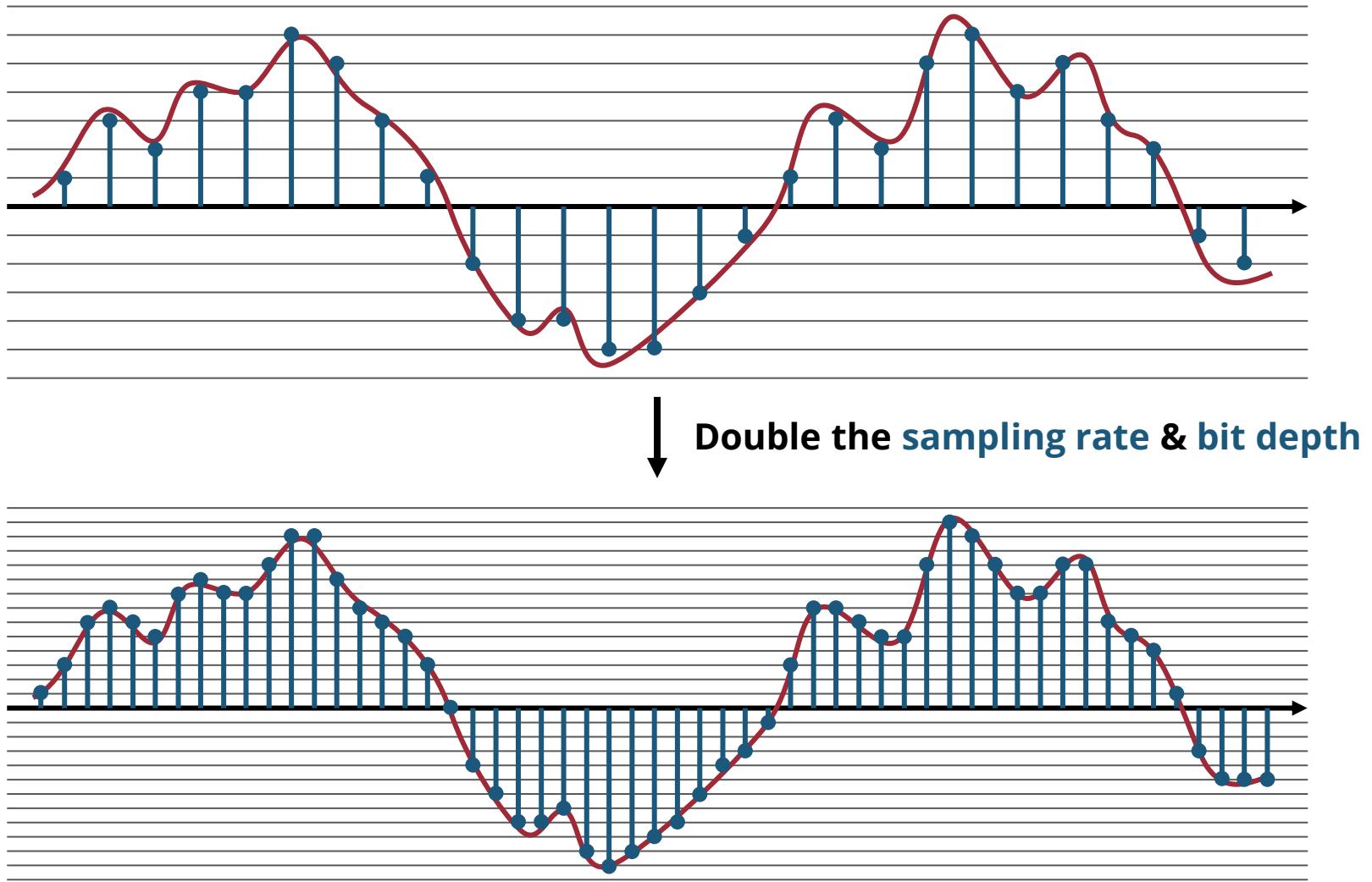
- **Definition:** Number of bits used to store each sample
 - How many bits used to store the amplitude
 - The higher the bit depth, the lower the amplitude distortion
- **Common bit depth**
 - **Chiptunes:** 8 bit
 - **CD:** 16 bit
 - **Modern audio interfaces & DAWs:** 24 bit, 32 bit



| Bit Depth

- **8 bit:** -128 to 127
- **16 bit:** -32,768 to 32,767
- **24 bit:** -8,388,608 to 8,388,607
- **32 bit:** 32-bit floating numbers

Resolution: Sampling Rate & Bit Depth



| Bit Depth ≠ Bit Rate

- **Bit Depth:** **Number of bits used to store each sample**
 - Example: **CD quality** is **16bit/44.1kHz**
- **Bit Rate:** **Amount of data transferred per second** (unit: bits/sec)
 - Example: **320K MP3** files → **320kbps** (320,000 bits per second)
 - Example: **YouTube** recommendation → **128 kbps** for mono and **384 kbps** for stereo
 - Determines the file size!

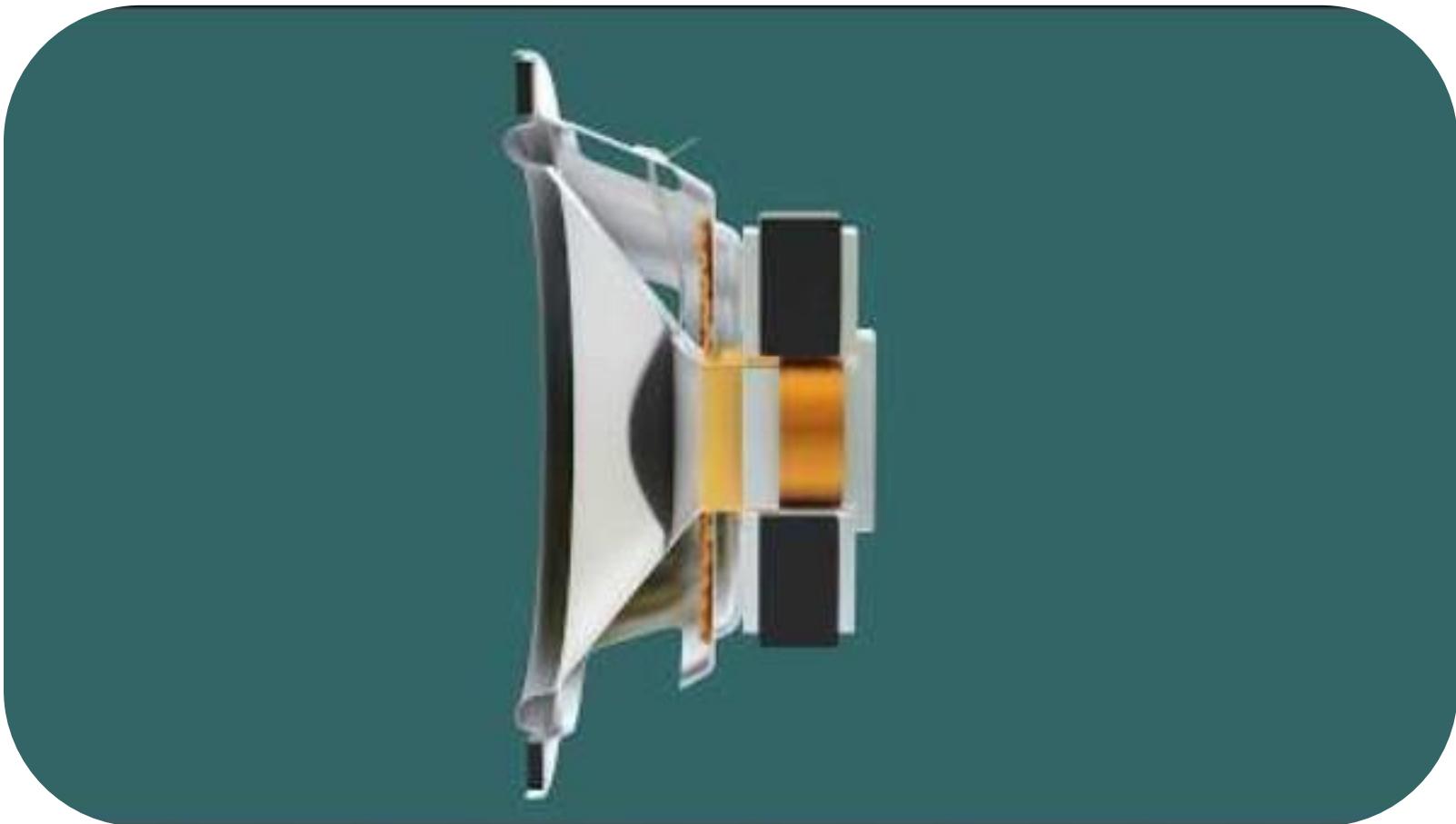
|  Reading: Microphones: Measuring Sound Pressure



youtu.be/d_crXXbuEKE



Reading: Speakers: Reproducing Sound Pressure

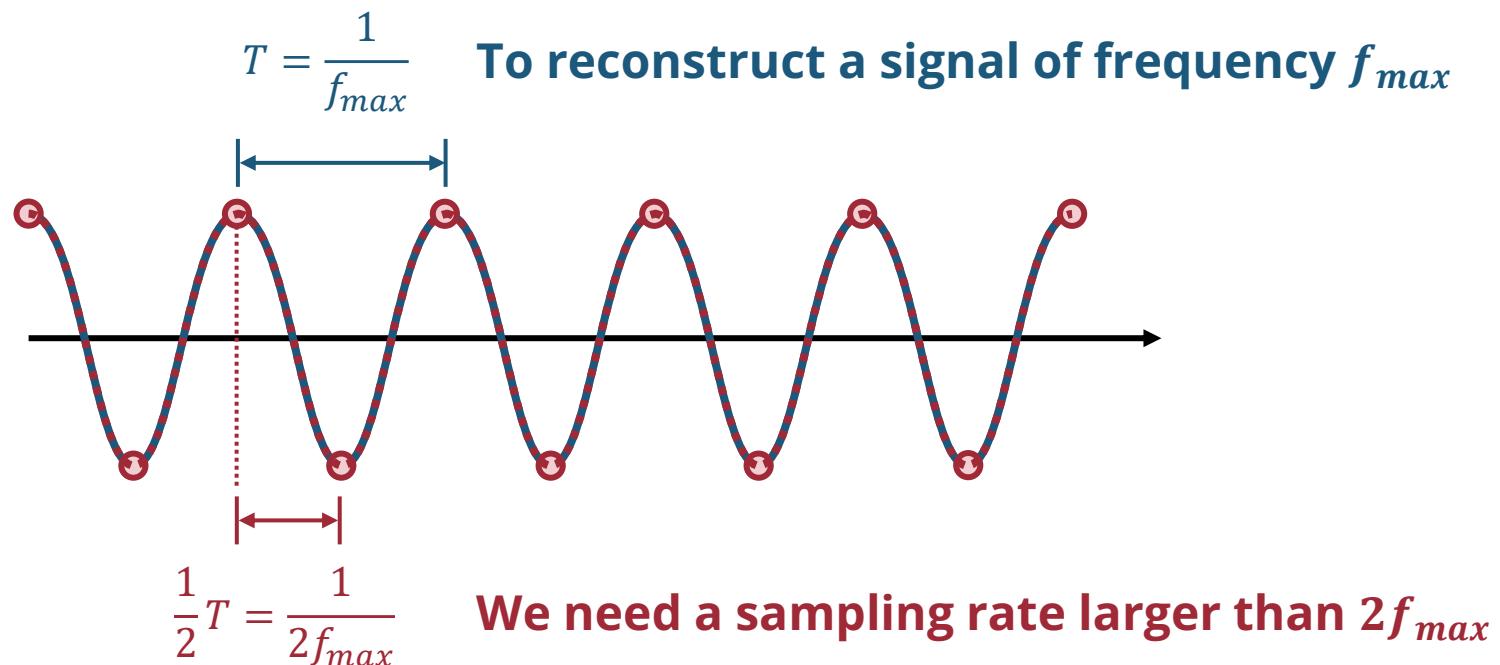


youtu.be/RxdFP31QYAg

Sampling Theorem

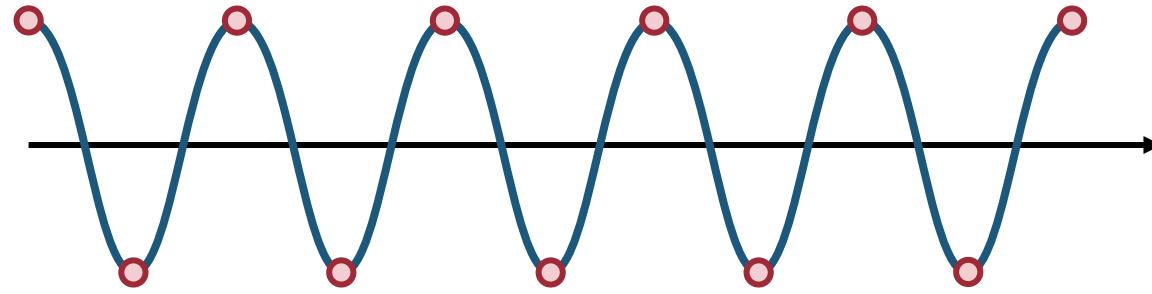
Nyquist–Shannon Sampling Theorem

- **Theorem:** If a signal contains no frequencies higher than f_{max} , then the signal can be perfectly reconstructed when sampled at a rate $f_s > 2f_{max}$
 - $2f_{max}$ is usually referred to as the **Nyquist rate**

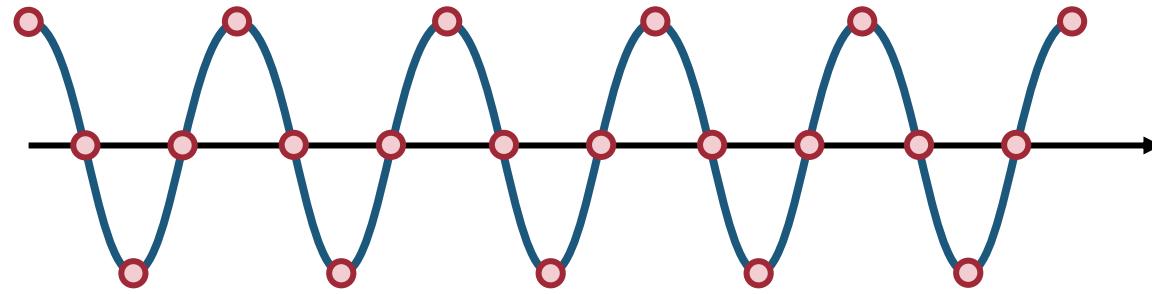


Sampling Theorem: Oversampling

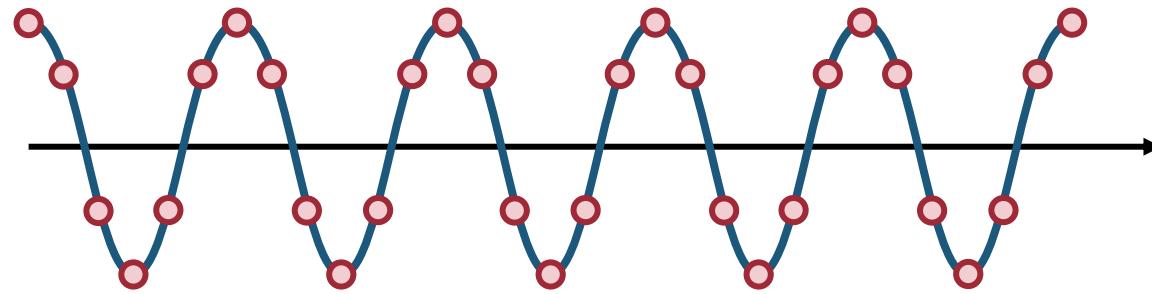
Critically sampled
 $(f_s = 2f_{max})$



Oversampled
 $(f_s = 4f_{max})$



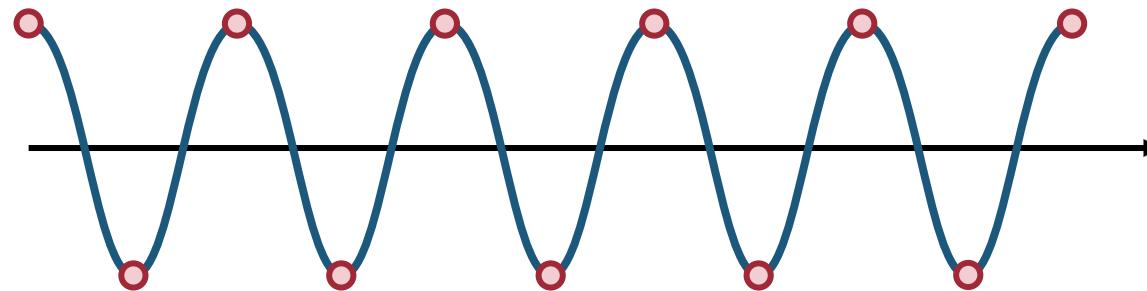
Oversampled
 $(f_s = 6f_{max})$



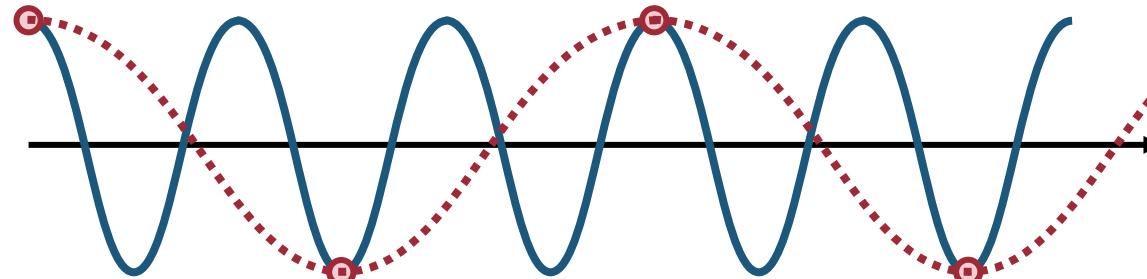
Reconstruction
is possible!

Sampling Theorem: Undersampling

Critically sampled
 $(f_s = 2f_{max})$

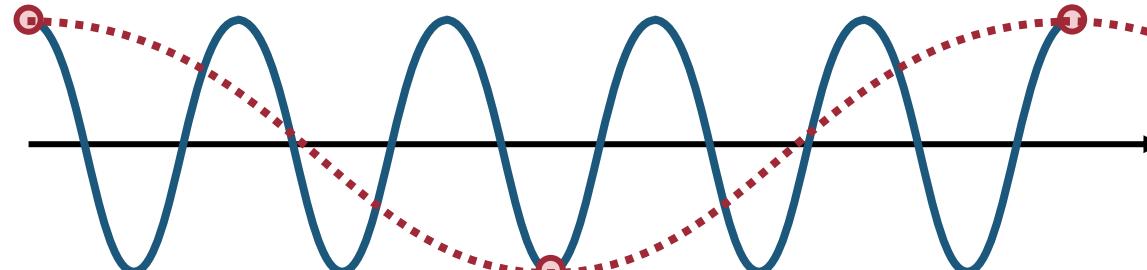


Undersampled
 $(f_s = \frac{2}{3}f_{max})$



Can only reconstruct frequency up to $\frac{1}{3}f_{max}$

Undersampled
 $(f_s = \frac{2}{5}f_{max})$



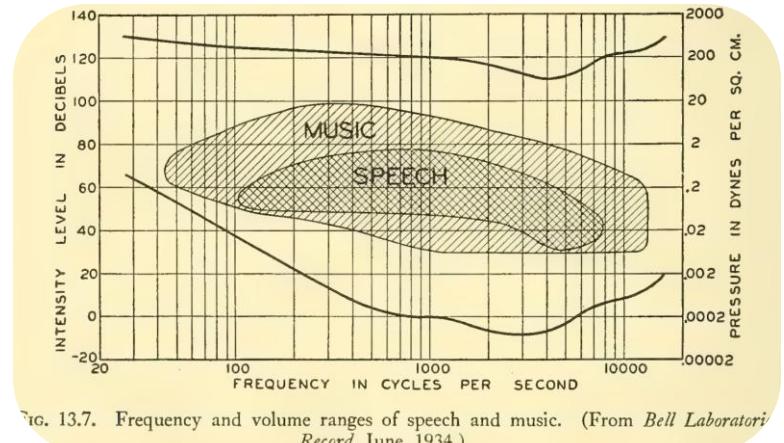
Can only reconstruct frequency up to $\frac{1}{3}f_{max}$



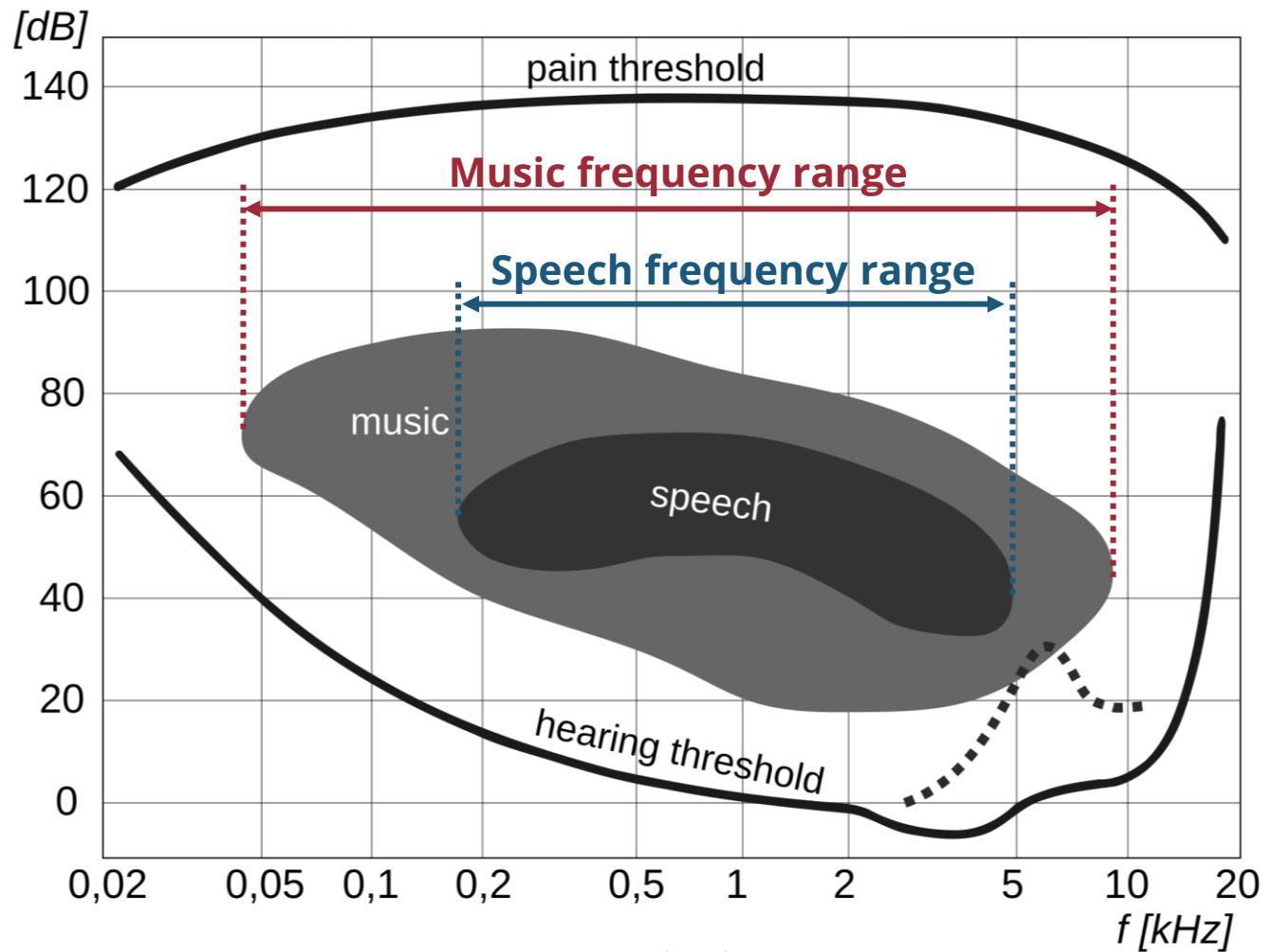
Sampling Theorem

- Telephone audio is sampled at **8 kHz**. What is the maximum frequency it can reconstruct?
 - **4 kHz**
- To cover the **human hearing range of 20 Hz to 20 kHz**, what is the minimum sampling rate required?
 - **40 kHz**

Sampling Rate & Frequency Range



(Source: Bell Laboratories Record 1934 & Olson 1947)



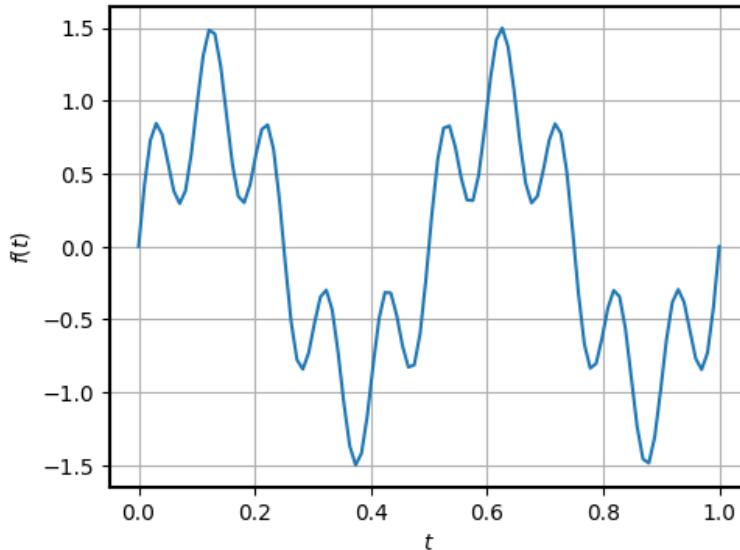
Bell Laboratories Record, 12(6):314, 1934.

