

PAT 464/564 (Winter 2026)

Generative AI for Music & Audio Creation

Lecture 4: Audio Processing Fundamentals

Instructor: Hao-Wen Dong

How do we process audio on a computer?

Four Representative Music Representations



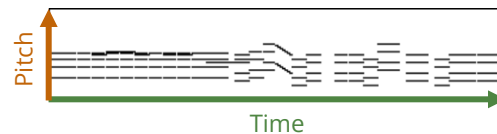
Symbolic music representations

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



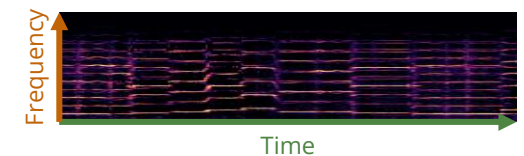
Audio-domain music representations

Time series-based



Waveform

Image-based



Spectrogram

Today's topic!

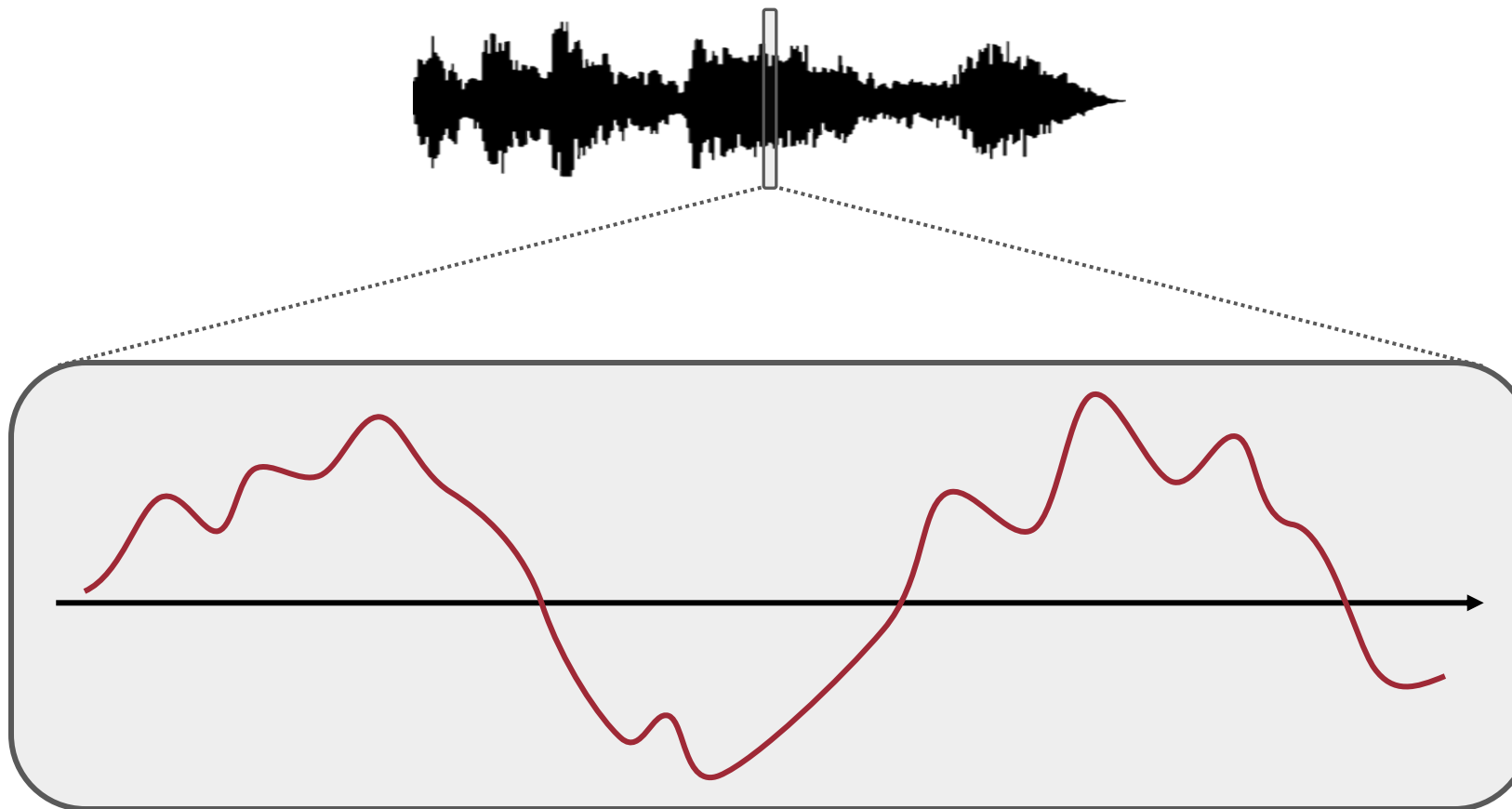
Digital Audio

| Digital Audio

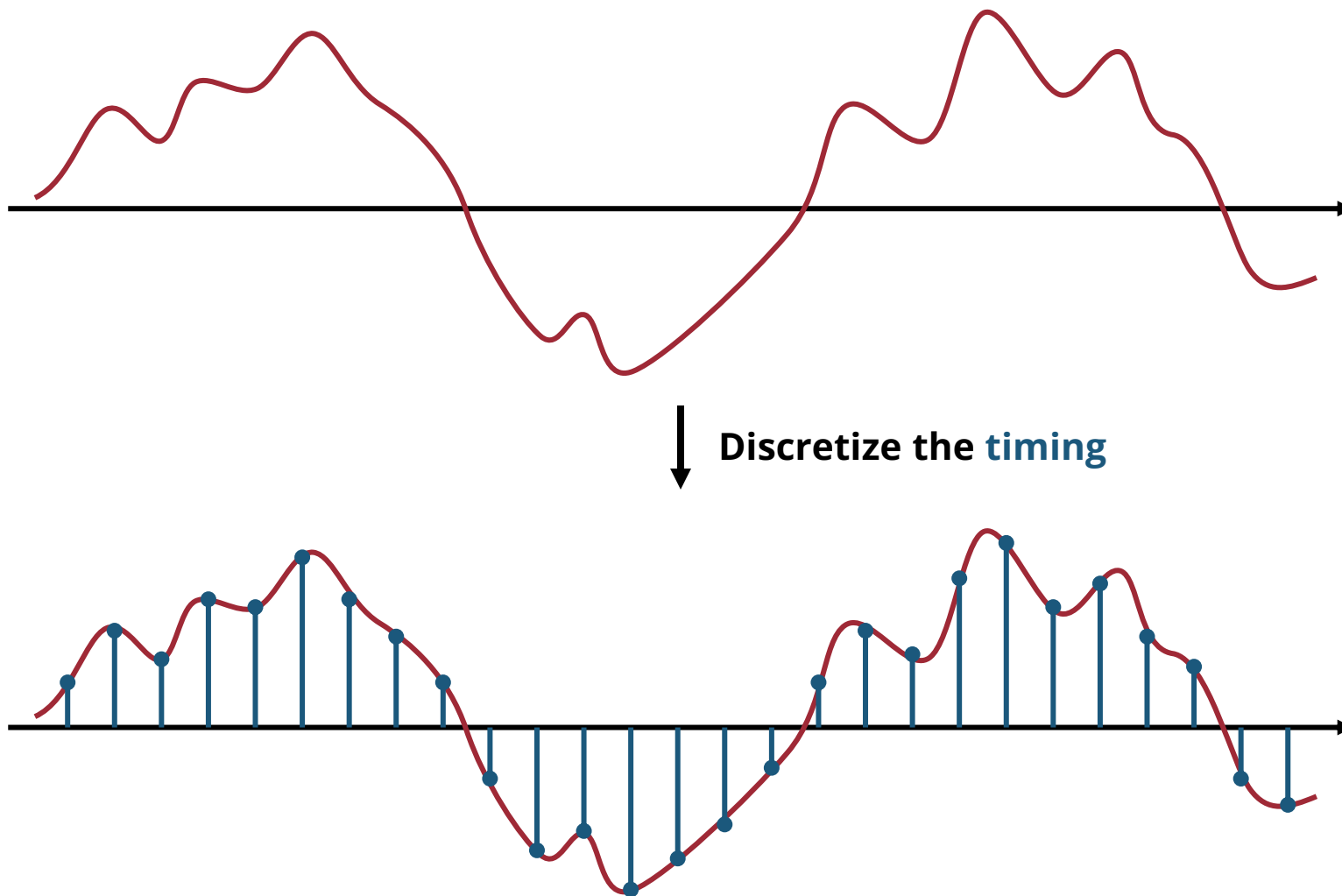


(Source: van den Oord et al., 2016)