Harry Ferdinand Olson, "Speech, Music and Hearing," *Elements of acoustical engineering Hardcover*, p. 326, 1947.
en.wikipedia.org/wiki/Hearing_range

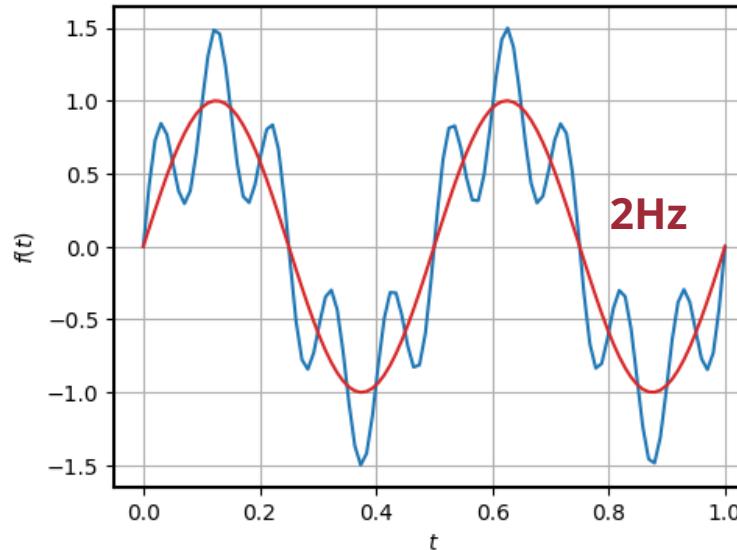
Spectral Analysis

Spectral Analysis

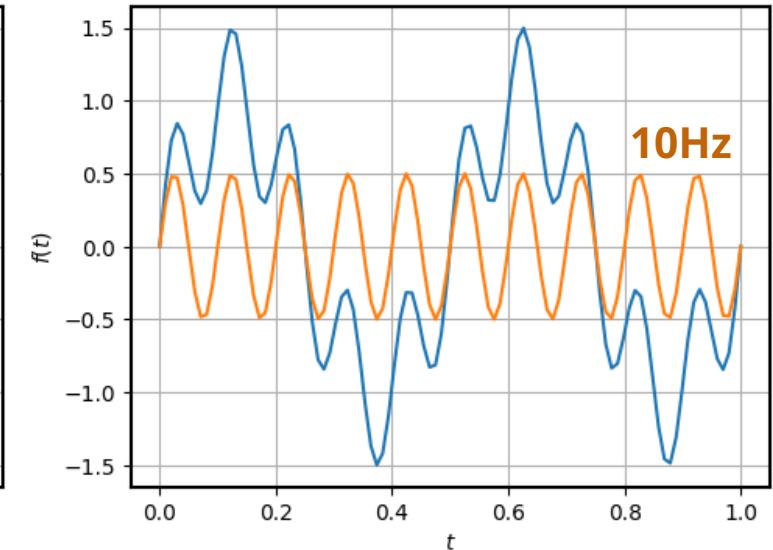
- **Goal:** Analyze the **frequency components** of a signal



$$\sin(2 \cdot 2\pi t) + \frac{1}{2} \sin(10 \cdot 2\pi t)$$

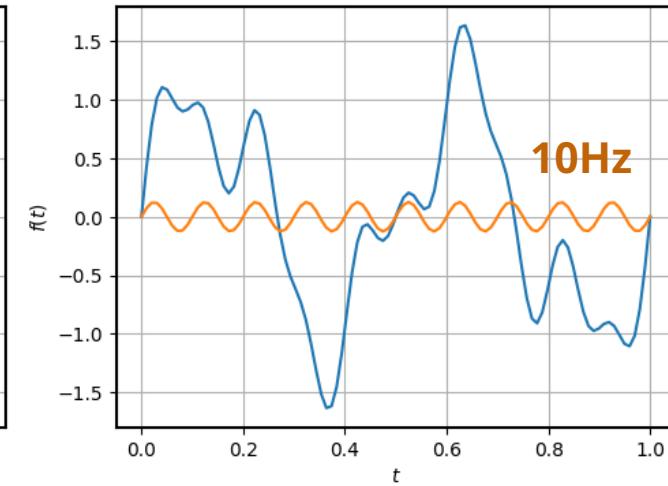
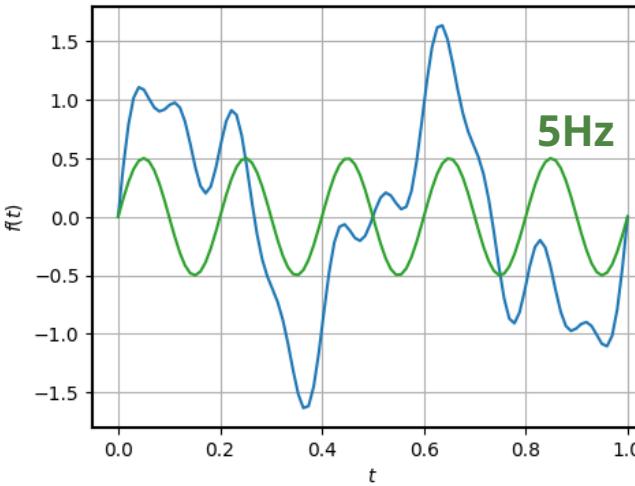
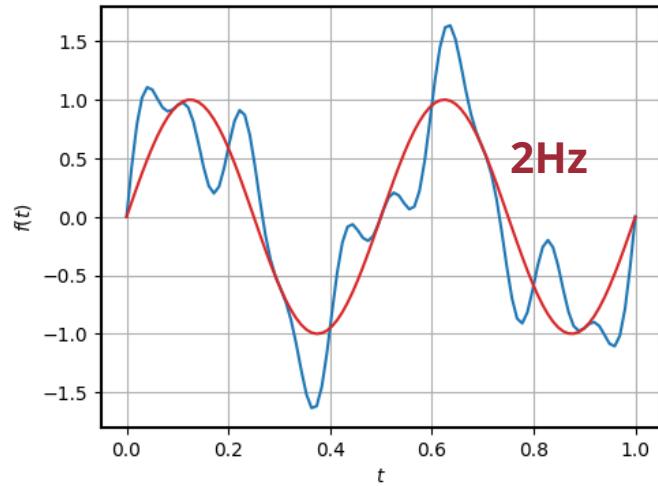
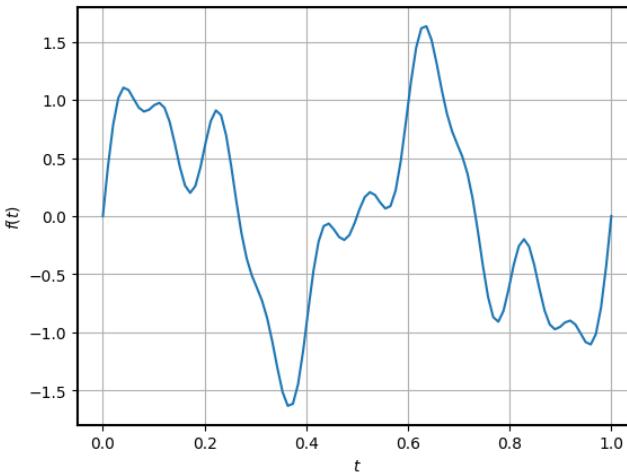


$$\sin(2 \cdot 2\pi t)$$



$$\frac{1}{2} \sin(10 \cdot 2\pi t)$$

Spectral Analysis



Fourier Transform

- **Intuition:** Decompose time-domain signals into **frequency components**
- Math formulation:

$$F(\omega) = \int_{-\infty}^{\infty} f(t) e^{-j\omega t} dt$$

Output spectrum Input signal
Frequency Sum over all t

The diagram illustrates the mathematical expression of the Fourier Transform. It features a green box labeled 'Output spectrum' containing the symbol $F(\omega)$. A blue box labeled 'Input signal' contains the expression $\int_{-\infty}^{\infty}$, followed by a red box with a thinking emoji containing $f(t)$, another red box containing $e^{-j\omega t}$, and a purple box containing dt . Arrows point from the text labels 'Frequency' and 'Sum over all t ' to their respective components in the equation.

| Demystifying Fourier Transform

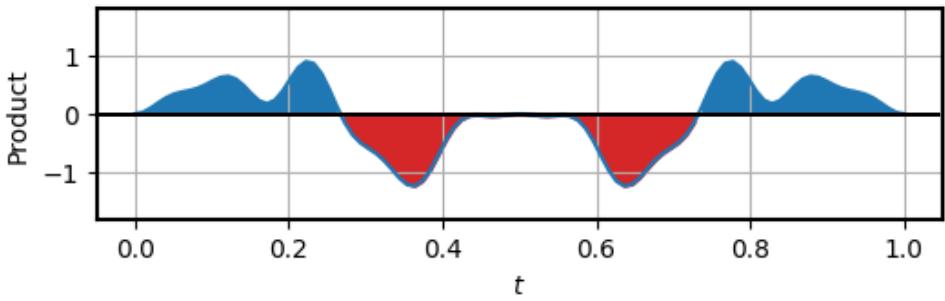
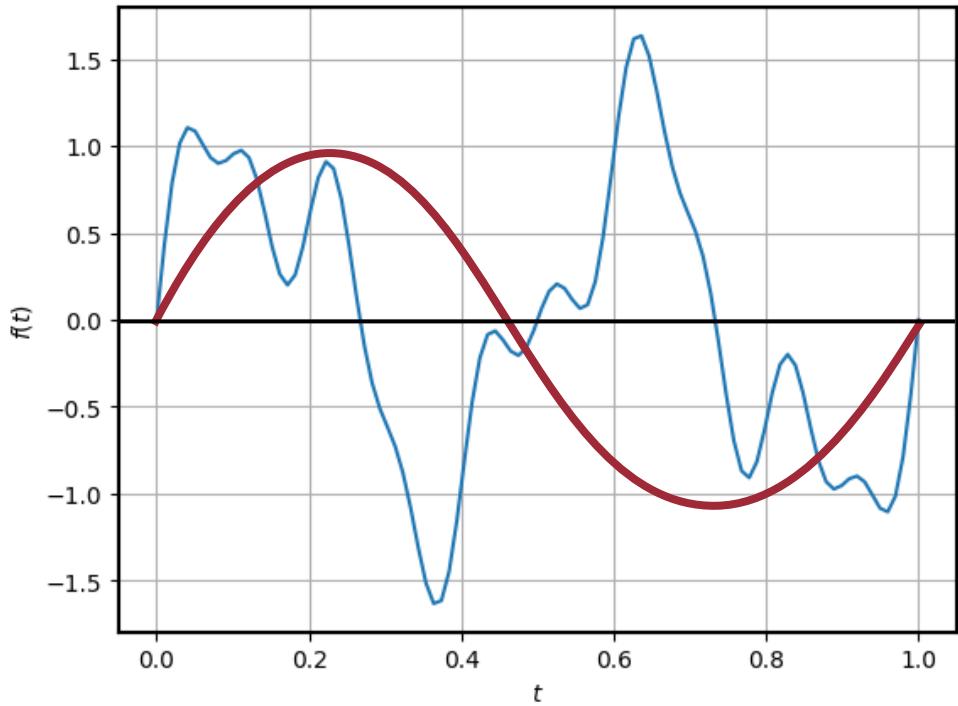
$$F(\omega) = \int_{-\infty}^{\infty} f(t) e^{-j\omega t} dt$$



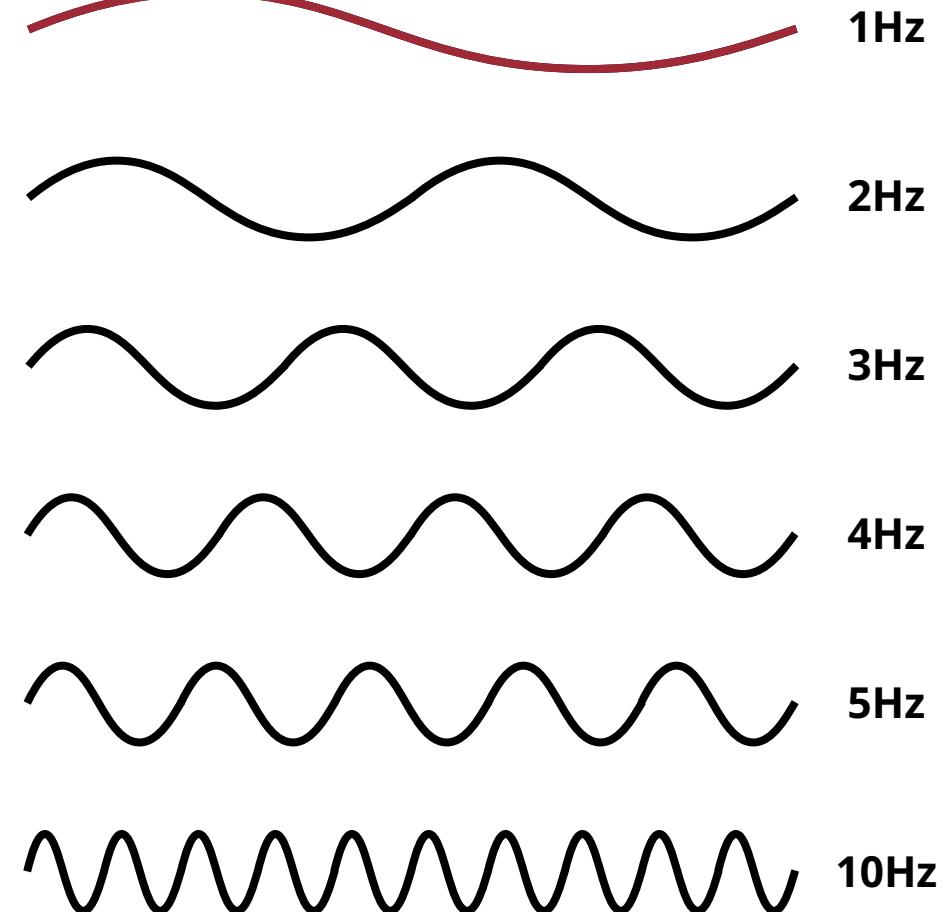
↓ Euler's formula
$$e^{-j\theta} = \cos \theta + j \sin \theta$$

$$F(\omega) = \int_{-\infty}^{\infty} f(t) \cos(-\omega t) + j f(t) \sin(-\omega t) dt$$

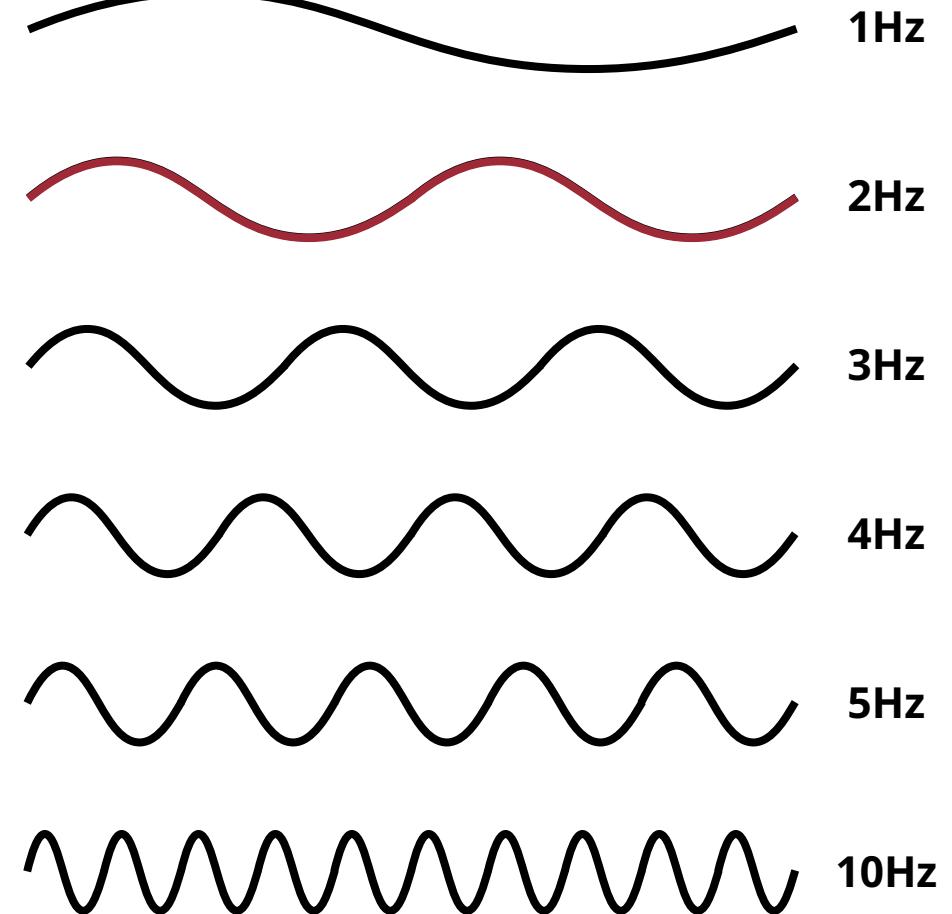
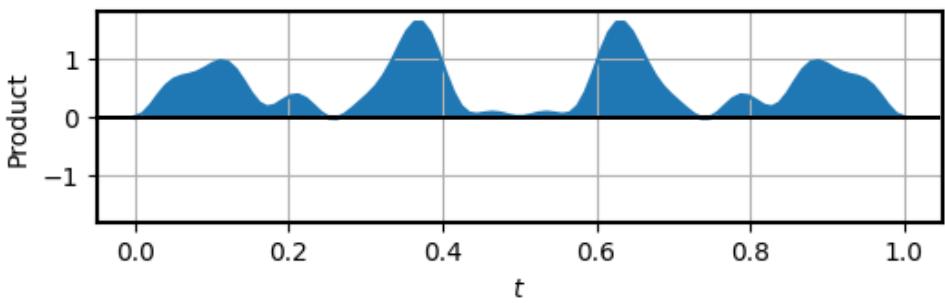
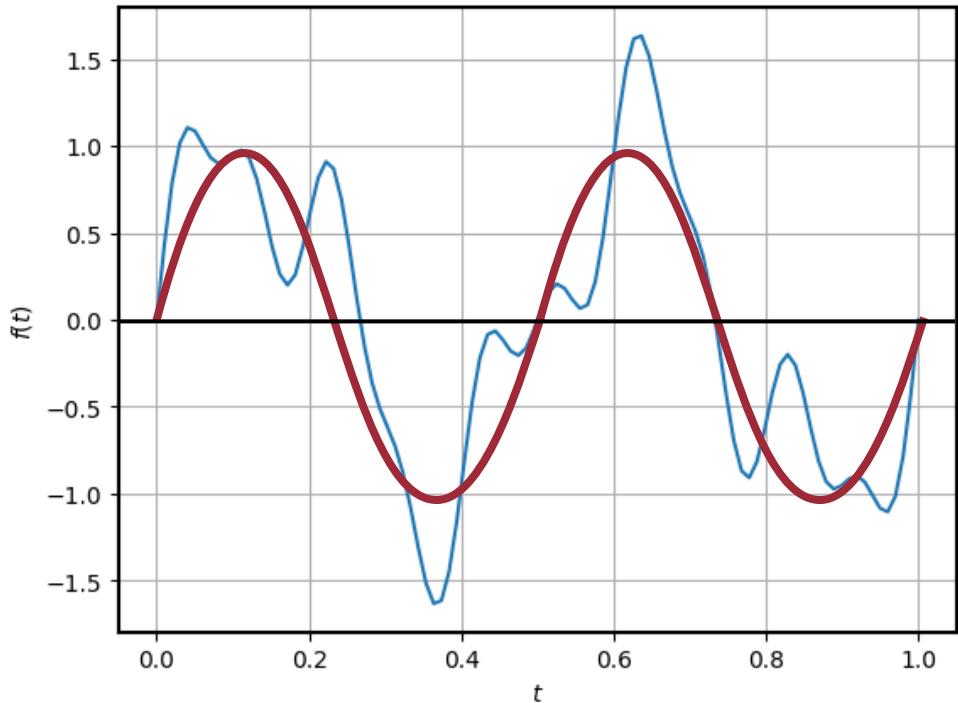
Demystifying Fourier Transform



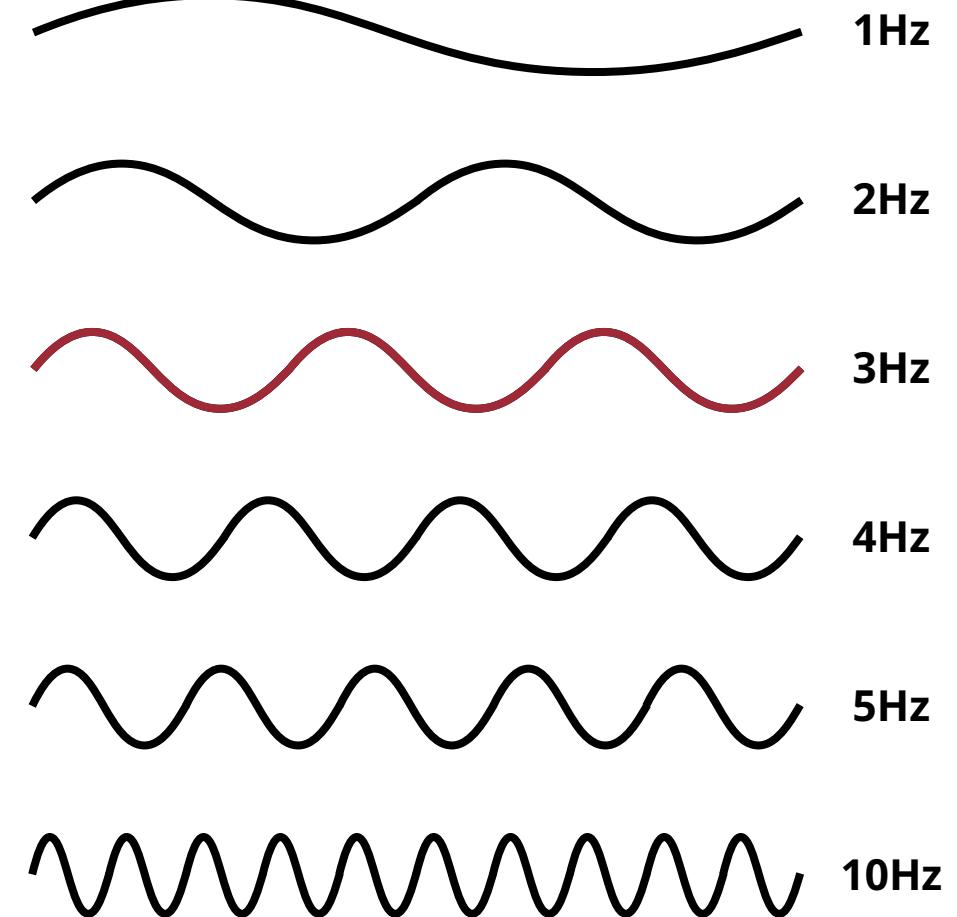
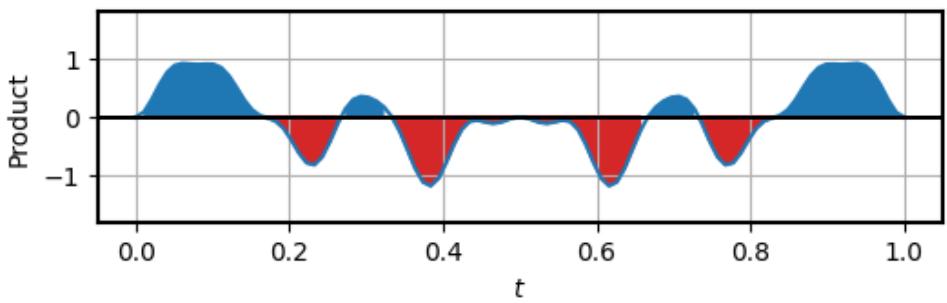
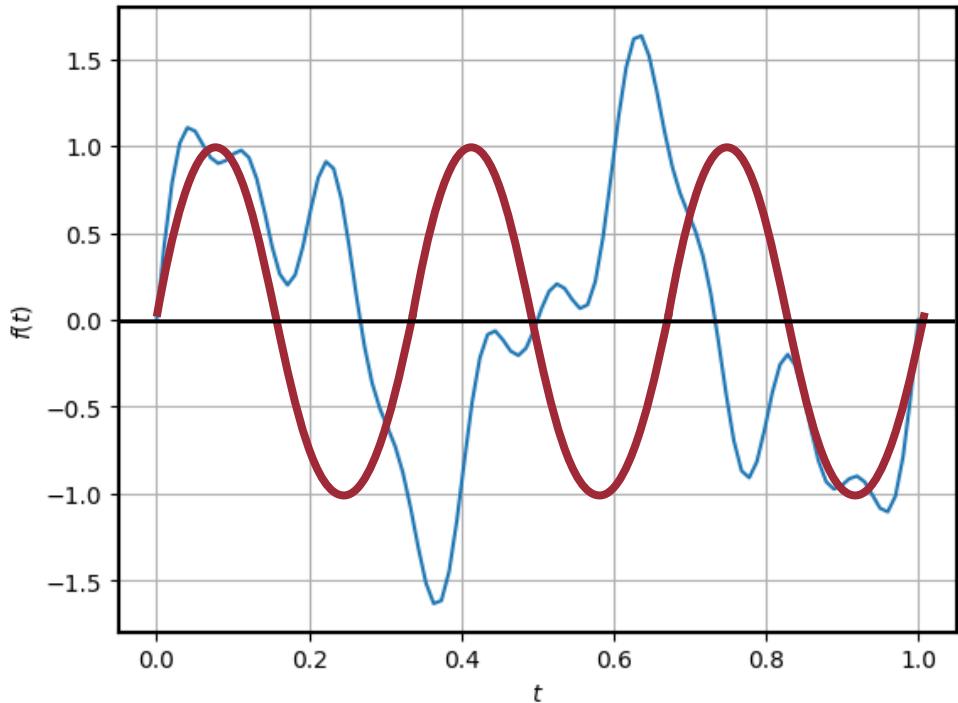
Candidate frequency components



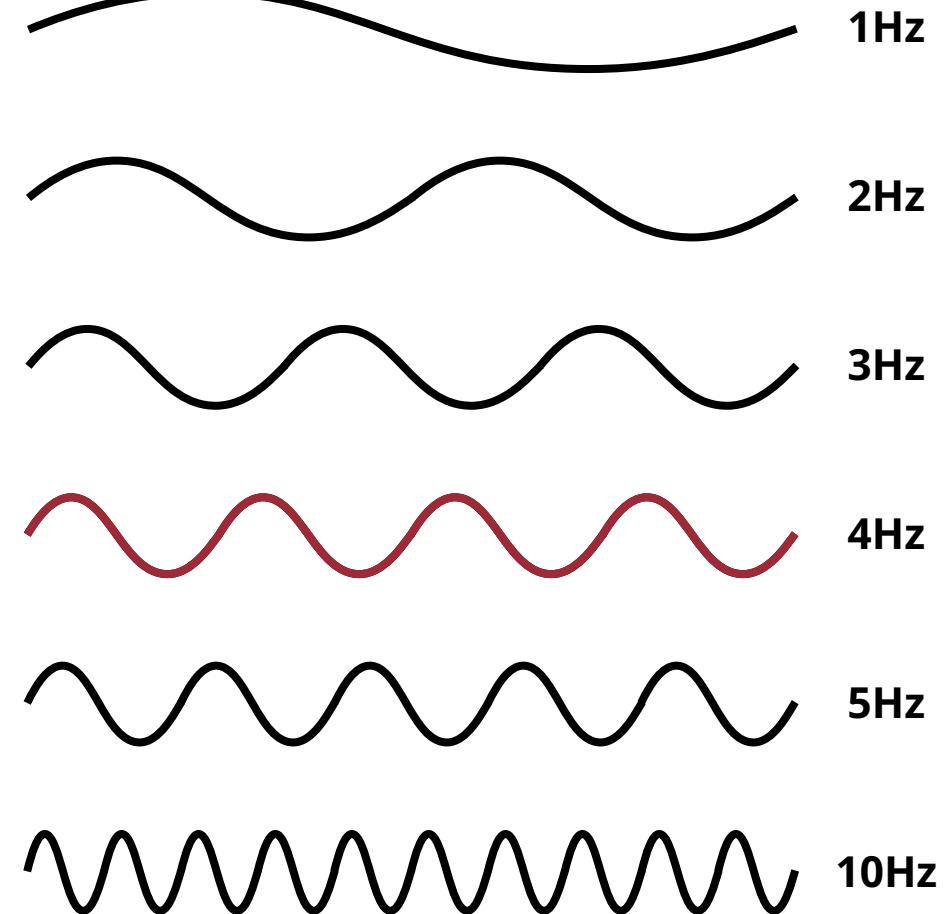
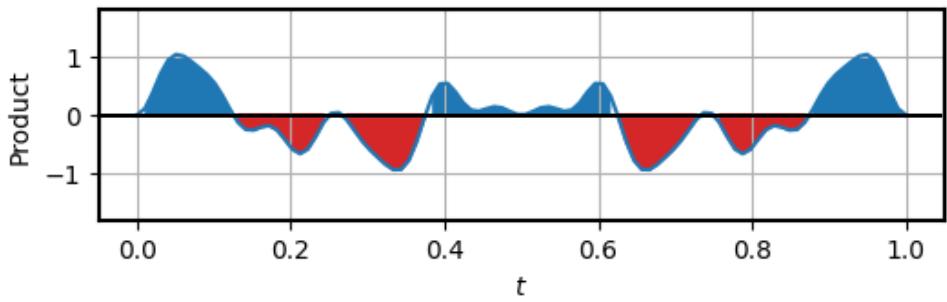
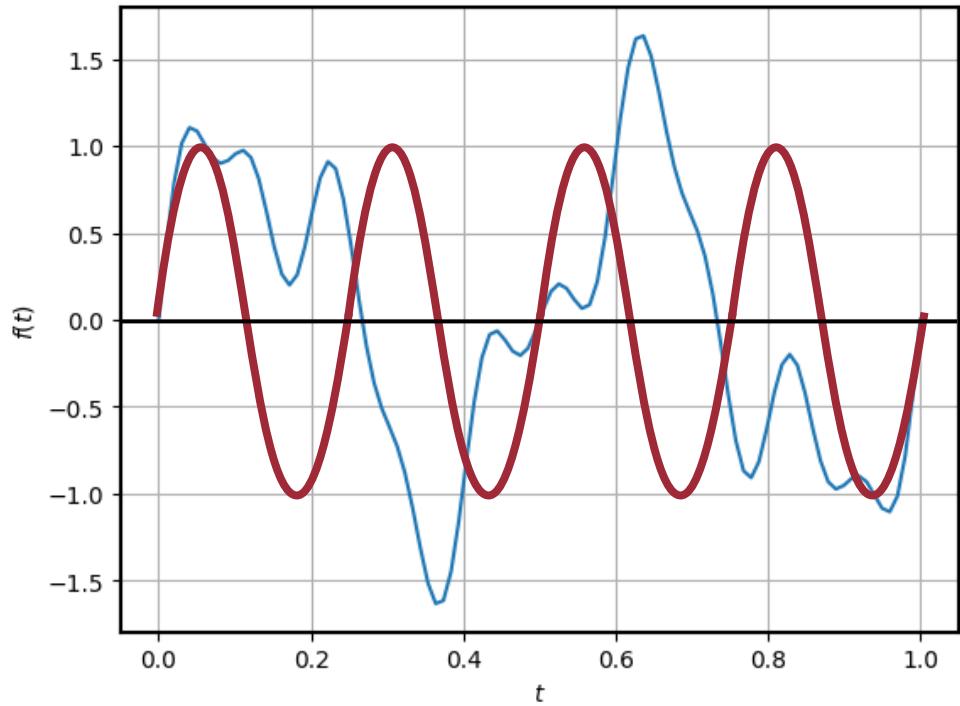
Demystifying Fourier Transform



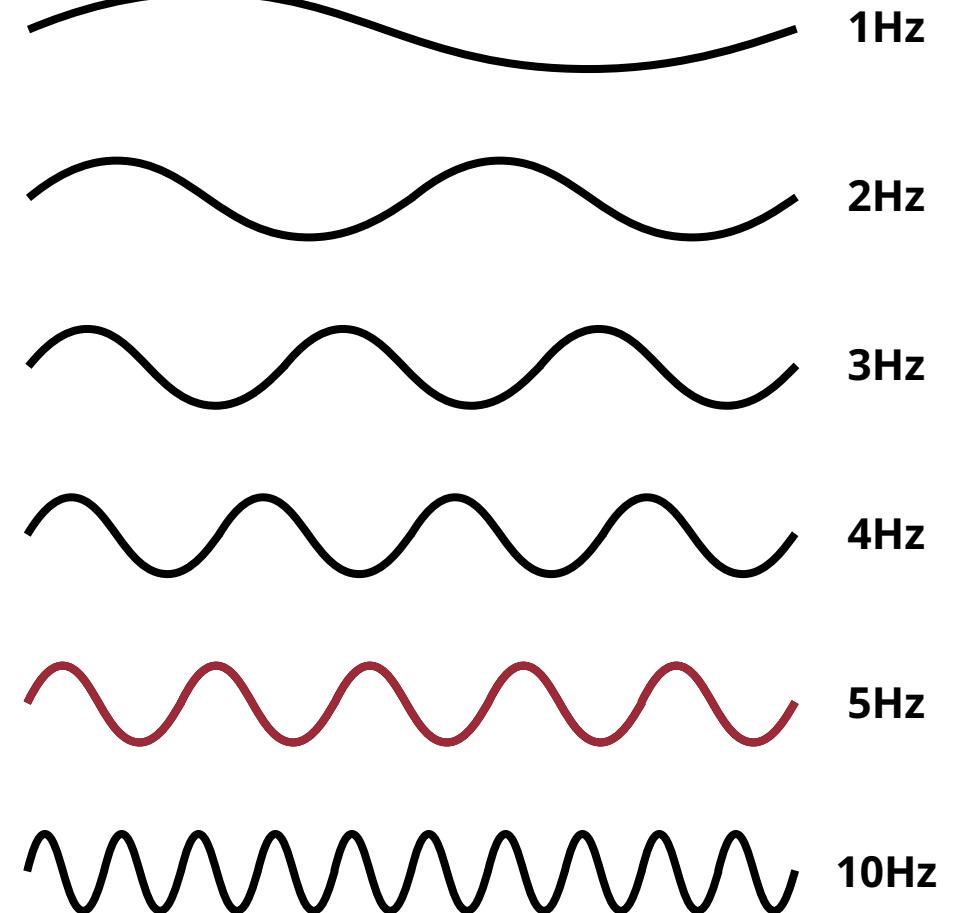
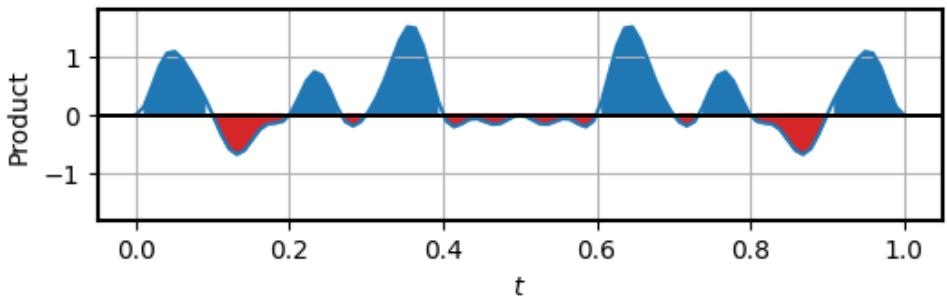
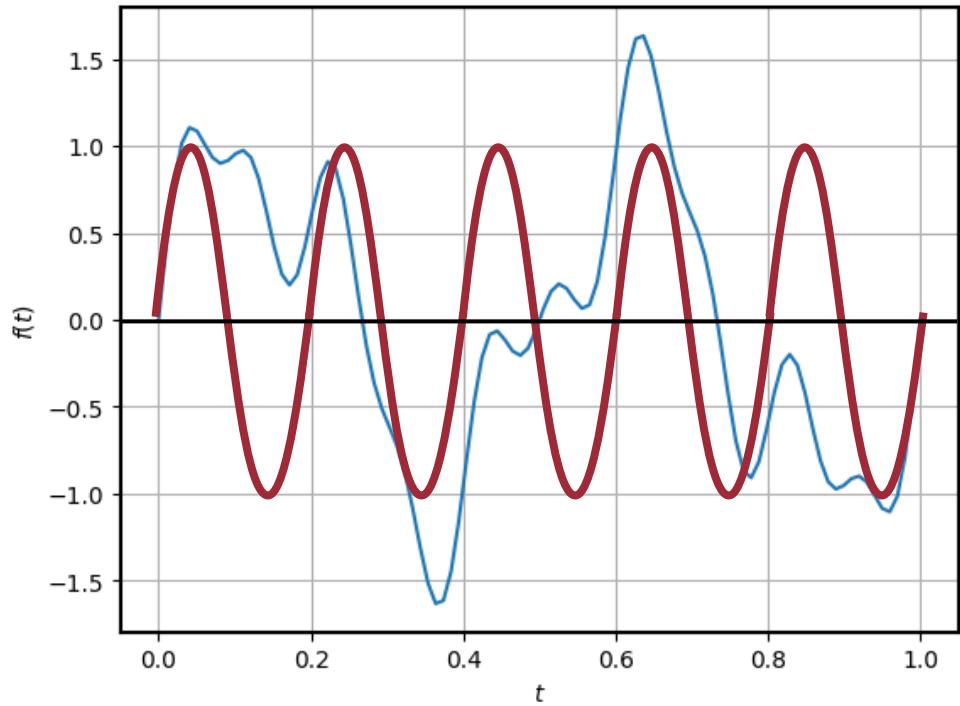
Demystifying Fourier Transform



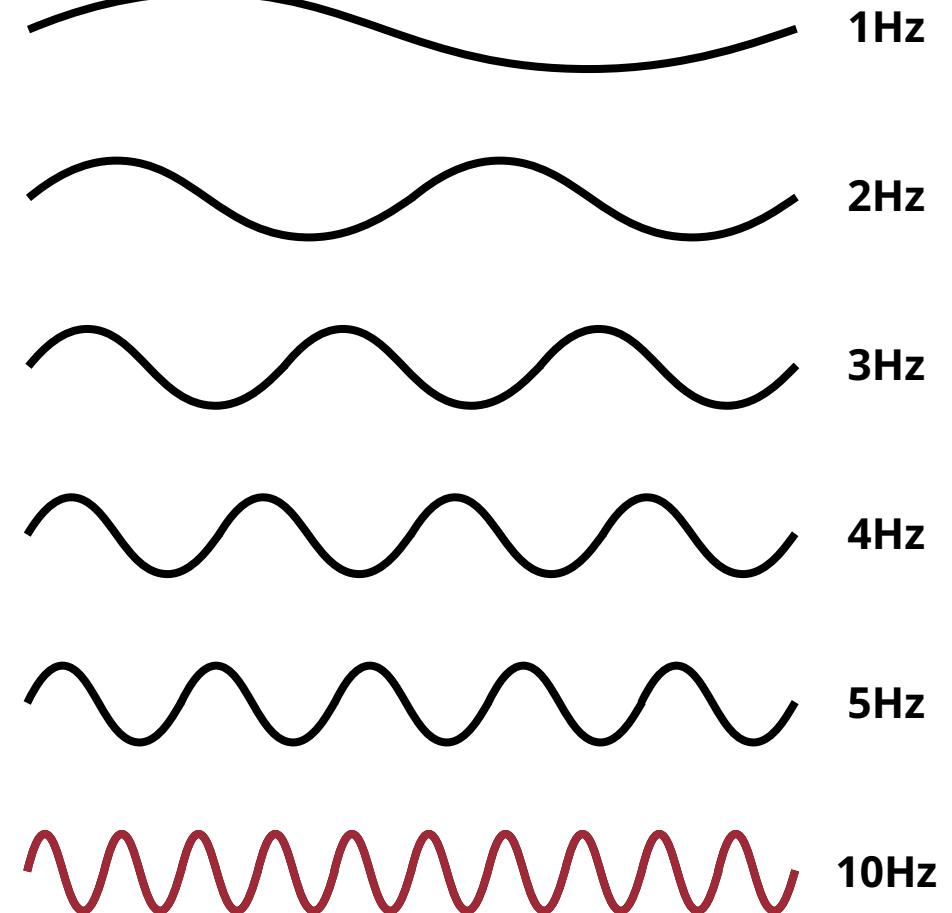
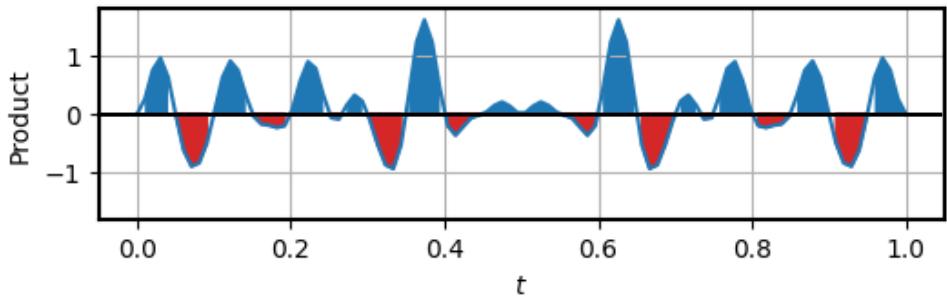
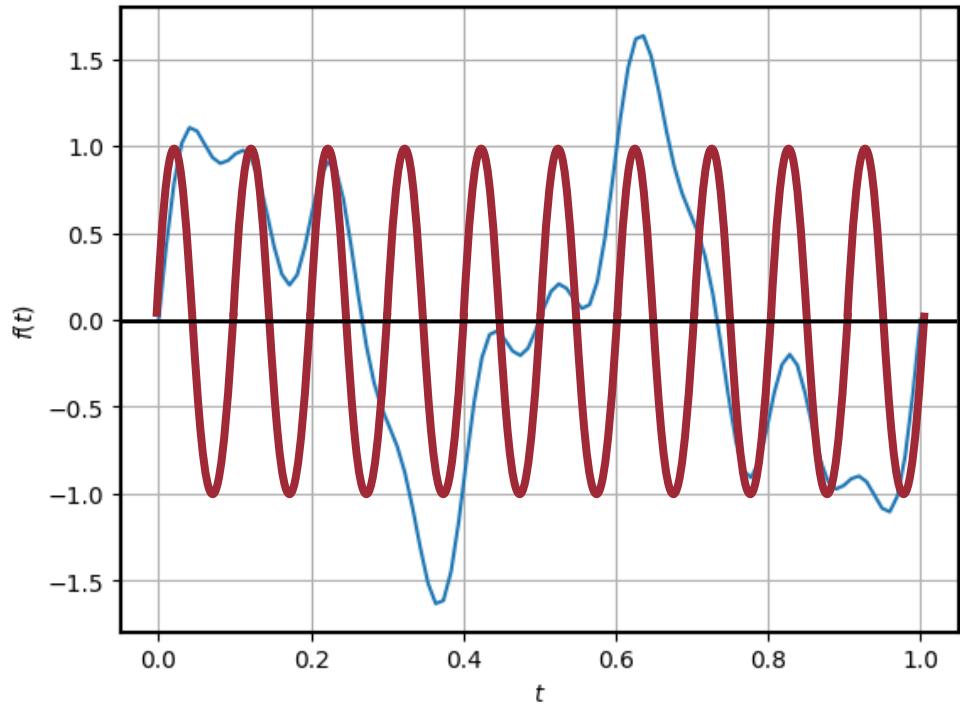
| Demystifying Fourier Transform



Demystifying Fourier Transform



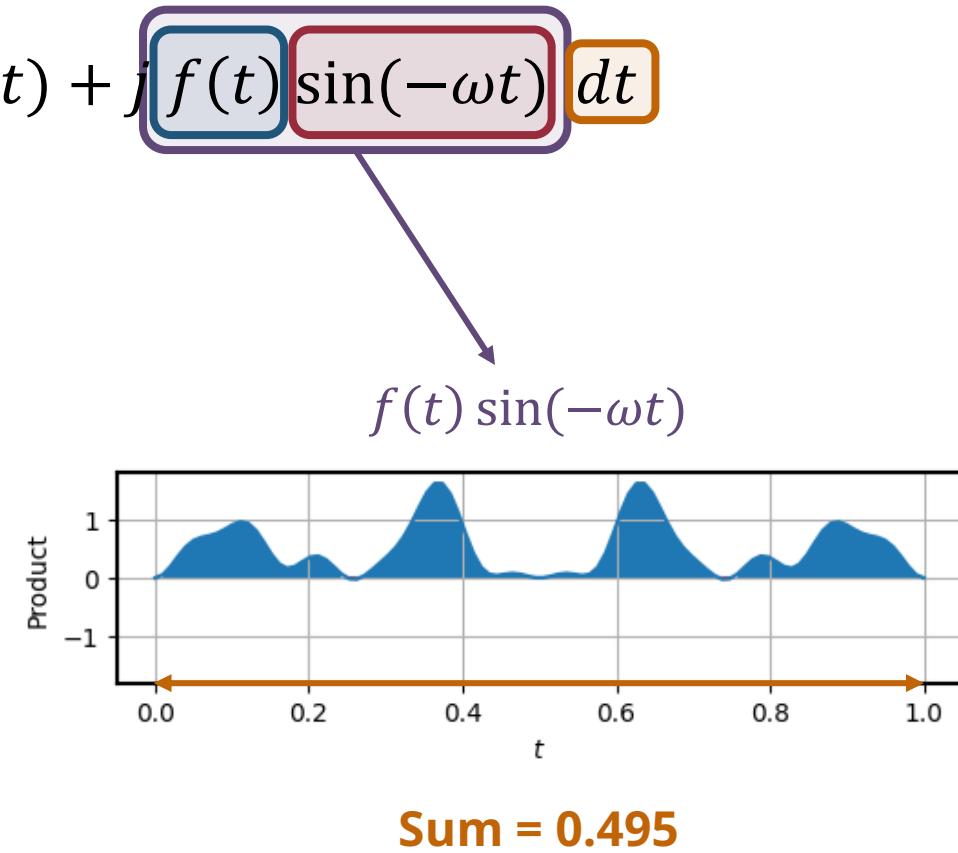
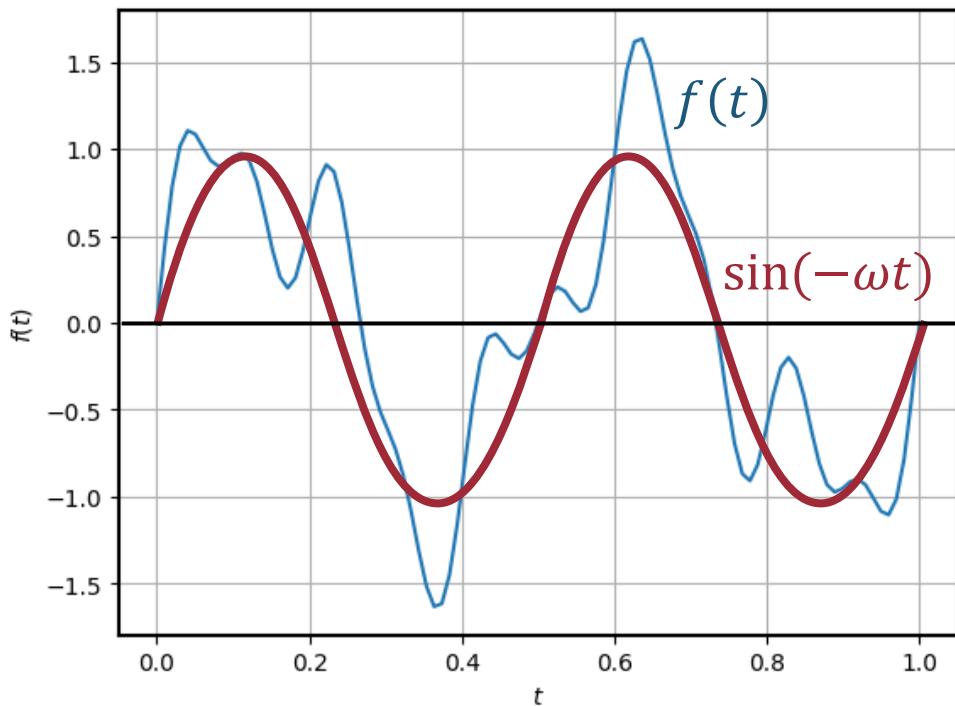
| Demystifying Fourier Transform



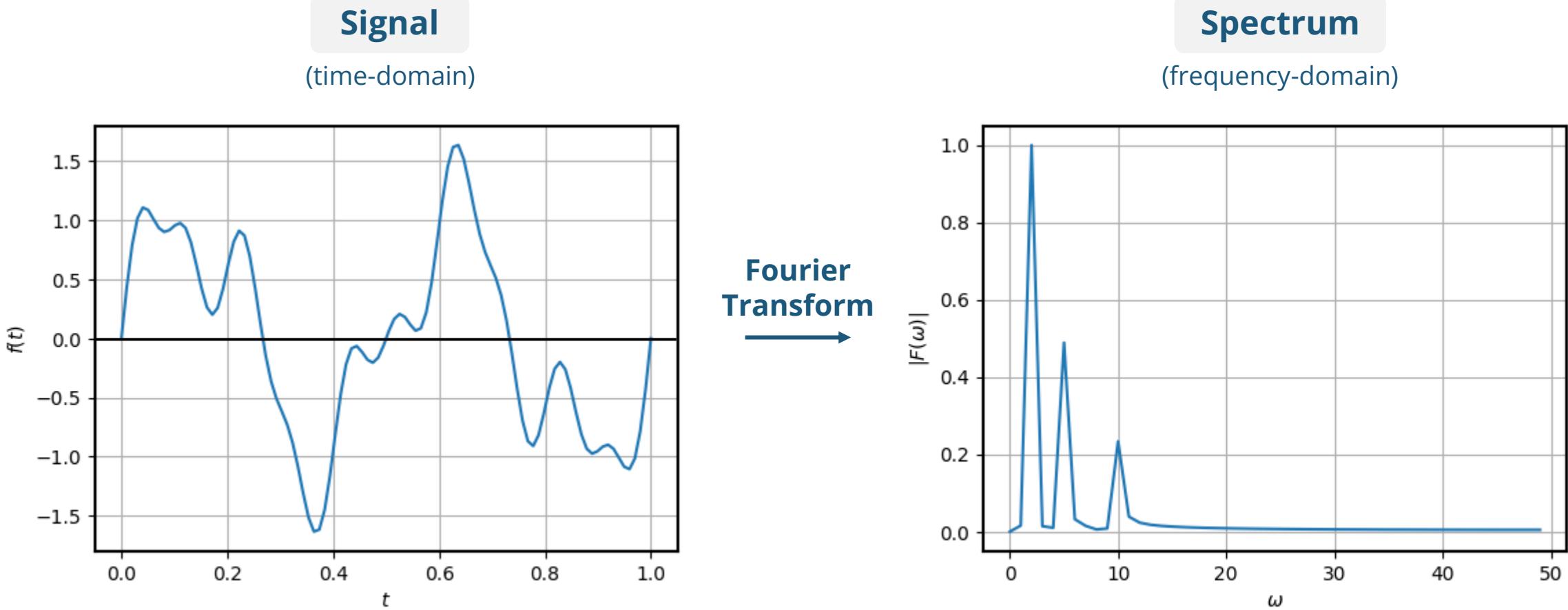
Demystifying Fourier Transform

$$F(\omega) = \int_{-\infty}^{\infty} f(t) \cos(-\omega t) + j \int_{-\infty}^{\infty} f(t) \sin(-\omega t) dt$$

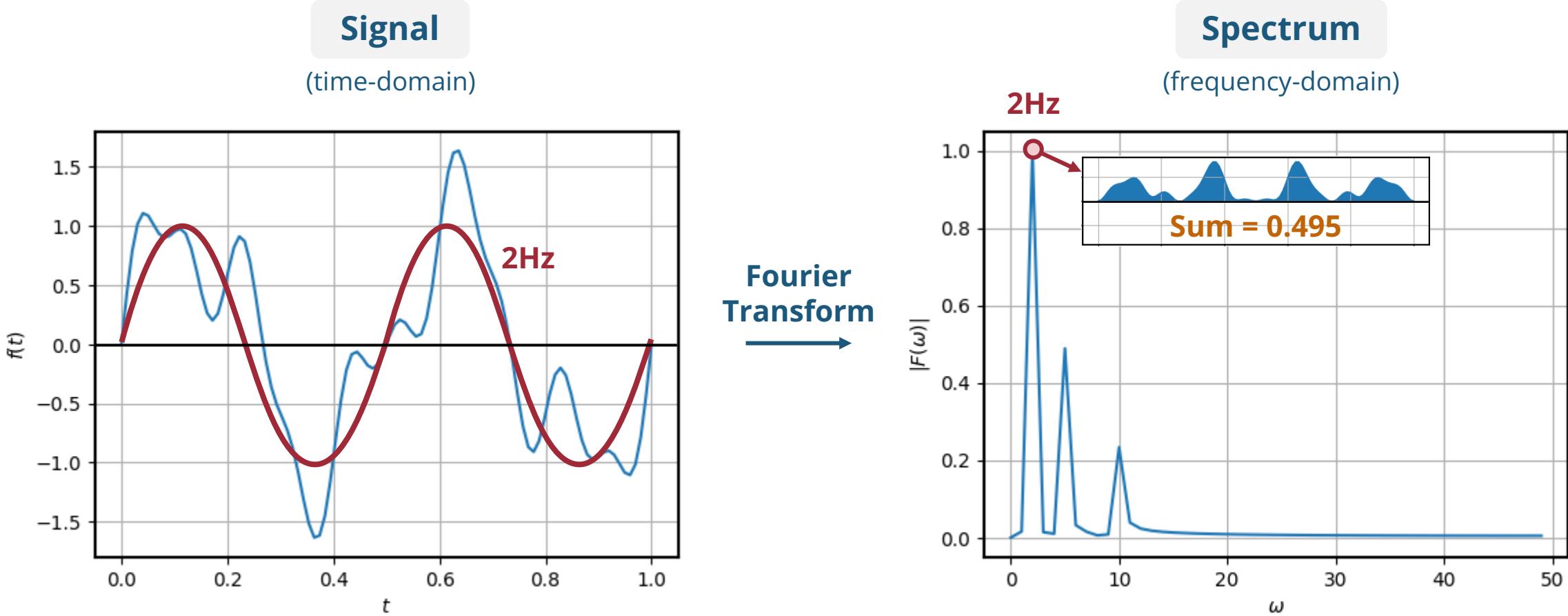
Sum over all t



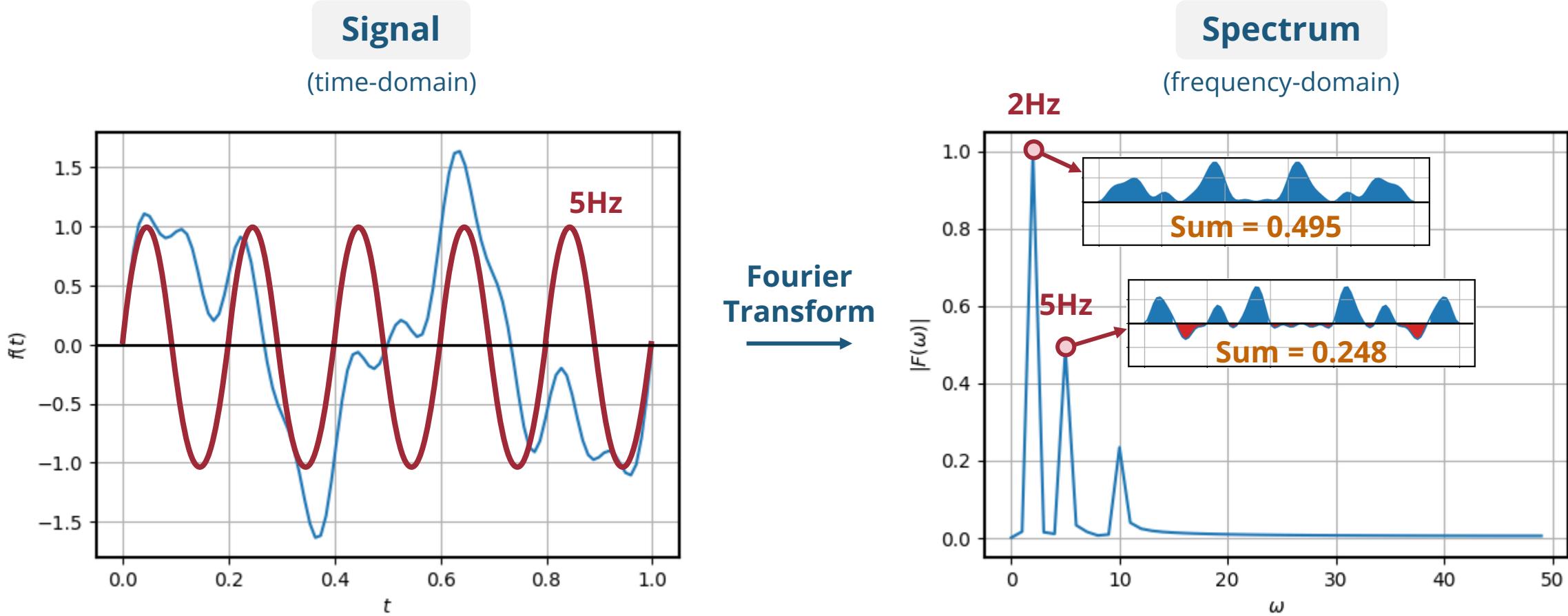
Demystifying Fourier Transform



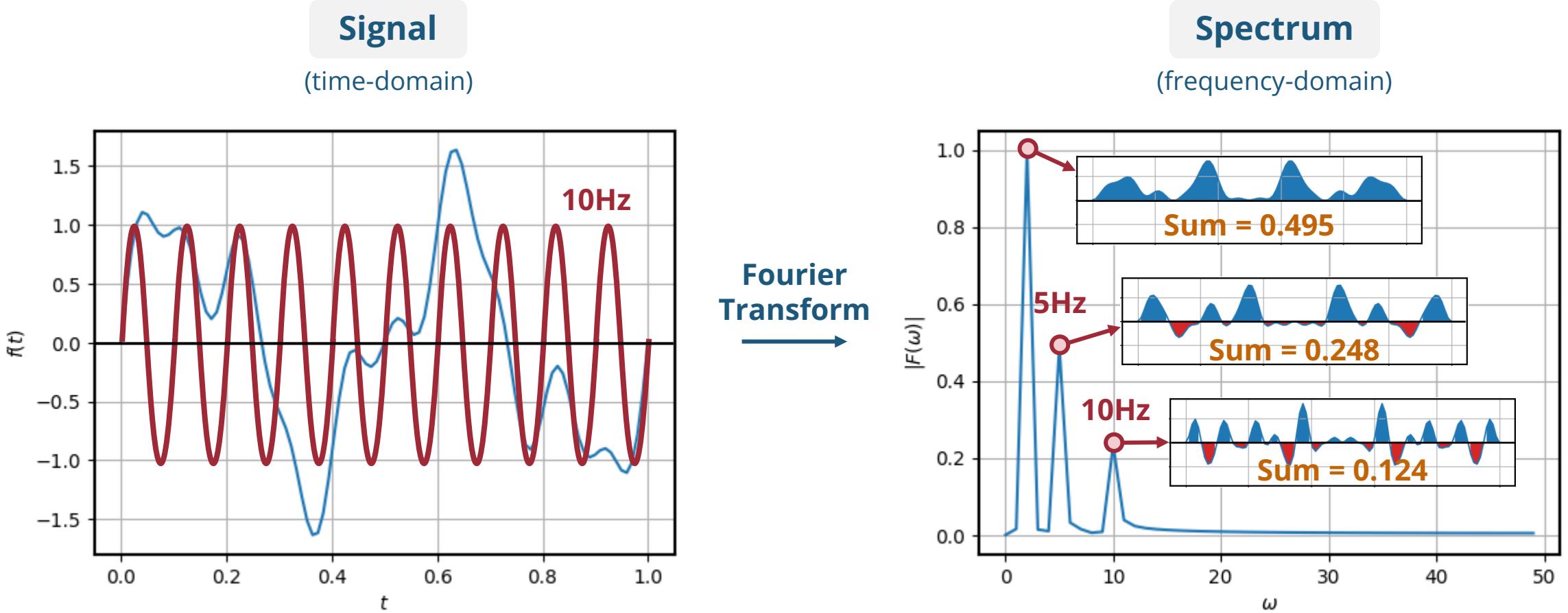
Demystifying Fourier Transform



Demystifying Fourier Transform



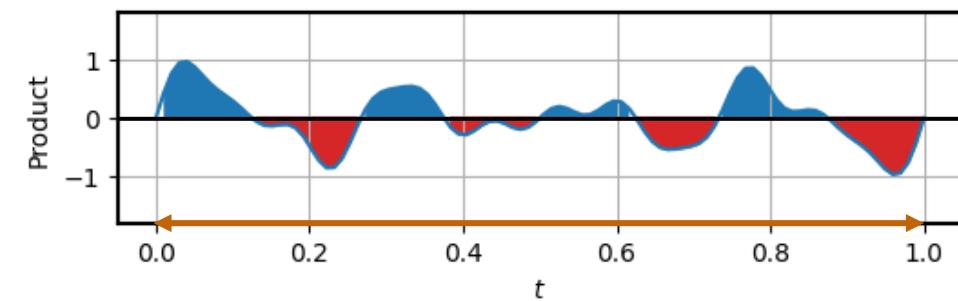
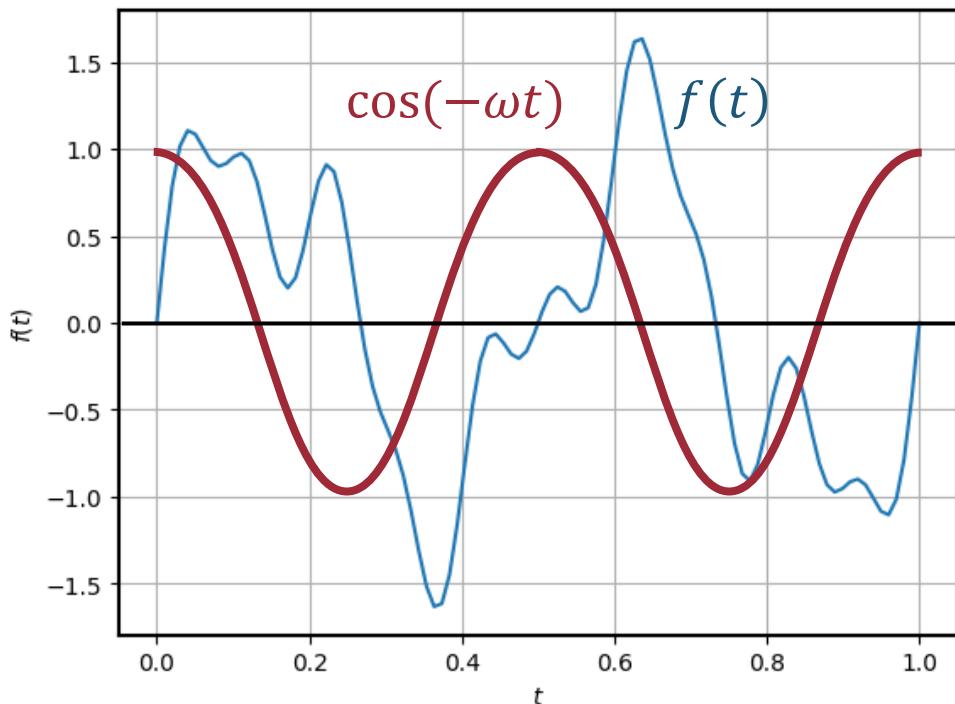
Demystifying Fourier Transform



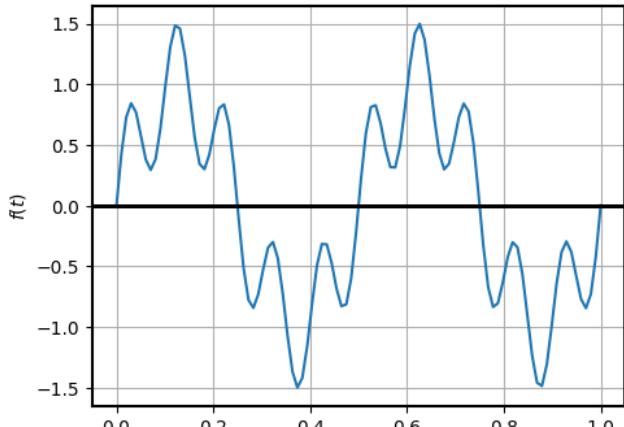
Demystifying Fourier Transform

$$F(\omega) = \int_{-\infty}^{\infty} f(t) \cos(-\omega t) + j f(t) \sin(-\omega t) dt$$

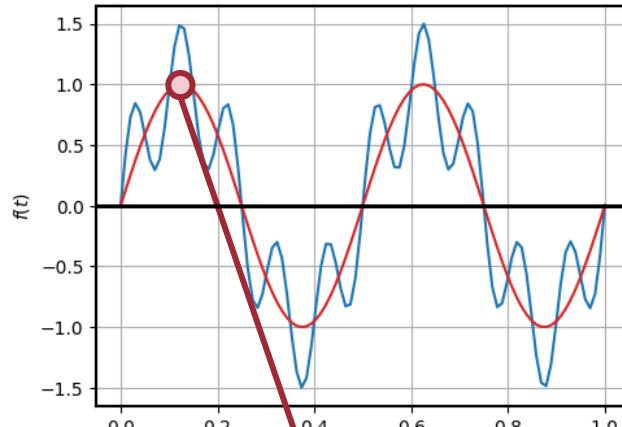
Sum over all t



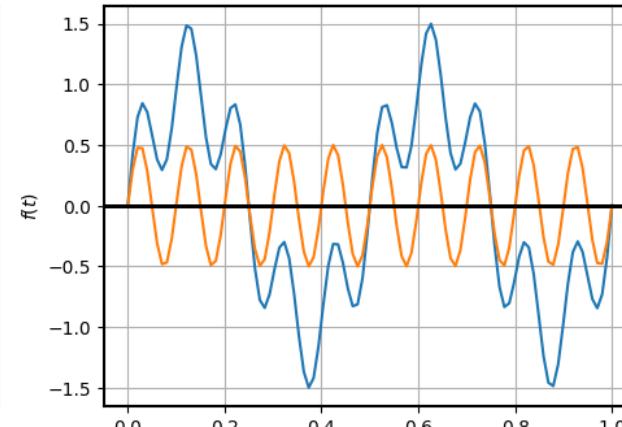
Phase



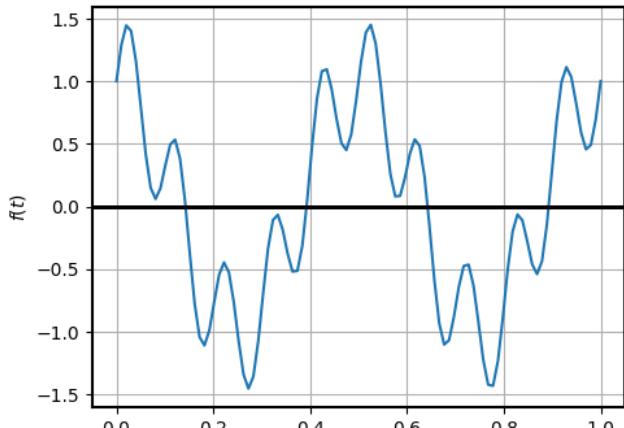
$$\sin(2 \cdot 2\pi t) + \frac{1}{2} \sin(10 \cdot 2\pi t)$$



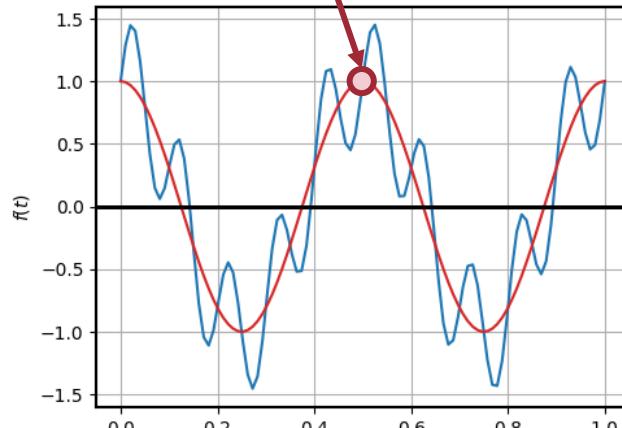
$$\sin(2 \cdot 2\pi t)$$



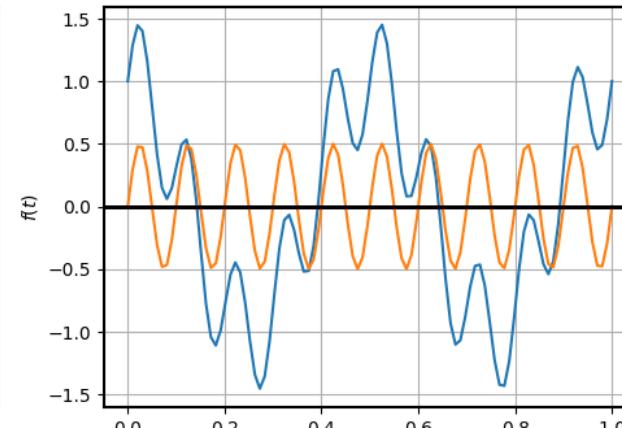
$$\frac{1}{2} \sin(10 \cdot 2\pi t)$$



$$\cos(2 \cdot 2\pi t) + \frac{1}{2} \sin(10 \cdot 2\pi t)$$



$$\cos(2 \cdot 2\pi t)$$

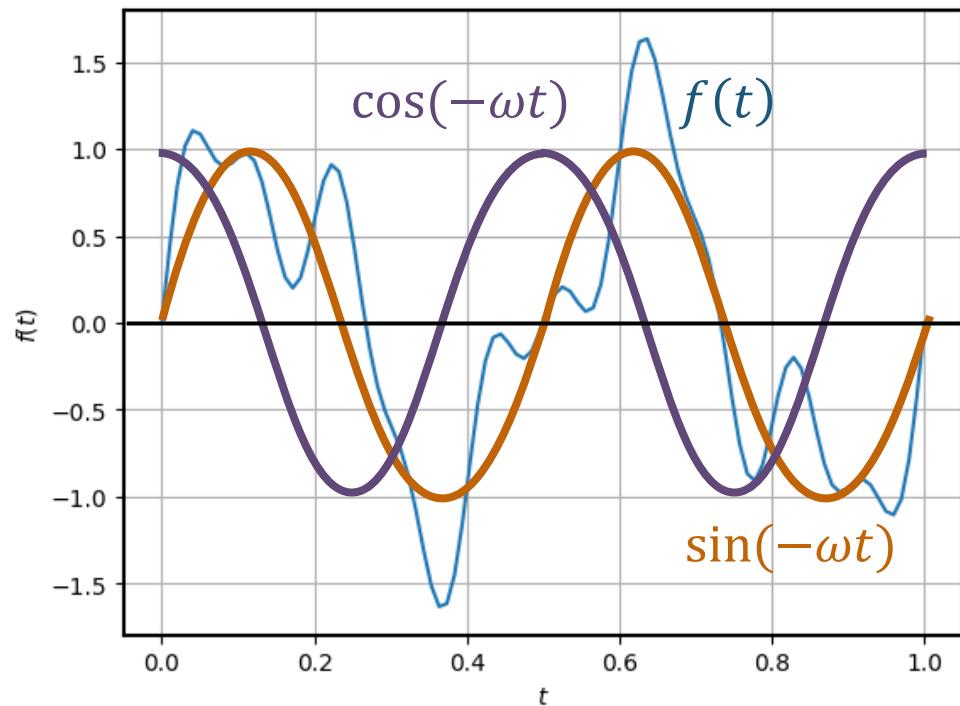


$$\frac{1}{2} \sin(10 \cdot 2\pi t)$$

| Demystifying Fourier Transform

$$F(\omega) = \int_{-\infty}^{\infty} f(t) \cos(-\omega t) + j f(t) \sin(-\omega t) dt$$

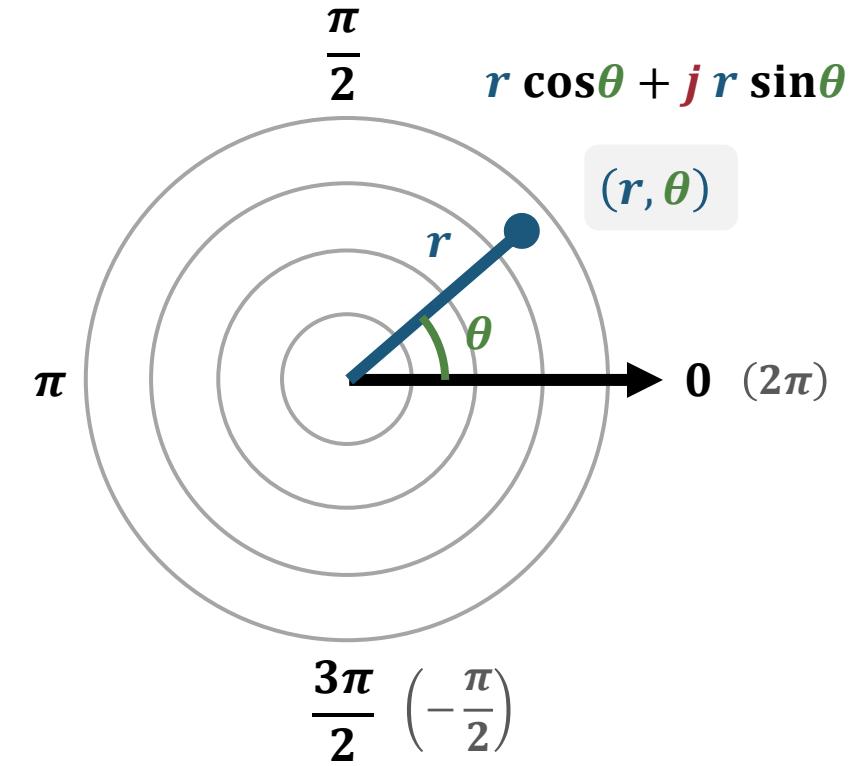
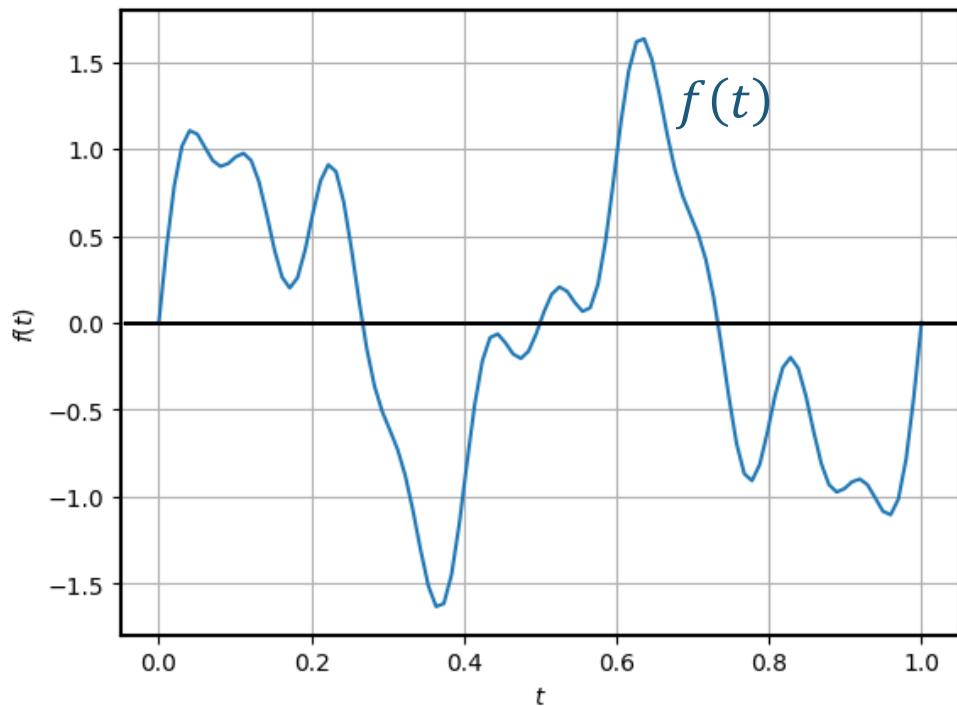
Real part **Imaginary part**



| Demystifying Fourier Transform

$$F(\omega) = \int_{-\infty}^{\infty} f(t) \cos(-\omega t) + j f(t) \sin(-\omega t) dt$$

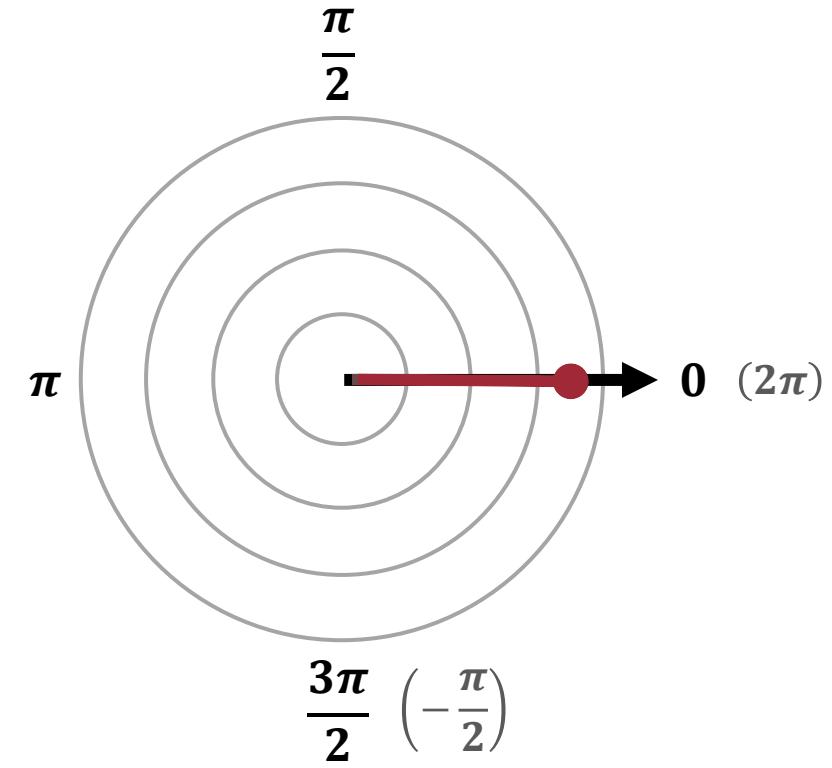
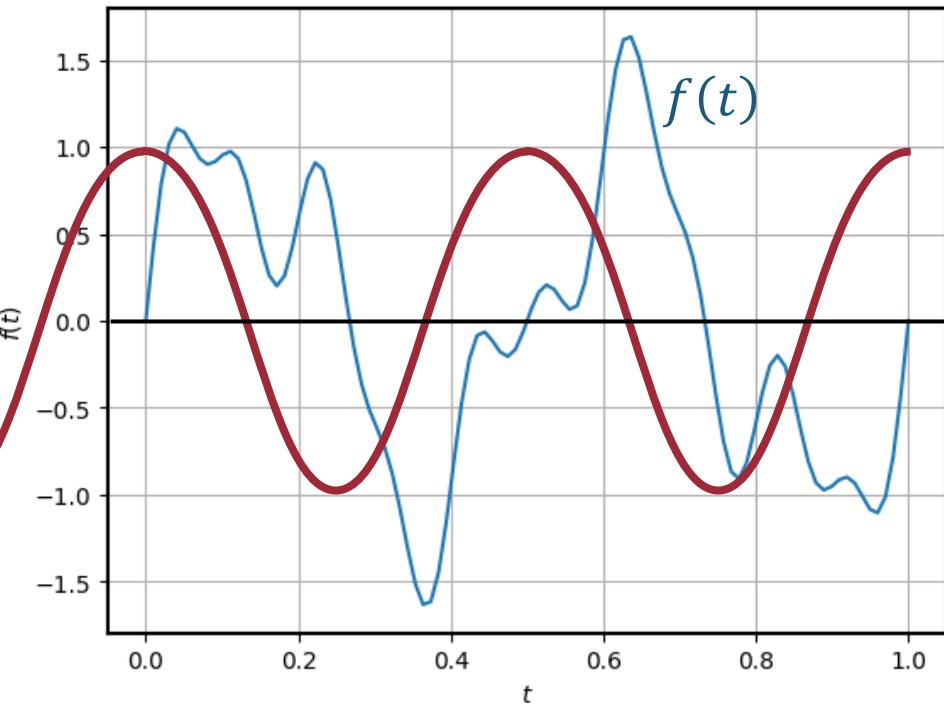
Real part Imaginary part



Demystifying Fourier Transform

$$F(\omega) = \int_{-\infty}^{\infty} f(t) \cos(-\omega t) + j f(t) \sin(-\omega t) dt$$

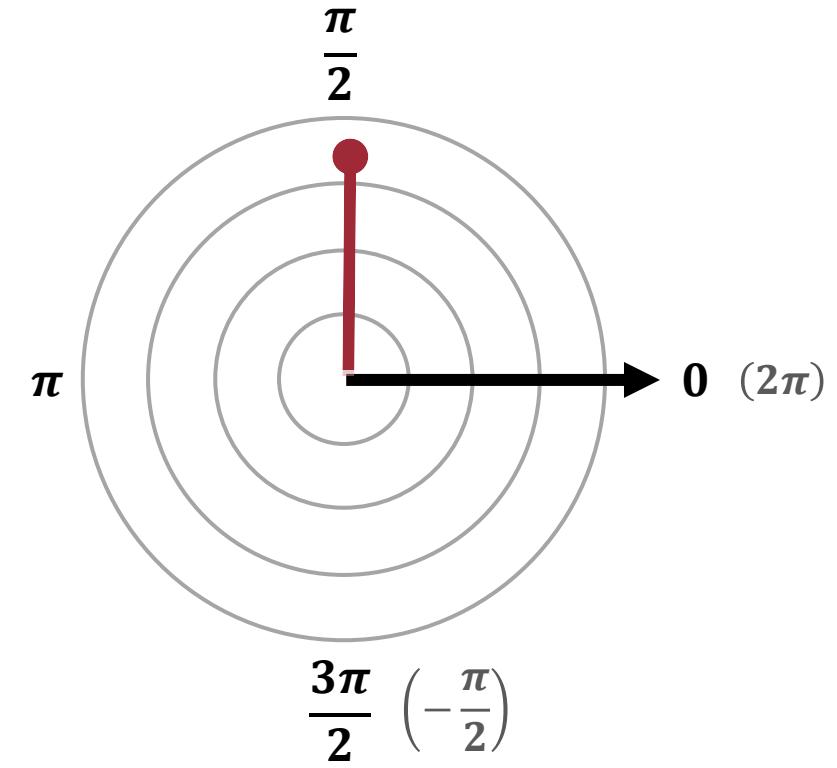
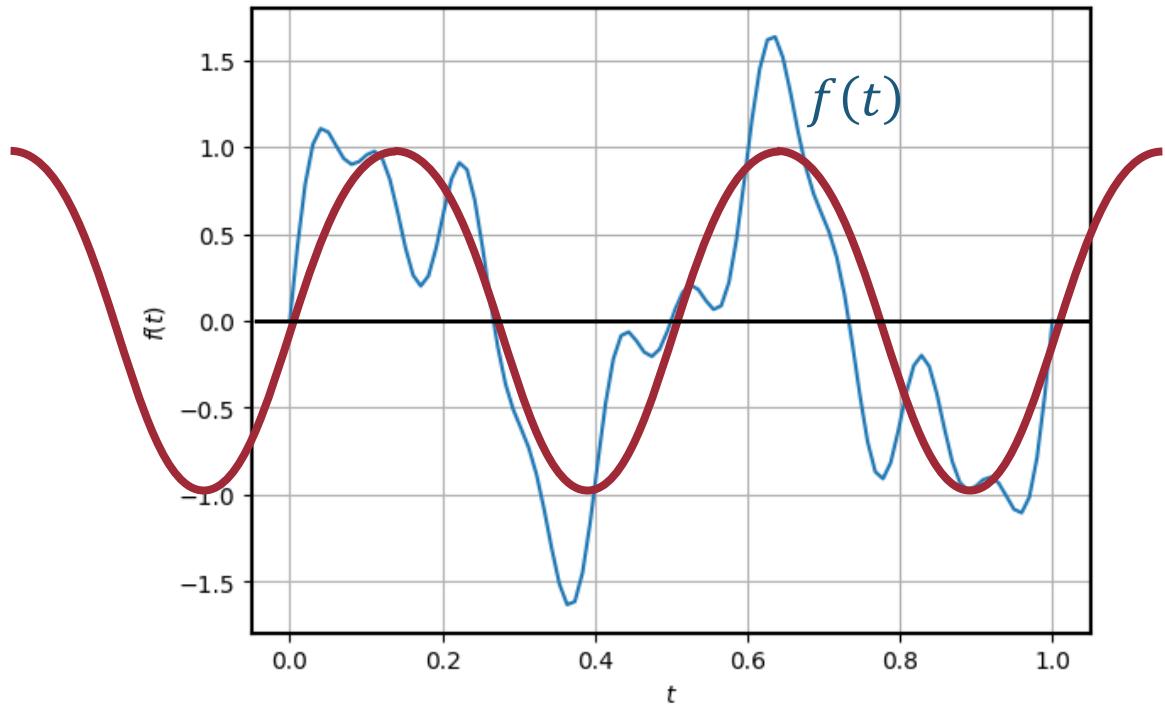
Real part **Imaginary part**



| Demystifying Fourier Transform

$$F(\omega) = \int_{-\infty}^{\infty} f(t) \cos(-\omega t) + j f(t) \sin(-\omega t) dt$$

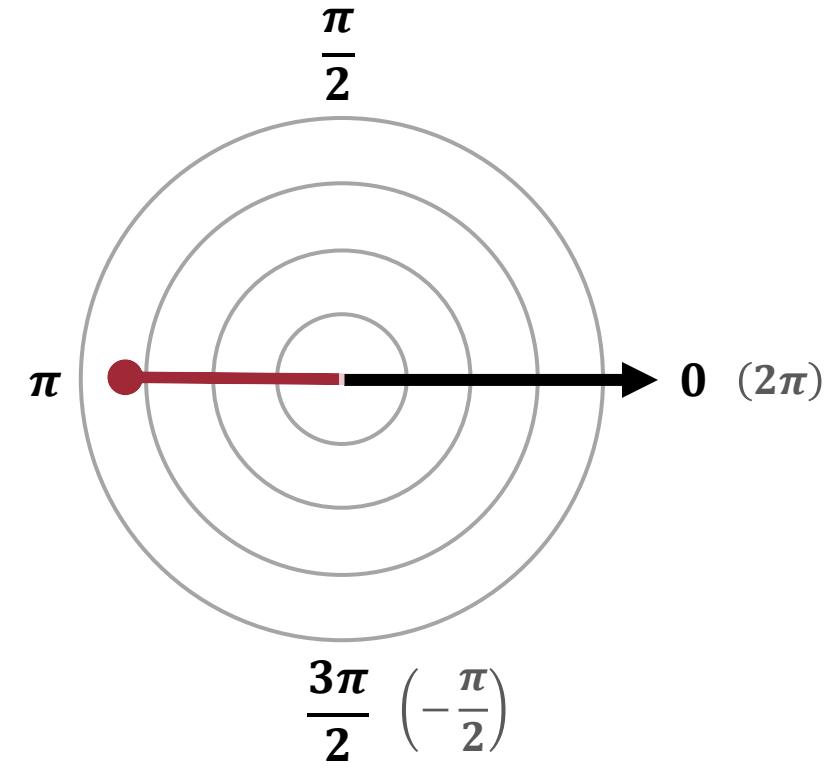
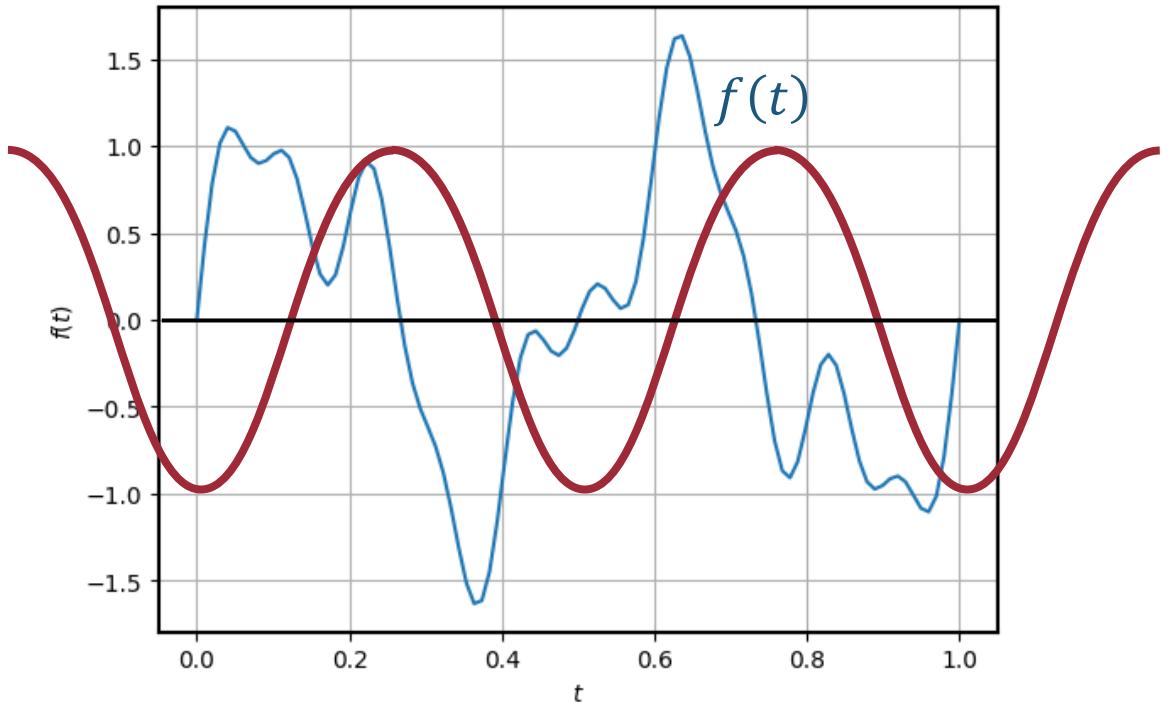
Real part **Imaginary part**



| Demystifying Fourier Transform

$$F(\omega) = \int_{-\infty}^{\infty} f(t) \cos(-\omega t) + j f(t) \sin(-\omega t) dt$$

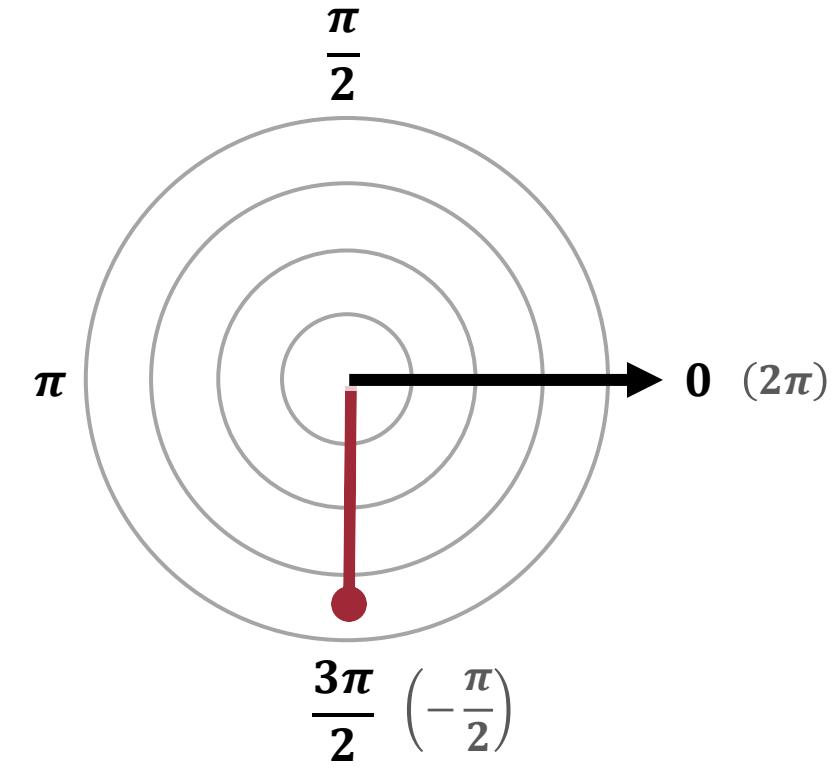
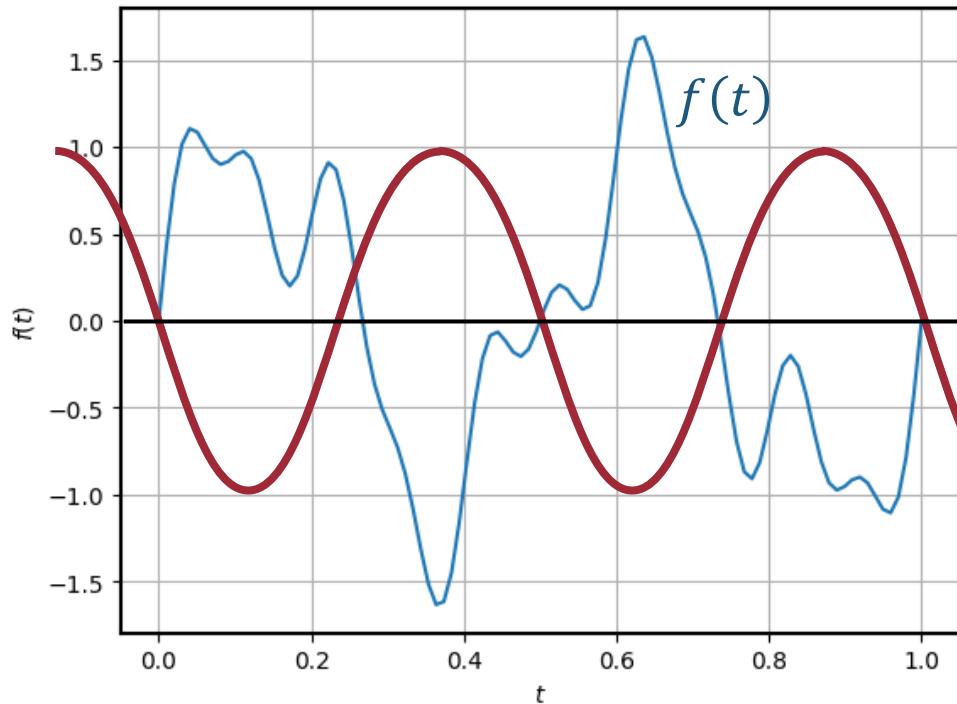
Real part Imaginary part



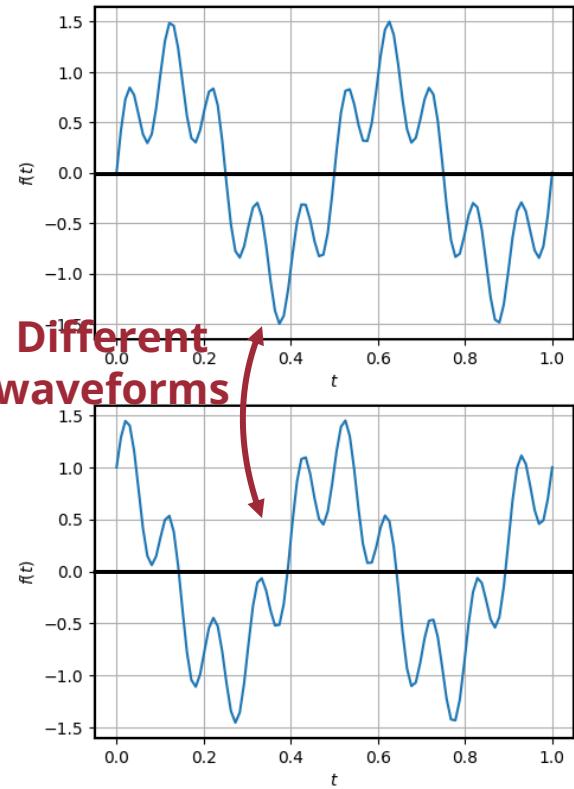
| Demystifying Fourier Transform

$$F(\omega) = \int_{-\infty}^{\infty} f(t) \cos(-\omega t) + j f(t) \sin(-\omega t) dt$$

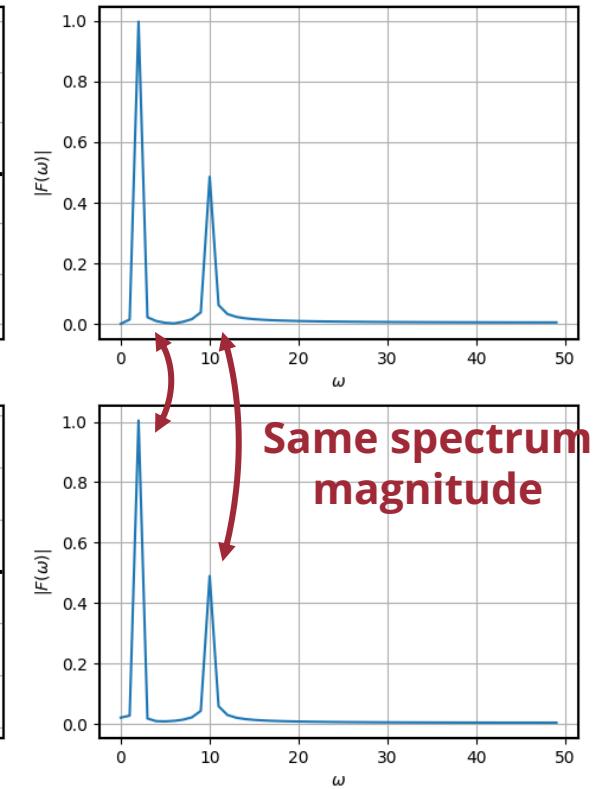
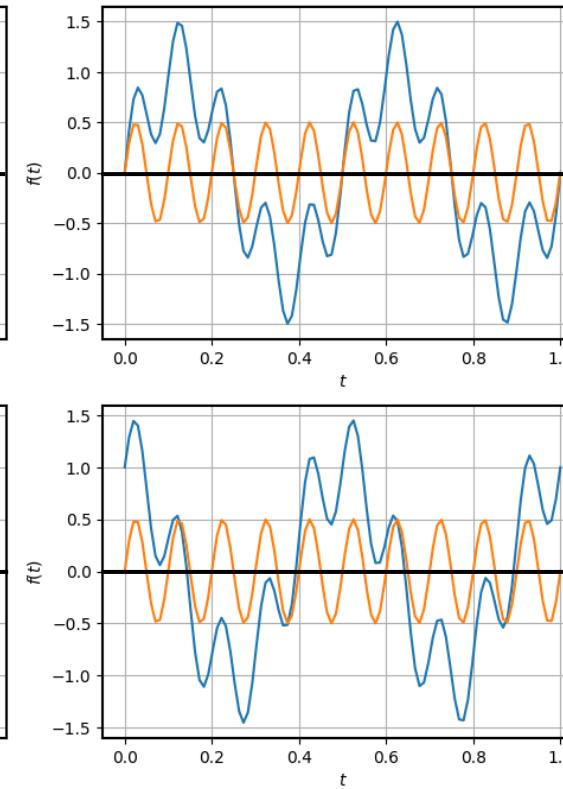
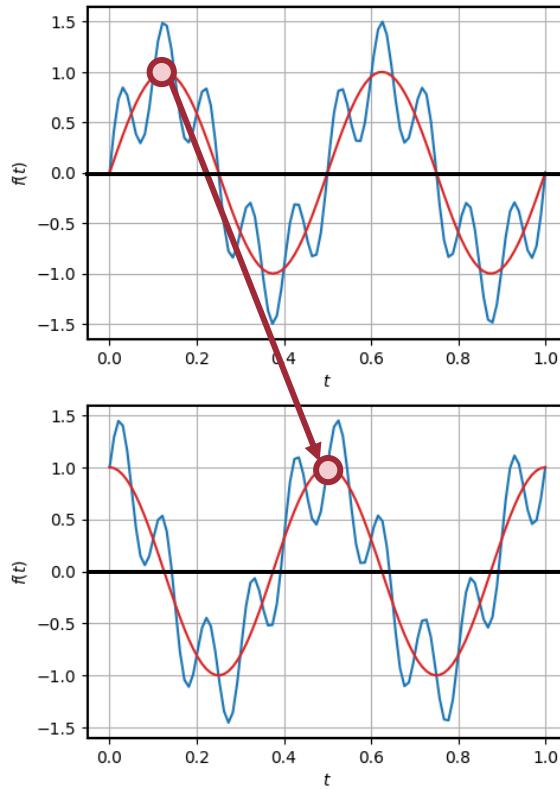
Real part **Imaginary part**



Magnitude & Phase



Different
waveforms

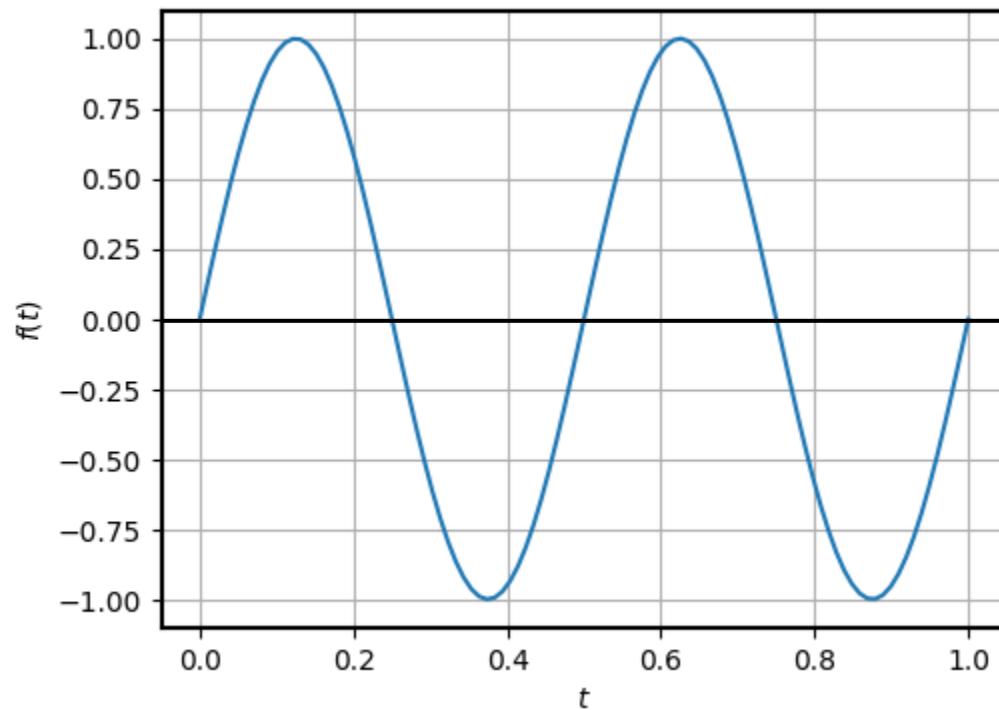


Same spectrum
magnitude

Example: A 2Hz Sine Wave

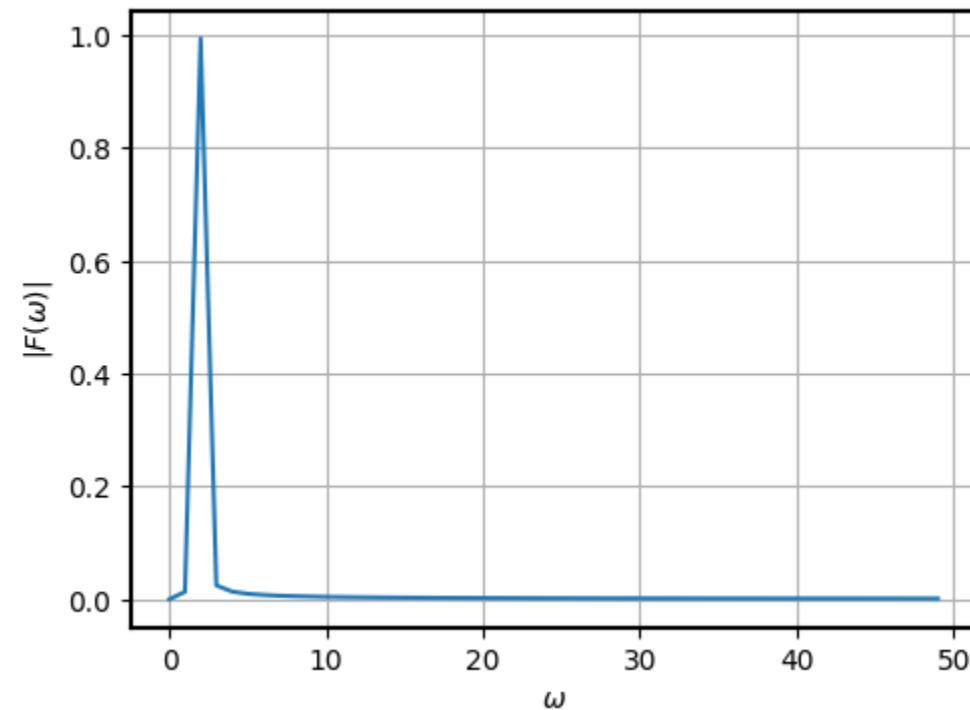
Signal

(time-domain)



Spectrum

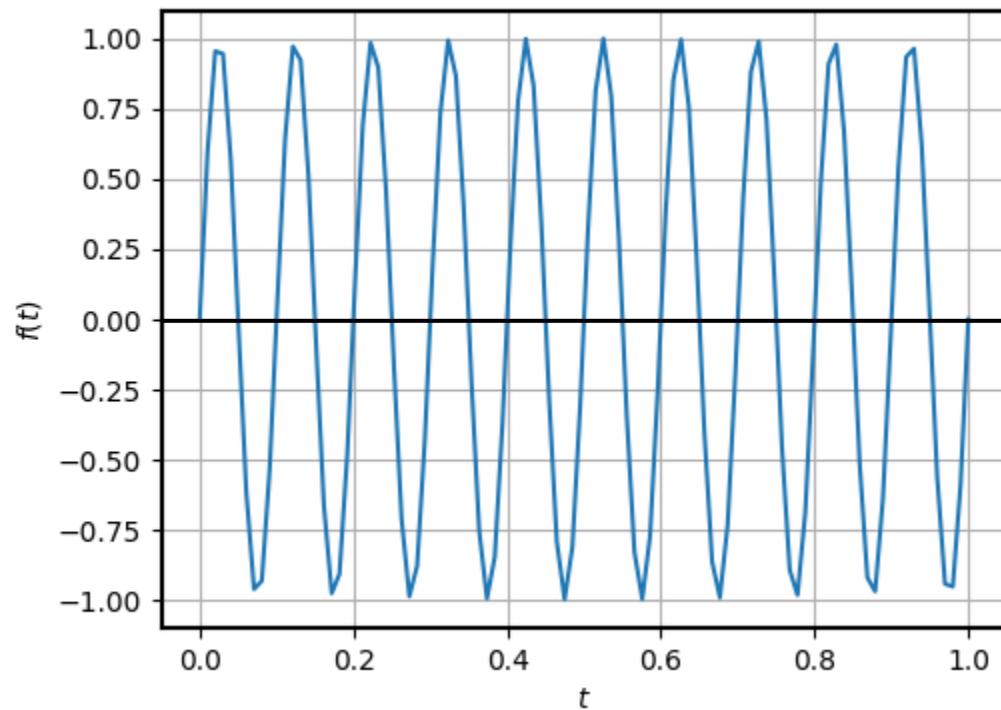
(frequency-domain)



| Example : A 10Hz Sine Wave

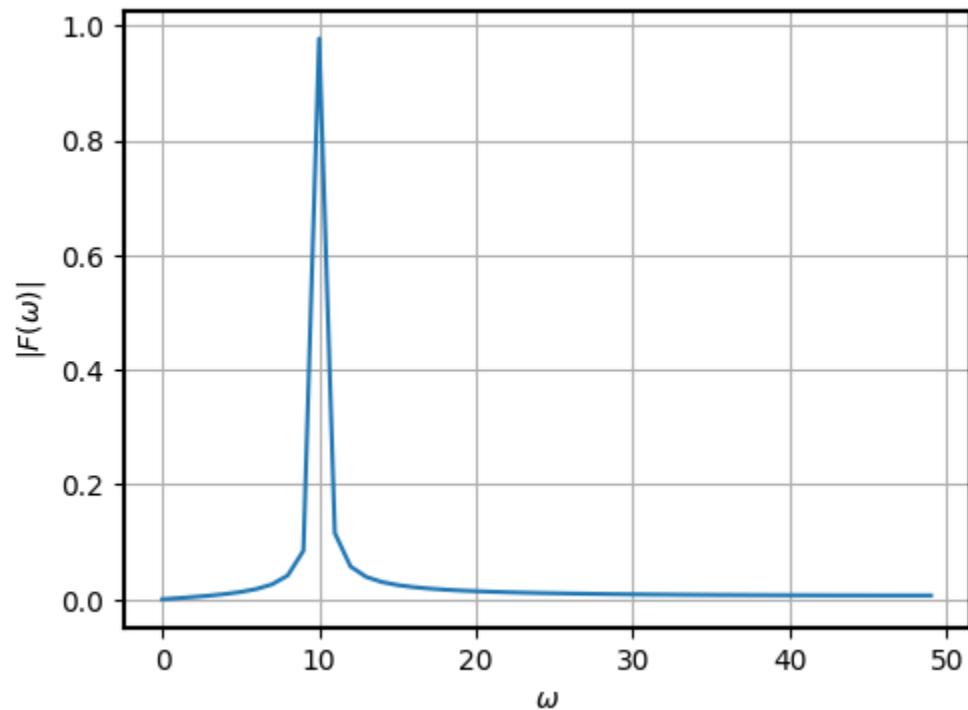
Signal

(time-domain)

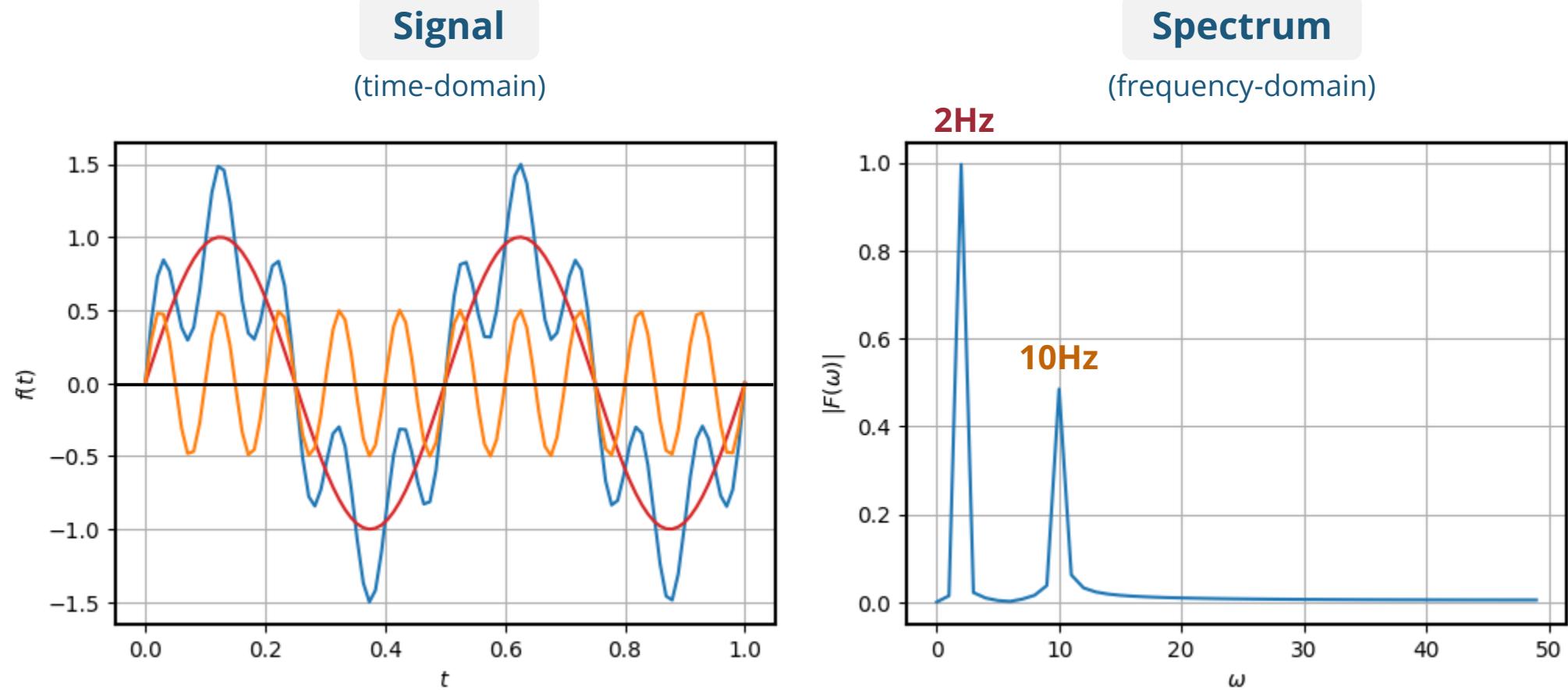


Spectrum

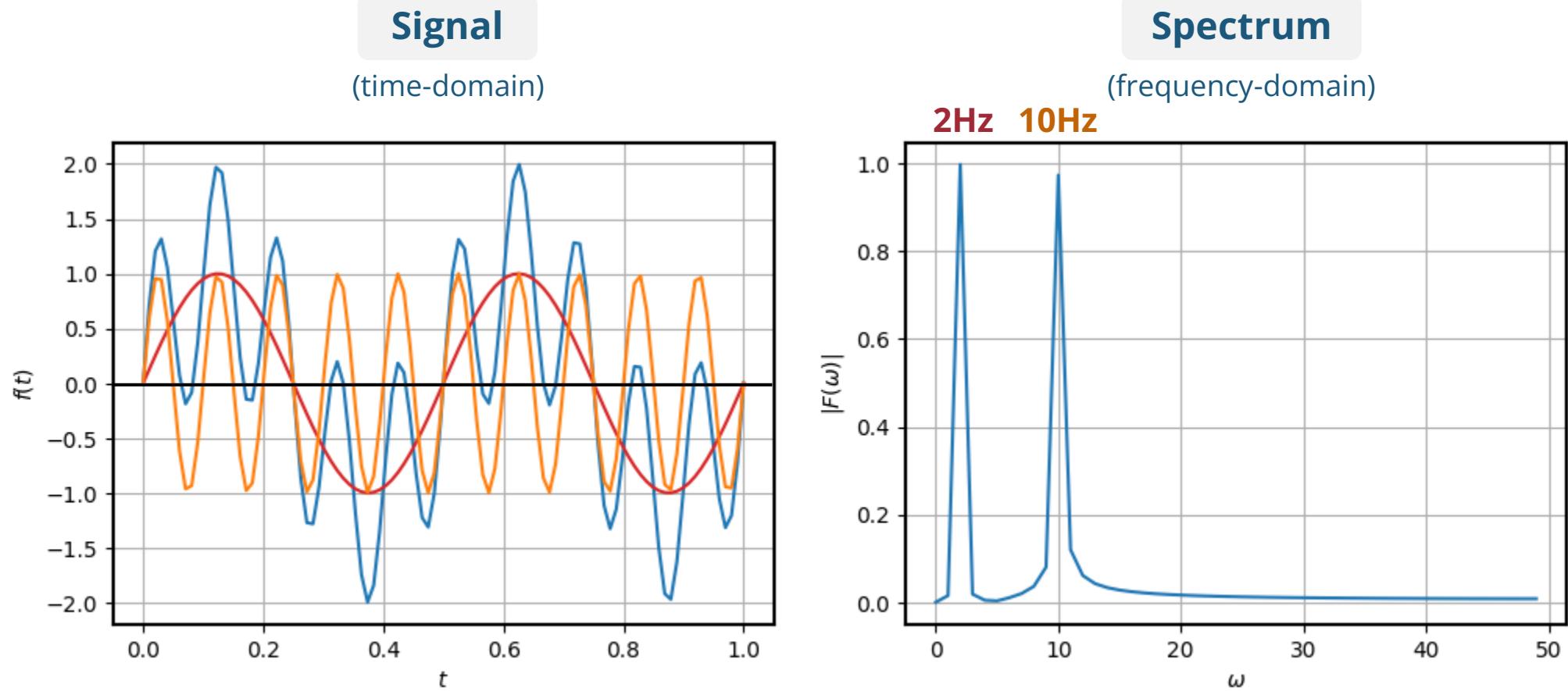
(frequency-domain)



| Example: Sum of a 10Hz & 2Hz Sine Waves



| How about this?



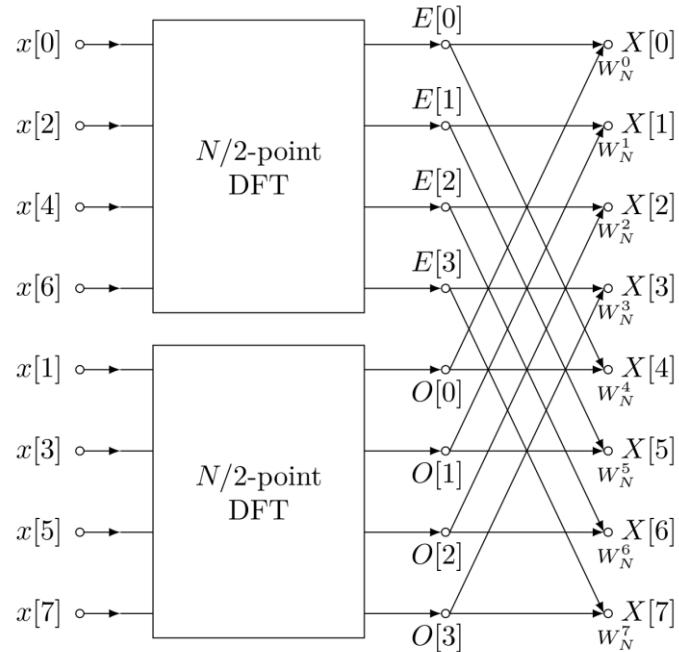
Discrete Fourier Transform (DFT)

- **Intuition:** Fourier transform with **discrete time and frequency**
 - Used for **digital audio** → we cannot achieve an infinite sampling rate...
- Math formulation:

$$X_k = \sum_{n=0}^{N-1} x_n e^{-j2\pi \frac{k}{N} n}$$

In Practice: Fast Fourier Transform (FFT)

- An efficient implementation of discrete Fourier transform
 - Reduce the complexity from $O(n^2)$ to $O(n \log n)$



(Source: Yangwenbo99 via Wikimedia)

Top 10 algorithms from the 20th century

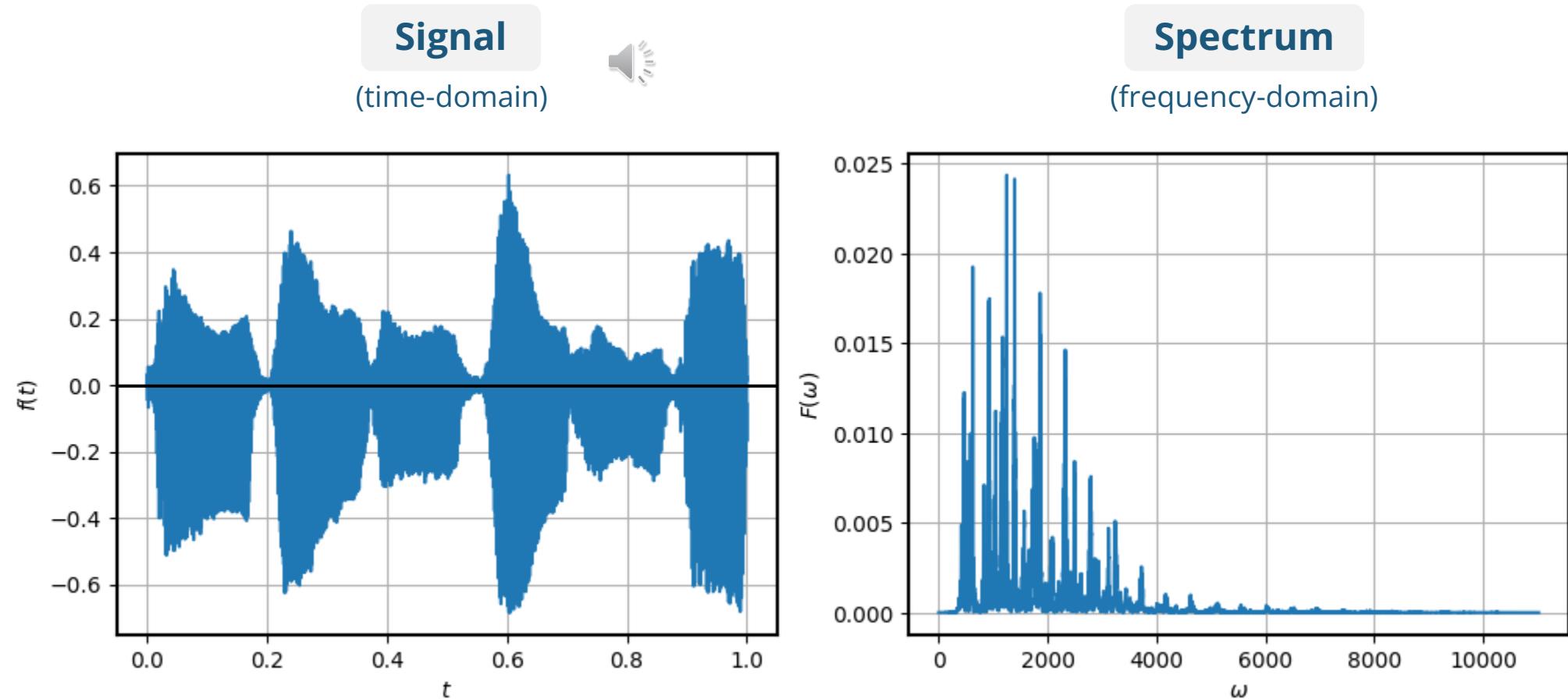
Computing
in SCIENCE & ENGINEERING



IEEE
COMPUTER SOCIETY
www.computer.org/cise

Time-Frequency Analysis

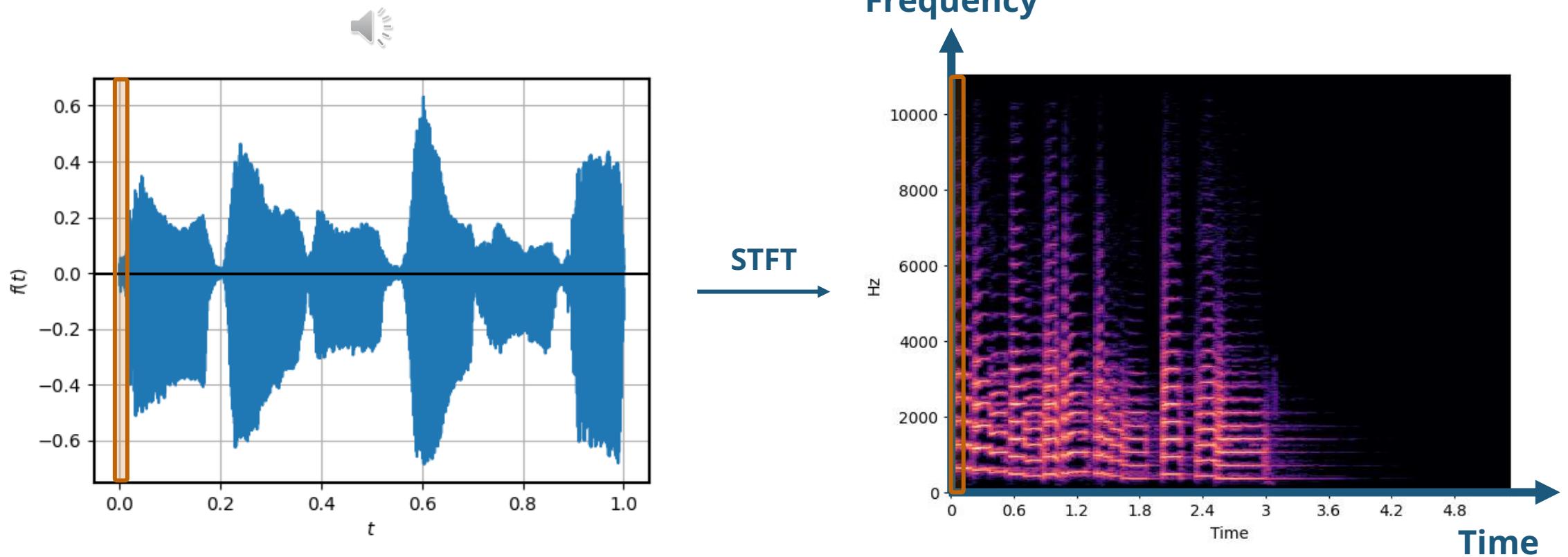
Fourier Transform of a Trumpet Sound



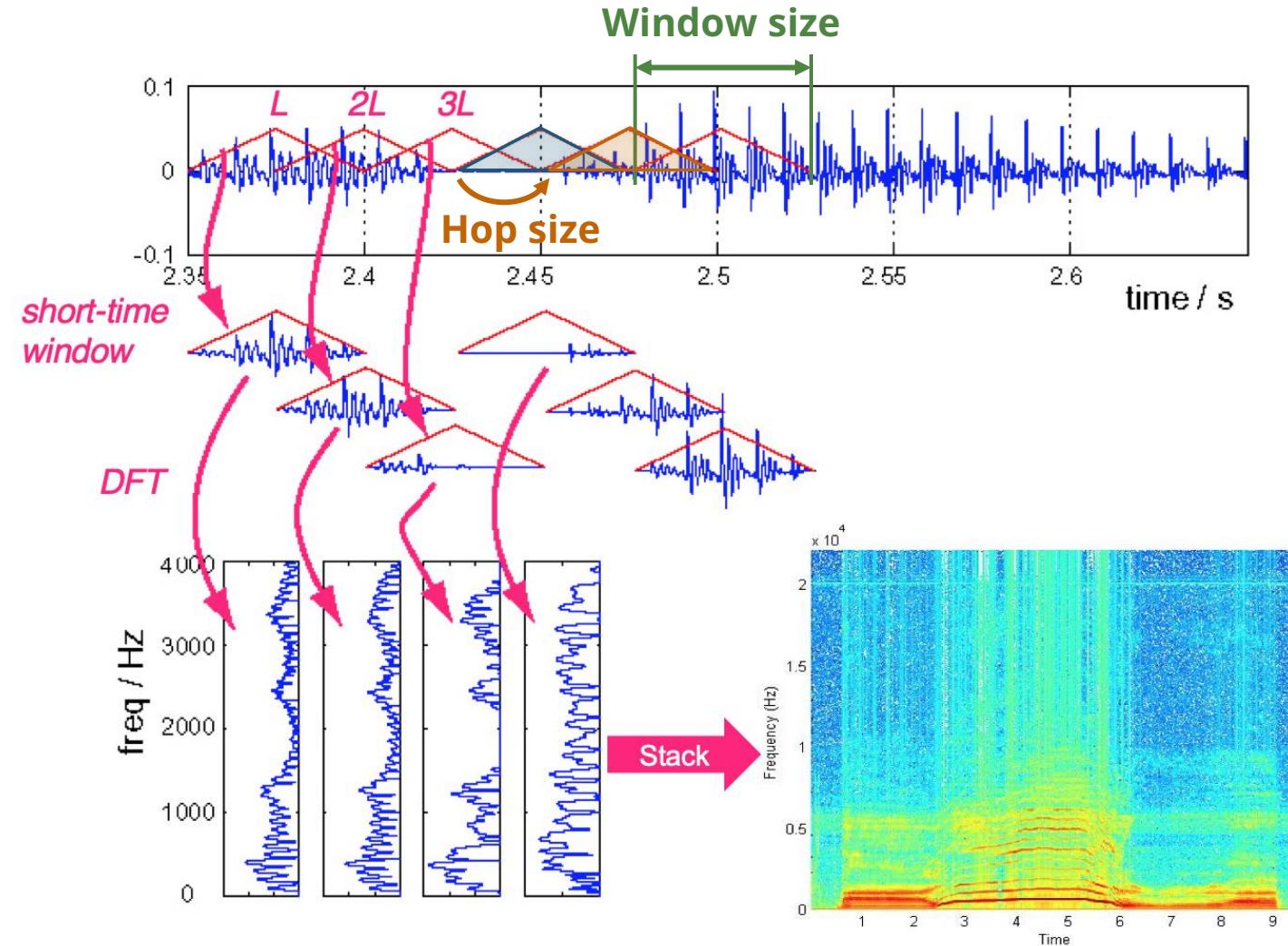
Fourier Transform cannot localize! 😰

| Short-Time Fourier Transform (STFT)

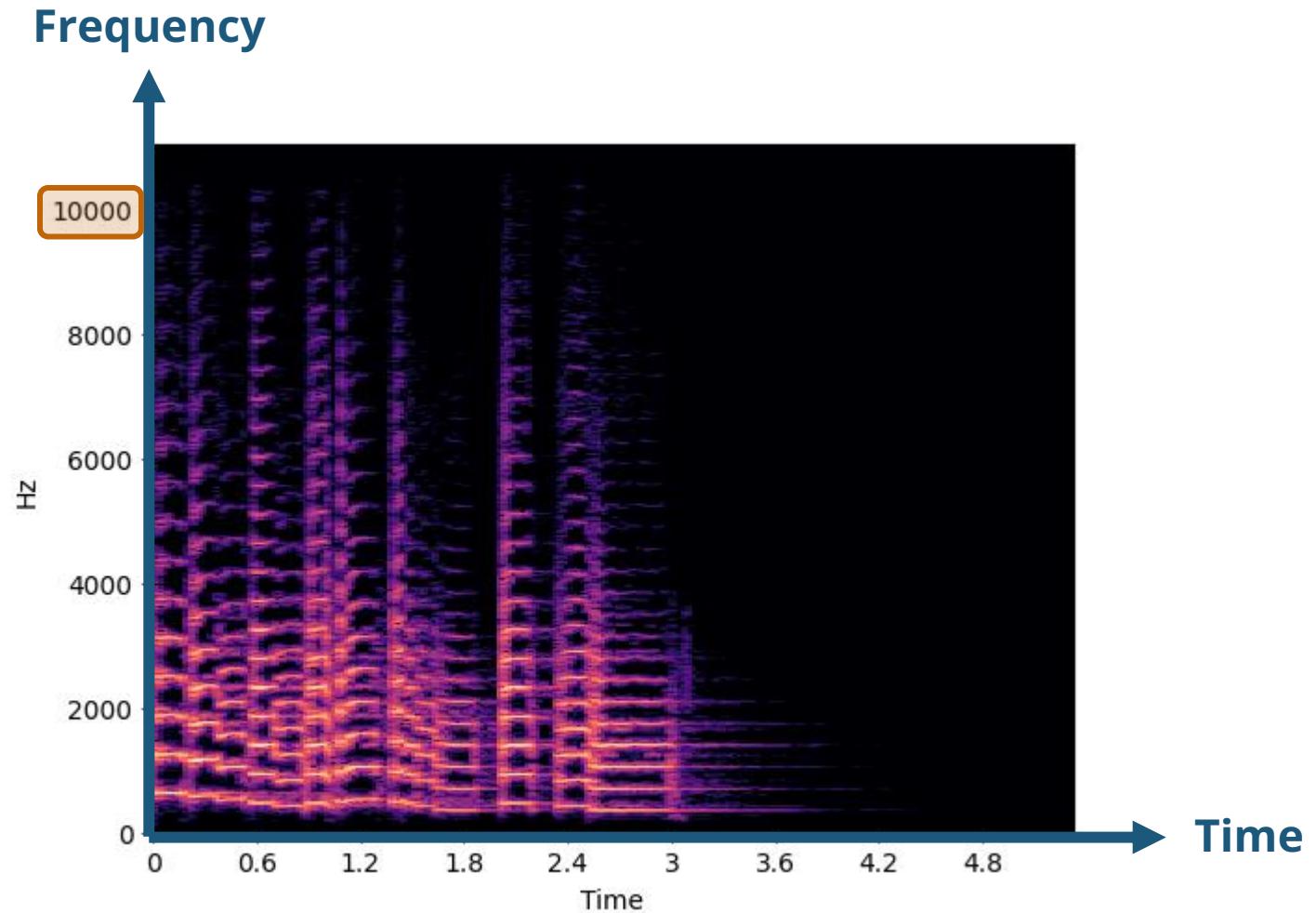
- **Intuition:** Slice the audio into chunks and apply Fourier transform



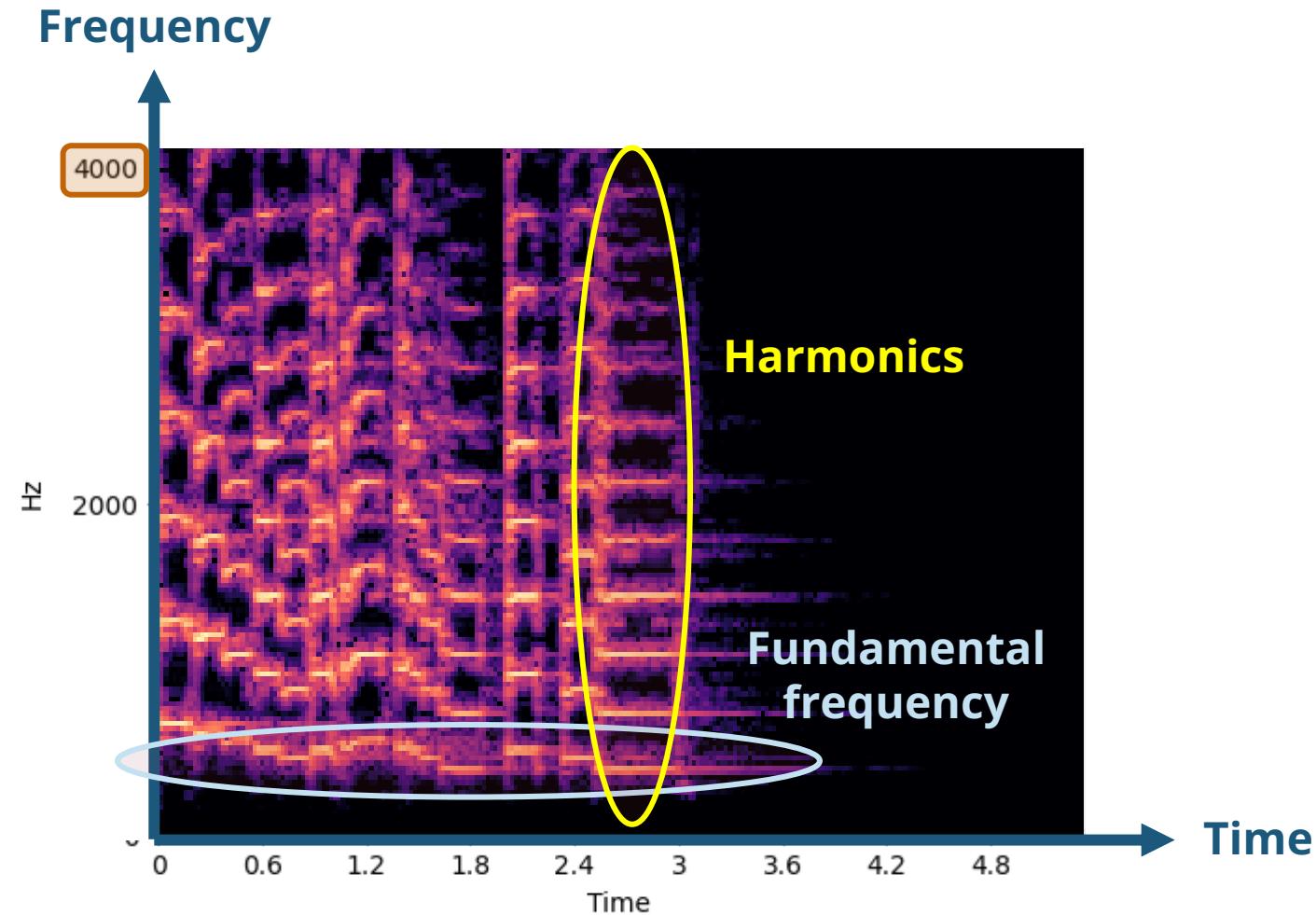
Short-Time Fourier Transform (STFT)



Spectrogram



Spectrogram



Timbre

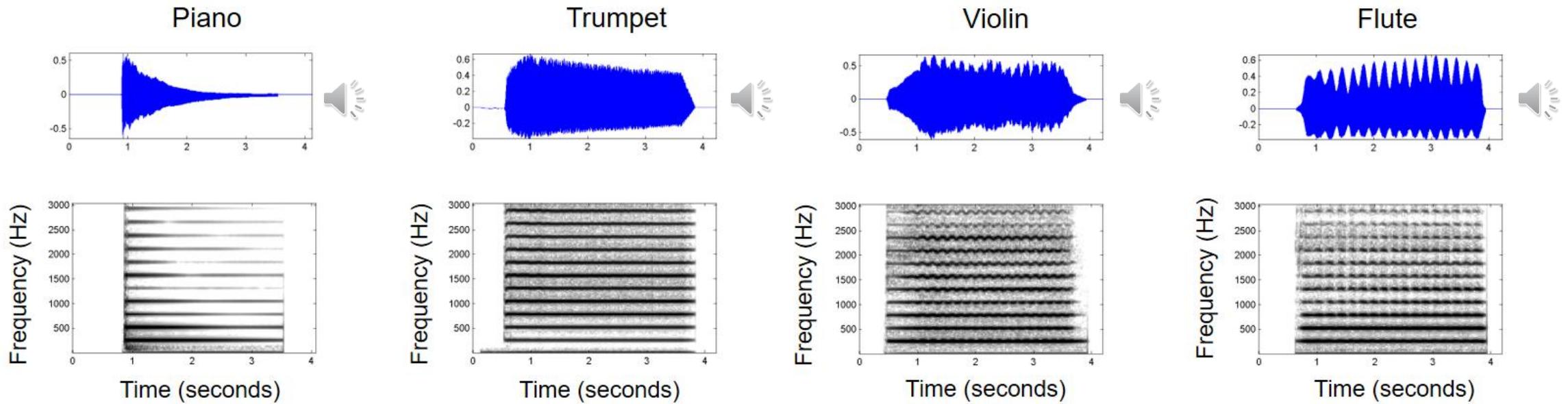
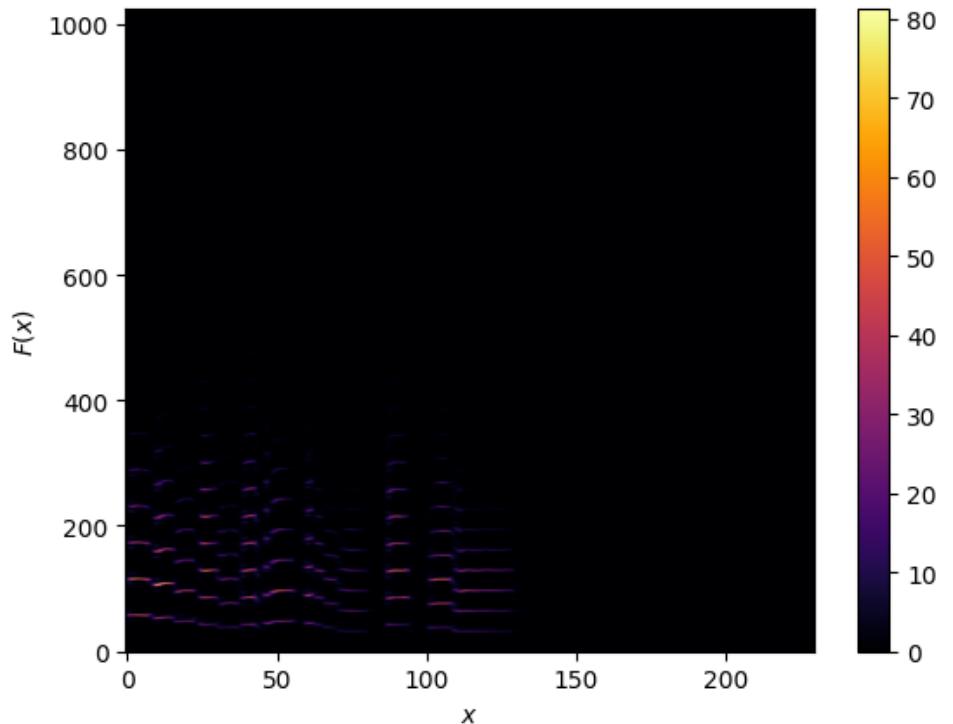


Figure 1.23 from [Müller, FMP, Springer 2015]

(Source: Müller et al., 2021)

| Example: `librosa.stft`

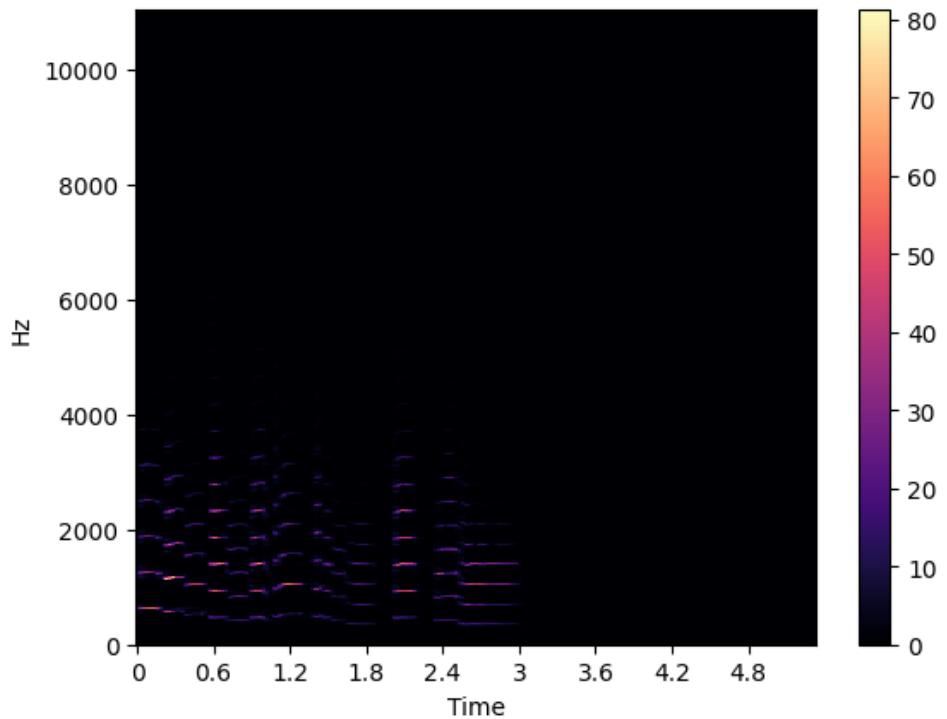


```
# Load the example audio in librosa
y, sr = librosa.load(librosa.example("trumpet"))

# Compute the spectrogram
S = np.abs(librosa.stft(y))

# Plot the spectrogram
im = plt.imshow(S, cmap="inferno", aspect="auto",
                 origin="lower")
plt.colorbar(im)
plt.xlabel("Time (sec)")
plt.ylabel("Frequency (Hz)")
plt.show()
```

| Example: `librosa.display.specshow`

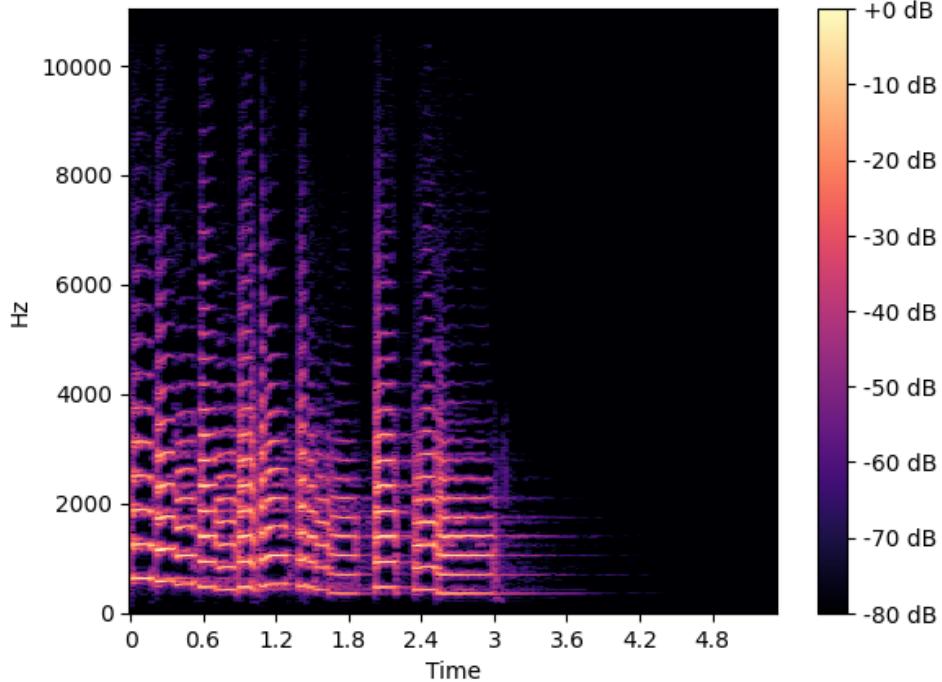


```
# Load the example audio in librosa
y, sr = librosa.load(librosa.example("trumpet"))

# Compute the spectrogram
S = np.abs(librosa.stft(y))

# Plot the spectrogram
im = librosa.display.specshow(S, x_axis="time",
                               y_axis="linear")
plt.colorbar(im)
plt.show()
```

| Example: `librosa.amplitude_to_db`



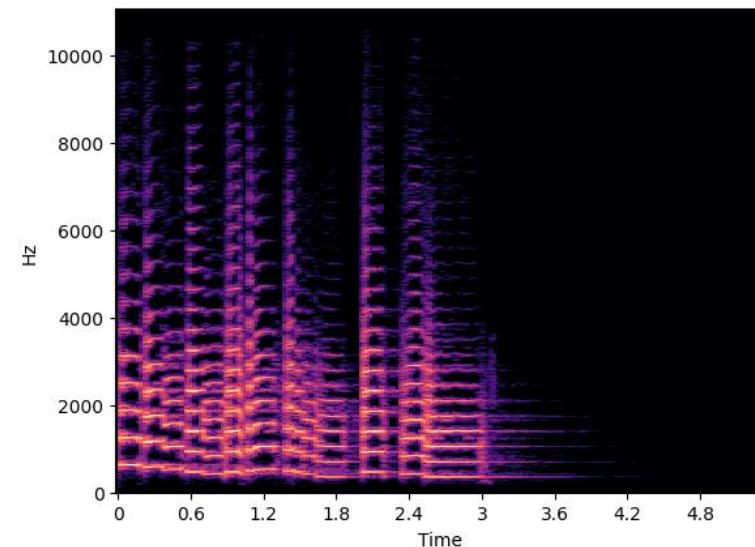
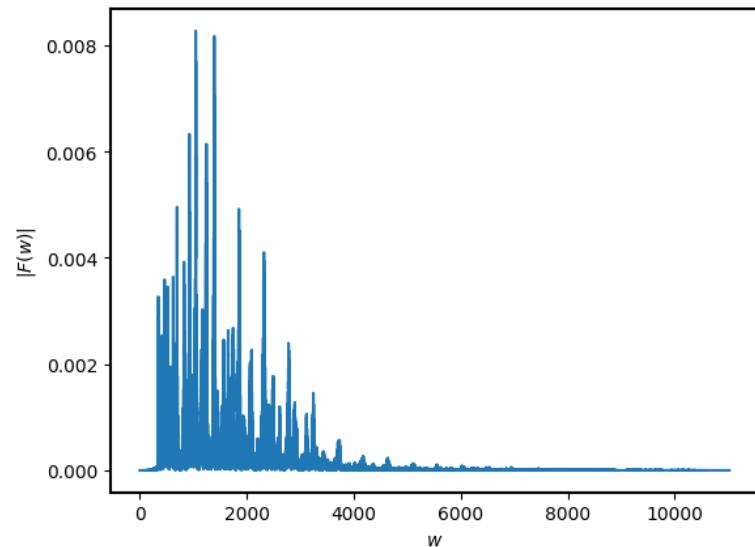
```
# Load the example audio in librosa
y, sr = librosa.load(librosa.example("trumpet"))

# Compute the spectrogram
S = np.abs(librosa.stft(y))
S_db = librosa.amplitude_to_db(S, ref=np.max)

# Plot the spectrogram
im = librosa.display.specshow(S_db, x_axis="time",
                               y_axis="linear")
plt.colorbar(im, format="%+2.0f dB")
plt.show()
```

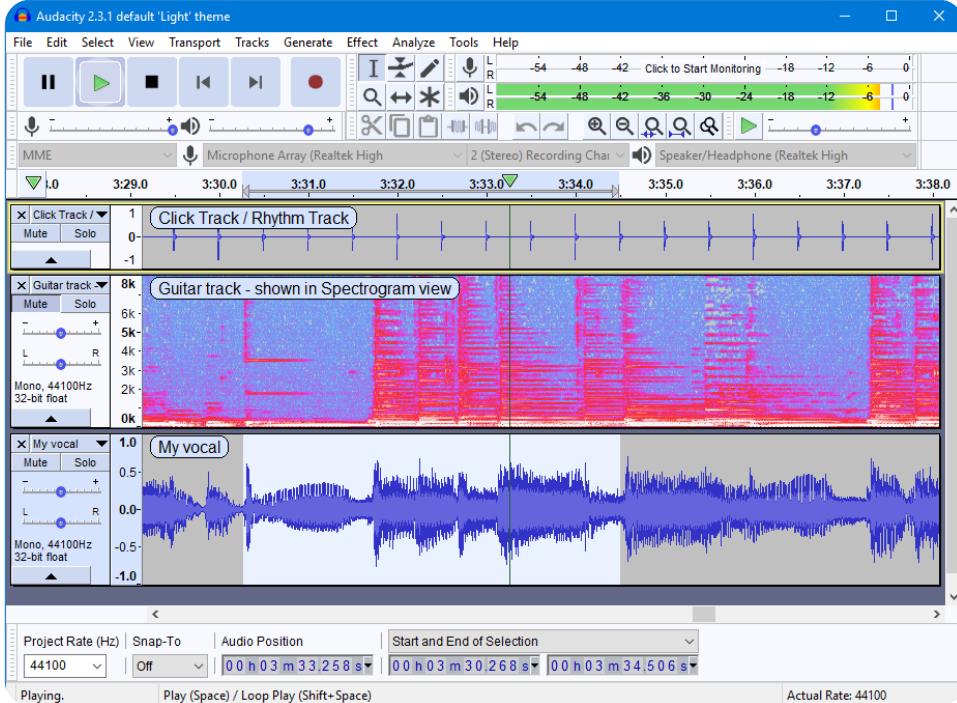
| 🔥 PA2: Spectral Analysis

- Instructions will be sent by **emails** and released on the **course website**
- Use **librosa** to process audio files
 - Fast Fourier transform (**FFT**)
 - Short-time Fourier transform (**STFT**)



Software

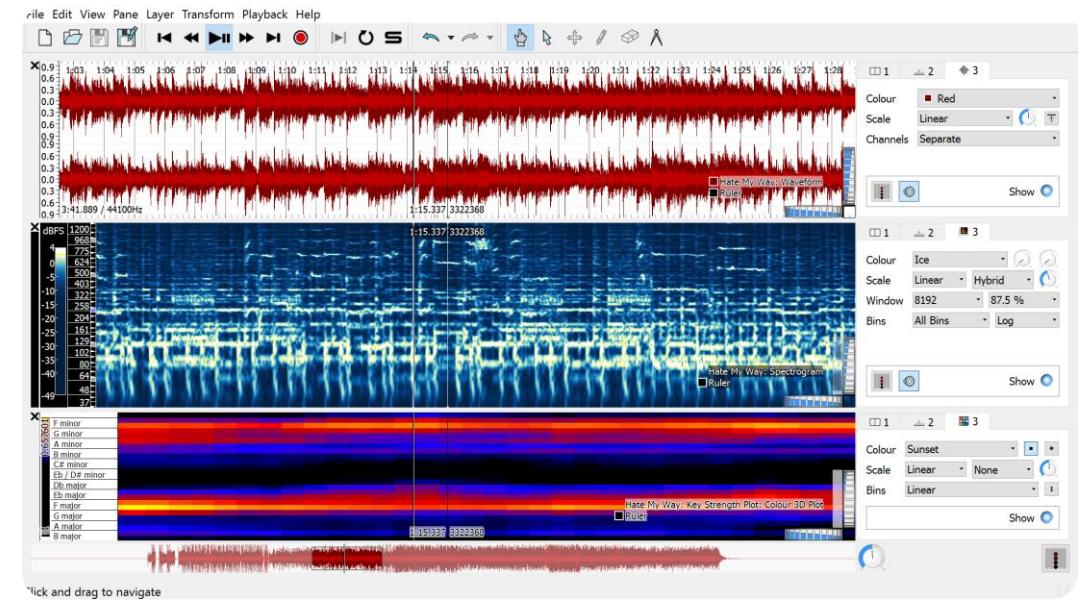
Audacity



(Source: audacity-2.3.1 via Internet Archive)

archive.org/details/audacity-2.3.1
sonicvisualiser.org

Sonic Visualiser



(Source: sonicvisualizer.org)

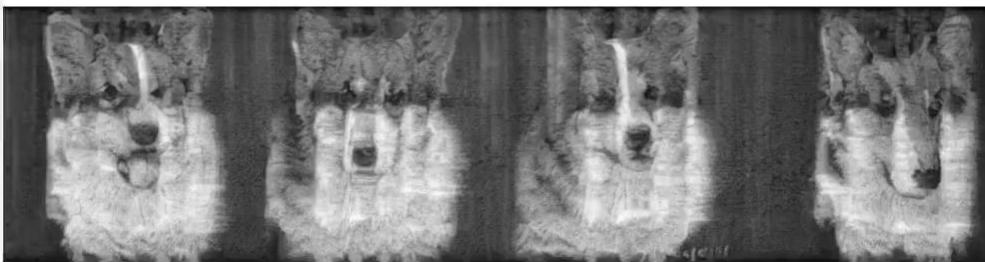
Images that Sound (Chen et al., 2024)

Using diffusion models to generate visual spectrograms that look like images but can also be played as sound.

Image prompt: a colorful photo of corgis



Audio prompt: dog barking



(Source: Chen et al., 2024)

Image prompt: a colorful photo of tigers



Audio prompt: tiger growling



(Source: Chen et al., 2024)

Images that Sound (Chen et al., 2024)

Using diffusion models to generate visual spectrograms that look like images but can also be played as sound.

Image prompt: a colorful photo of an auto racing game



Audio prompt: a race car passing by and disappearing

(Source: Chen et al., 2024)

Image prompt: a colorful photo of a castle with bell towers

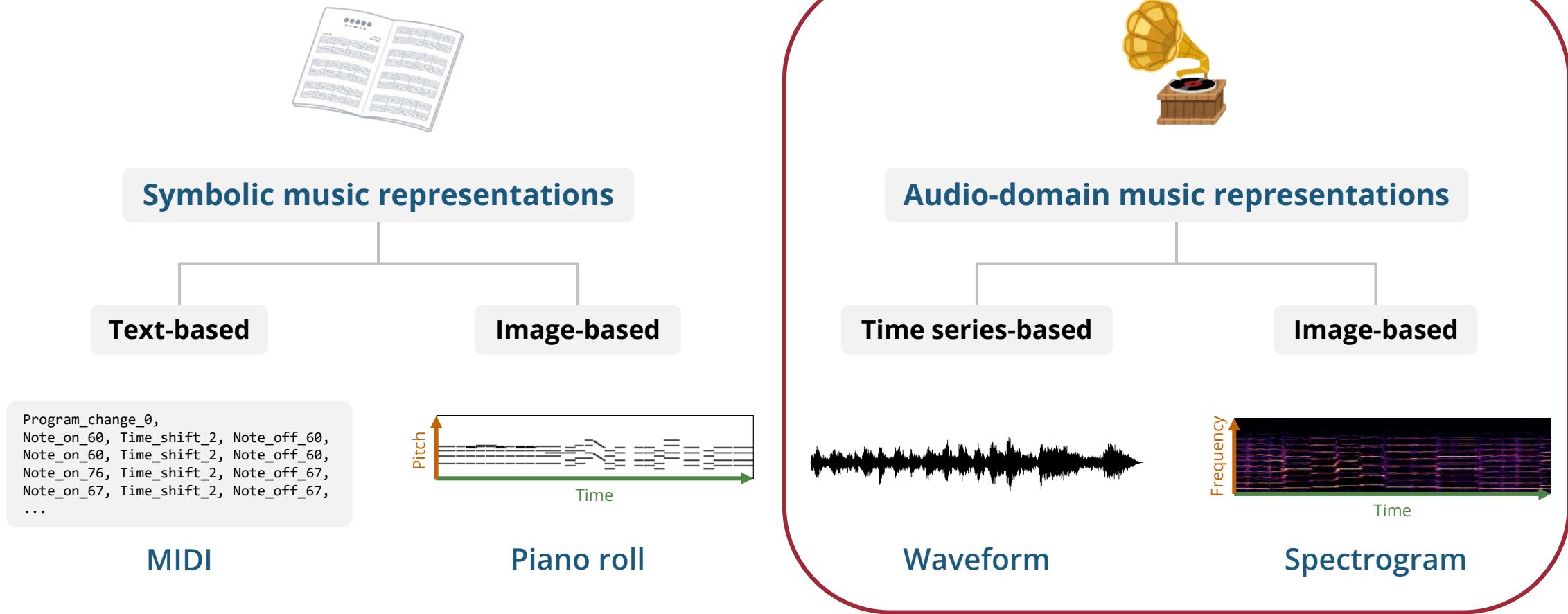


Audio prompt: bell ringing

(Source: Chen et al., 2024)

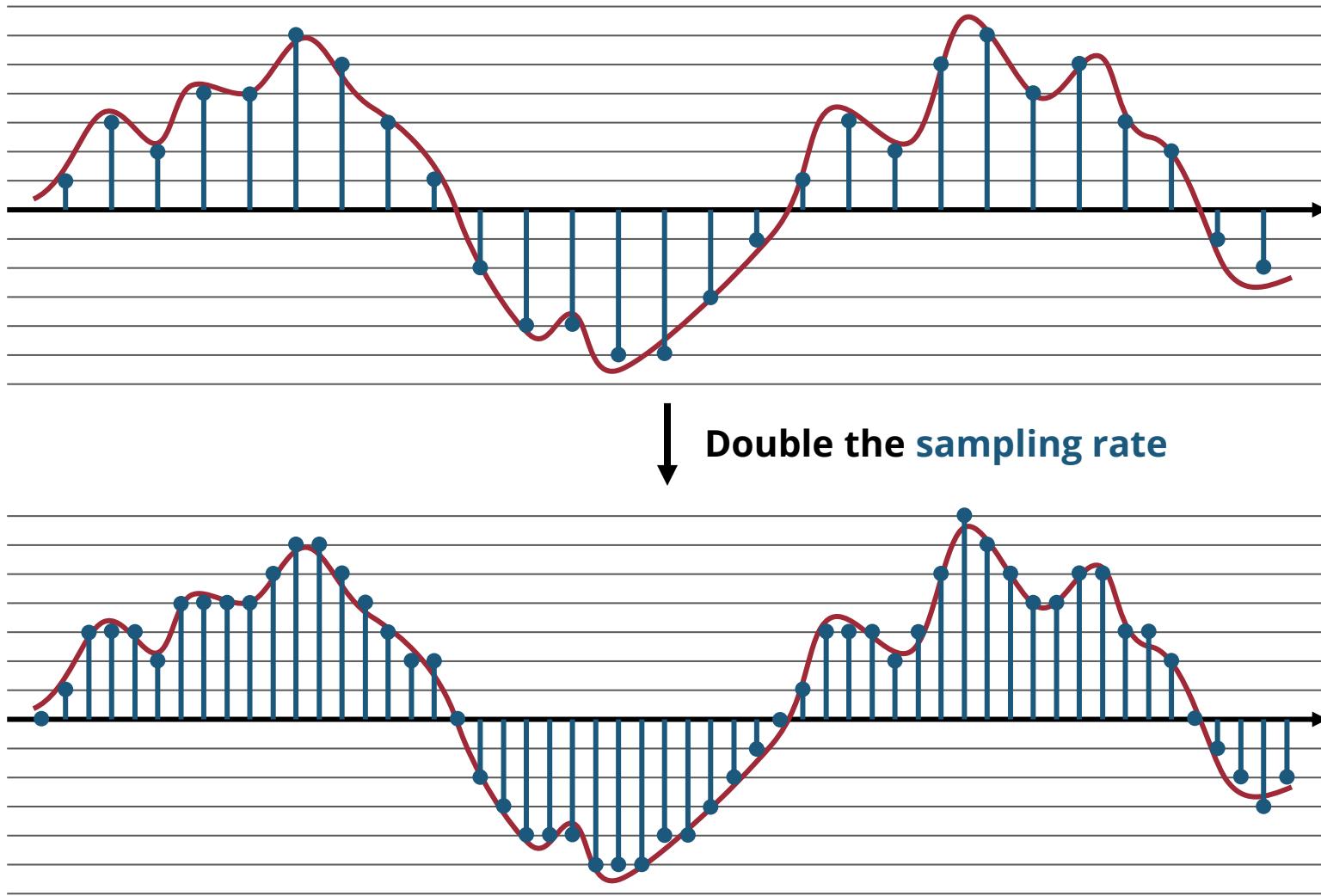
Recap

Four Representative Music Representations

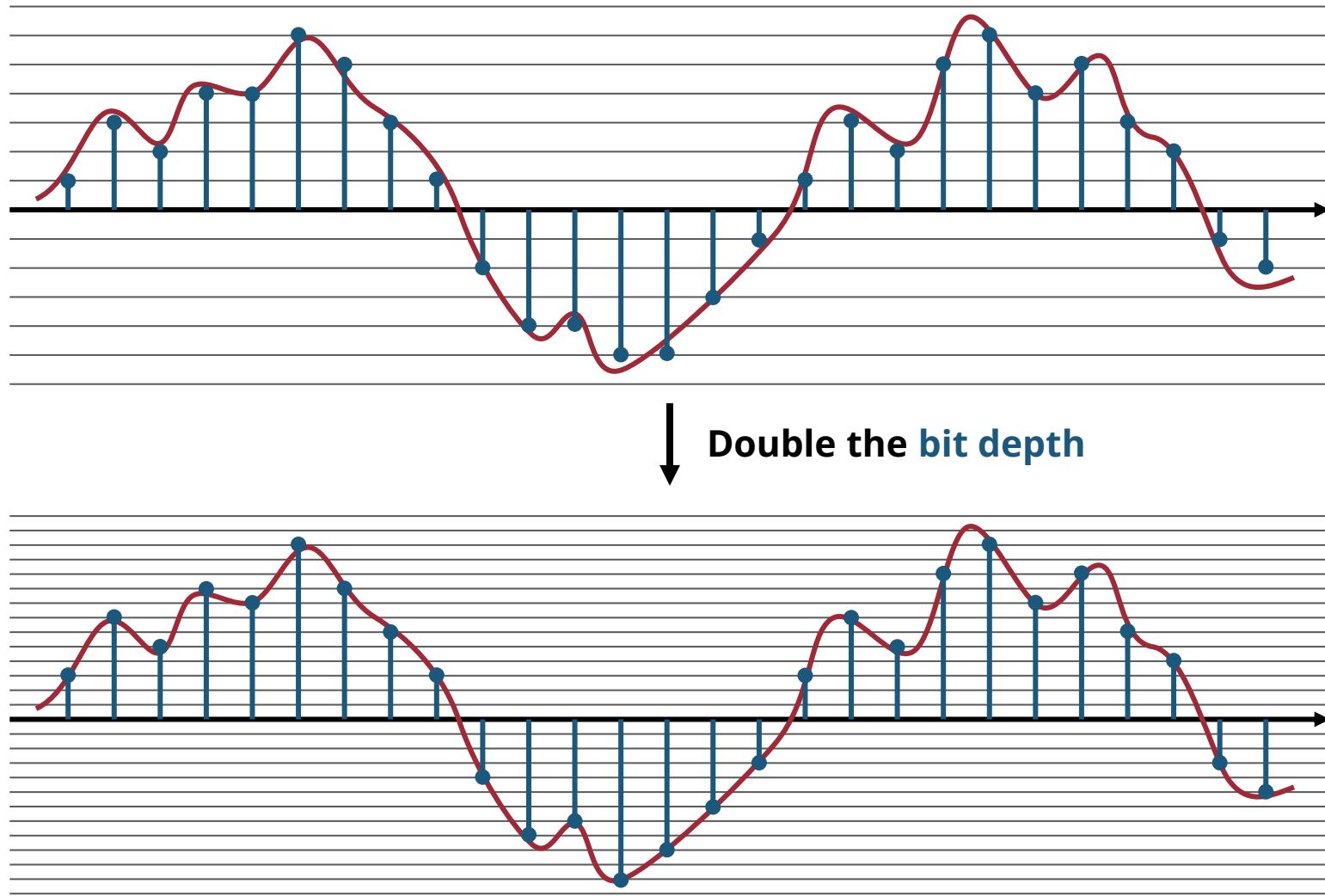


Today's topic!

Resolution: Sampling Rate

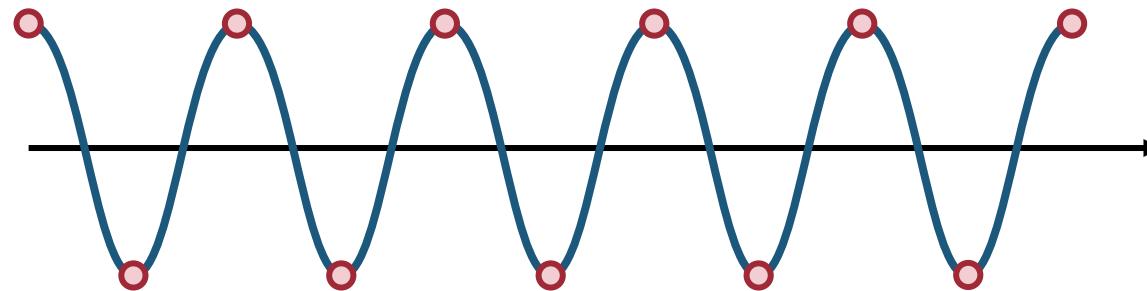


Resolution: Bit Depth

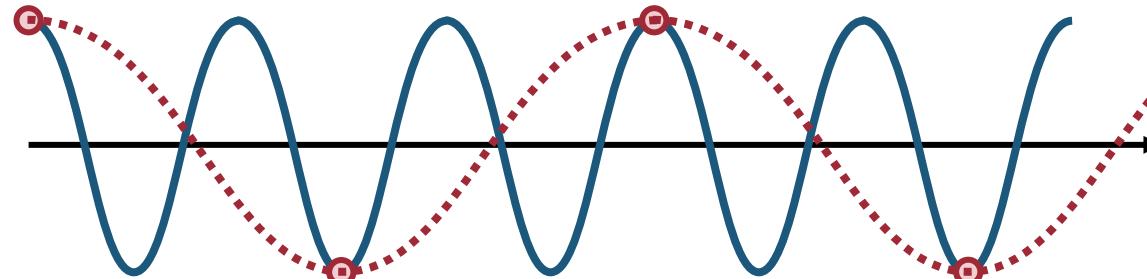


Sampling Theorem: Undersampling

Critically sampled
 $(f_s = 2f_{max})$

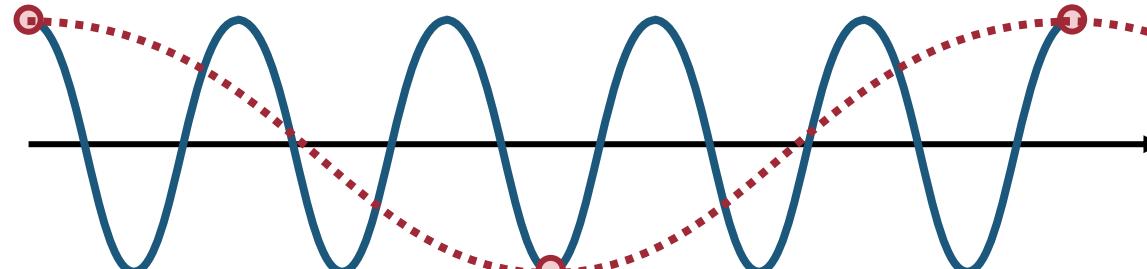


Undersampled
 $(f_s = \frac{2}{3}f_{max})$



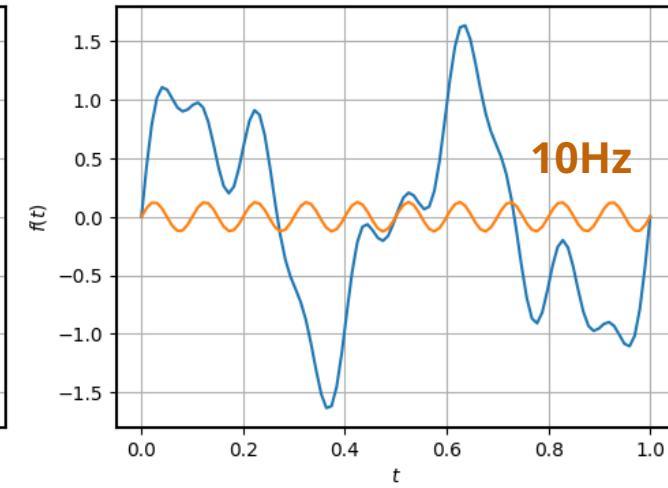
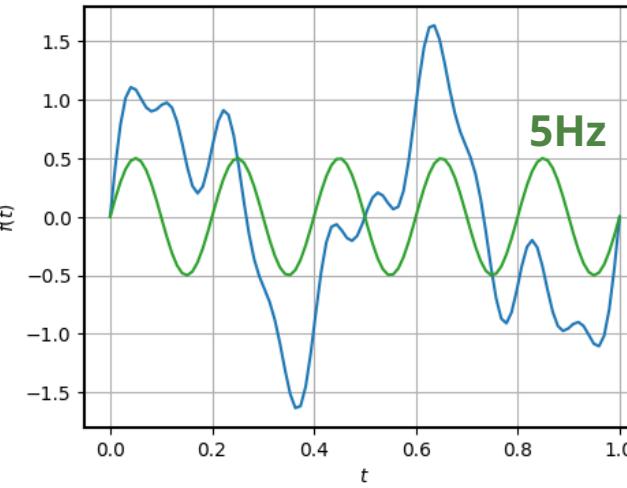
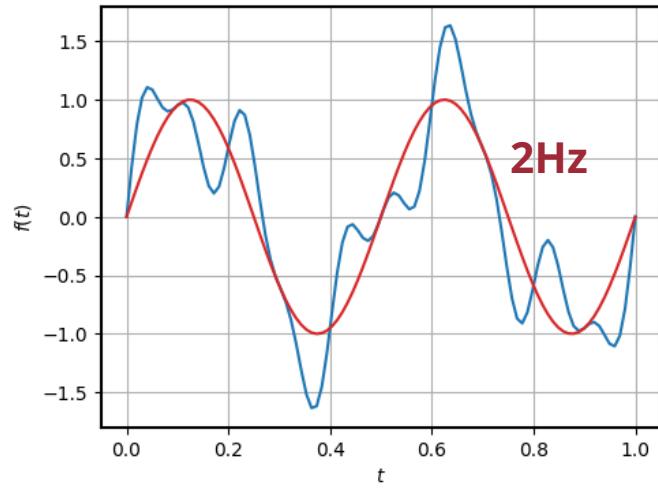
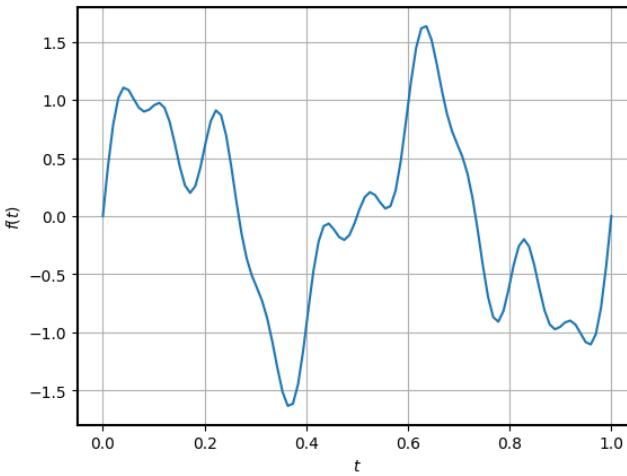
Can only reconstruct frequency up to $\frac{1}{3}f_{max}$

Undersampled
 $(f_s = \frac{2}{5}f_{max})$

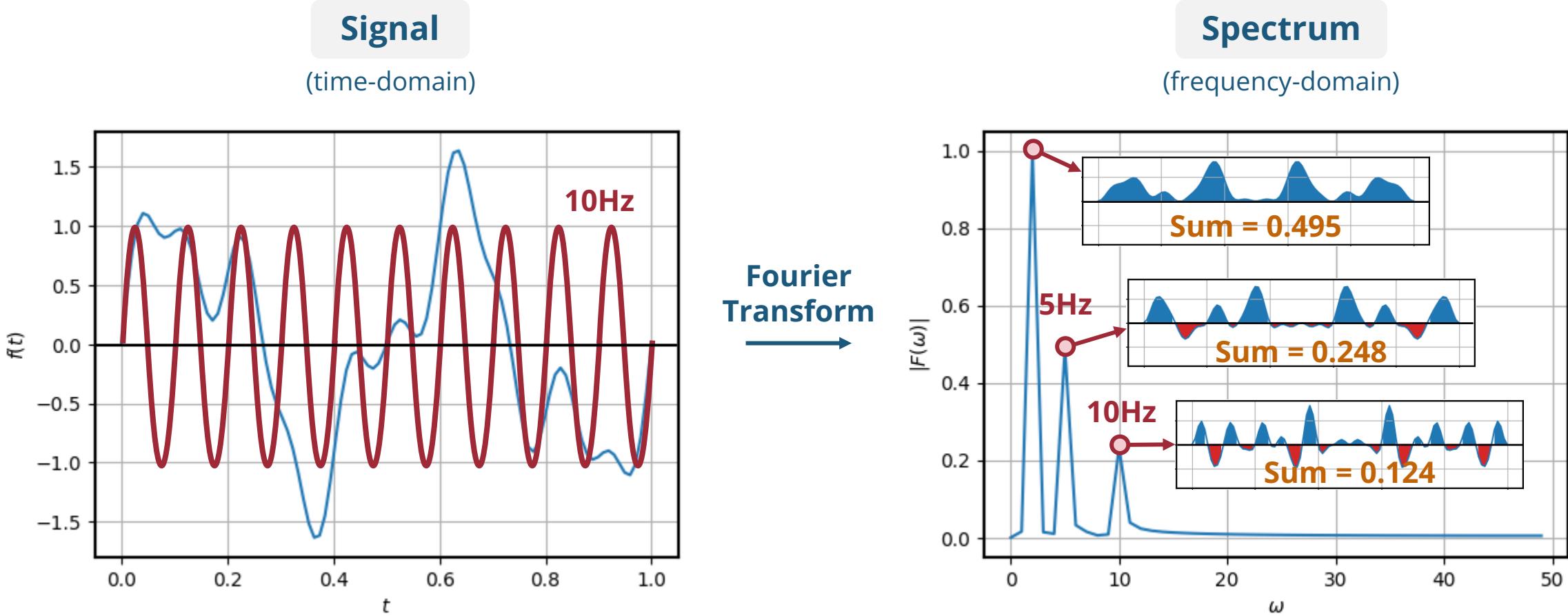


Can only reconstruct frequency up to $\frac{1}{3}f_{max}$

Spectral Analysis

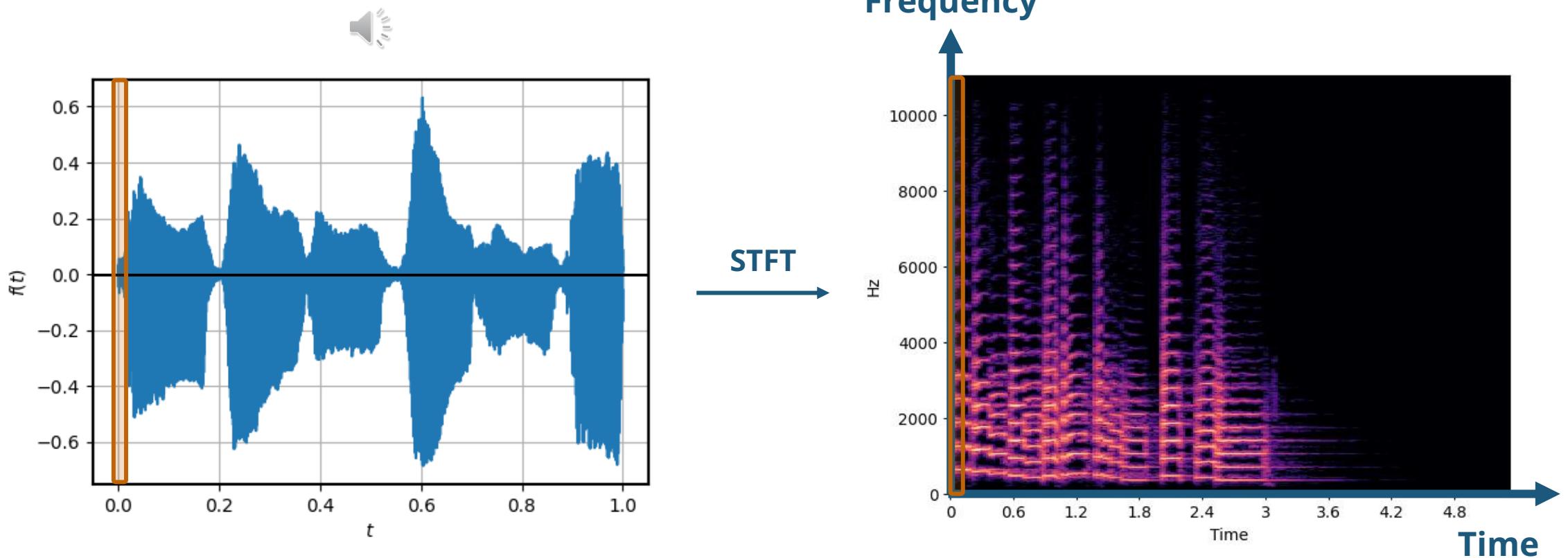


Demystifying Fourier Transform

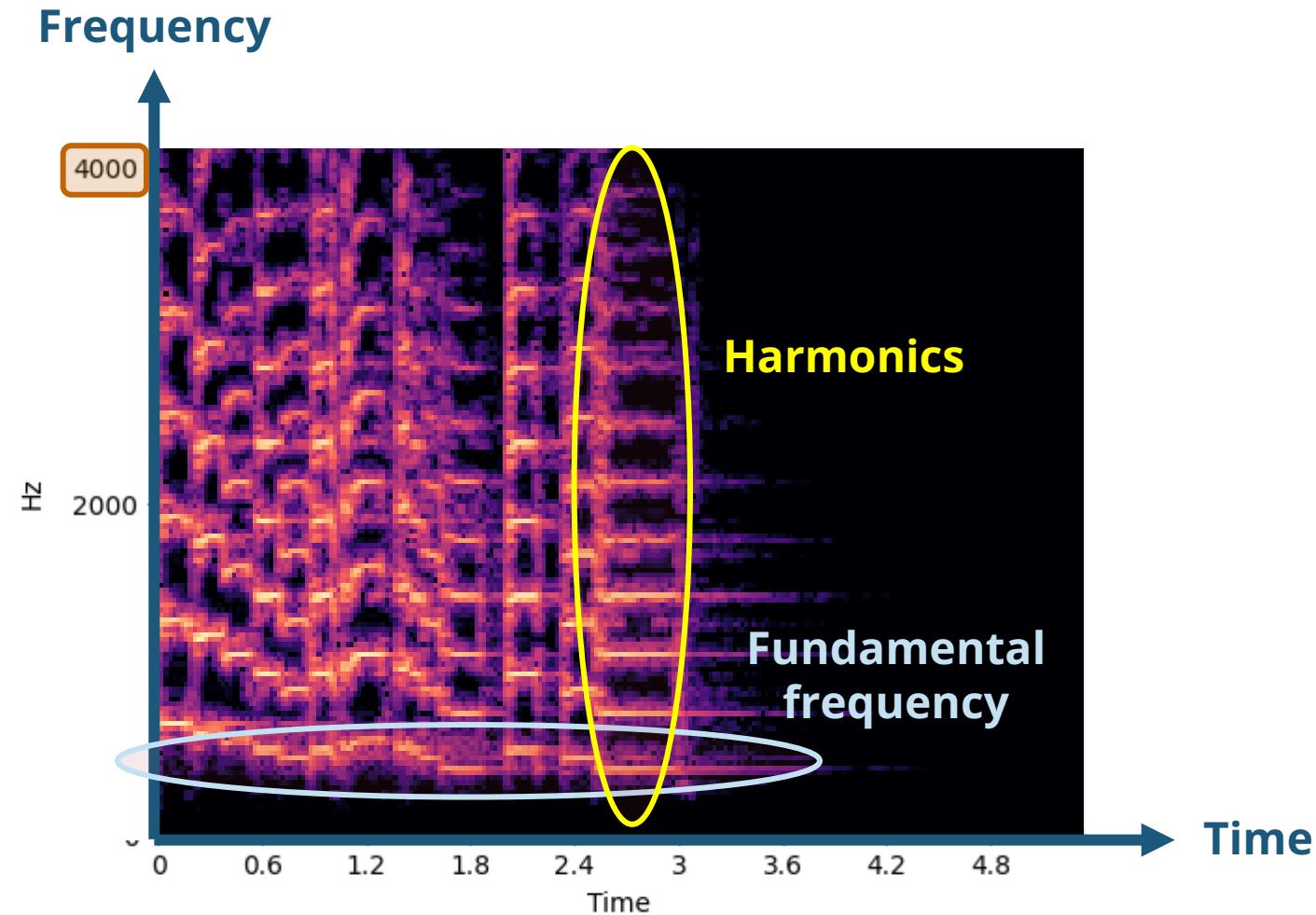


| Short-Time Fourier Transform (STFT)

- **Intuition:** Slice the audio into chunks and apply Fourier transform



Spectrogram



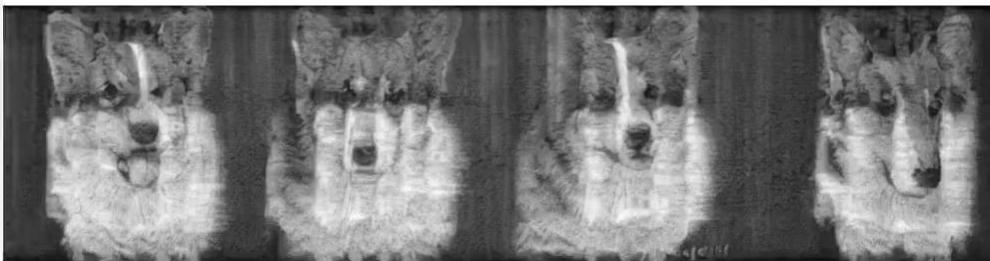
Images that Sound (Chen et al., 2024)

Using diffusion models to generate visual spectrograms that look like images but can also be played as sound.

Image prompt: a colorful photo of corgis



Audio prompt: dog barking



(Source: Chen et al., 2024)

Image prompt: a colorful photo of tigers



Audio prompt: tiger growling



(Source: Chen et al., 2024)

Next Lecture

Deep Learning Fundamentals

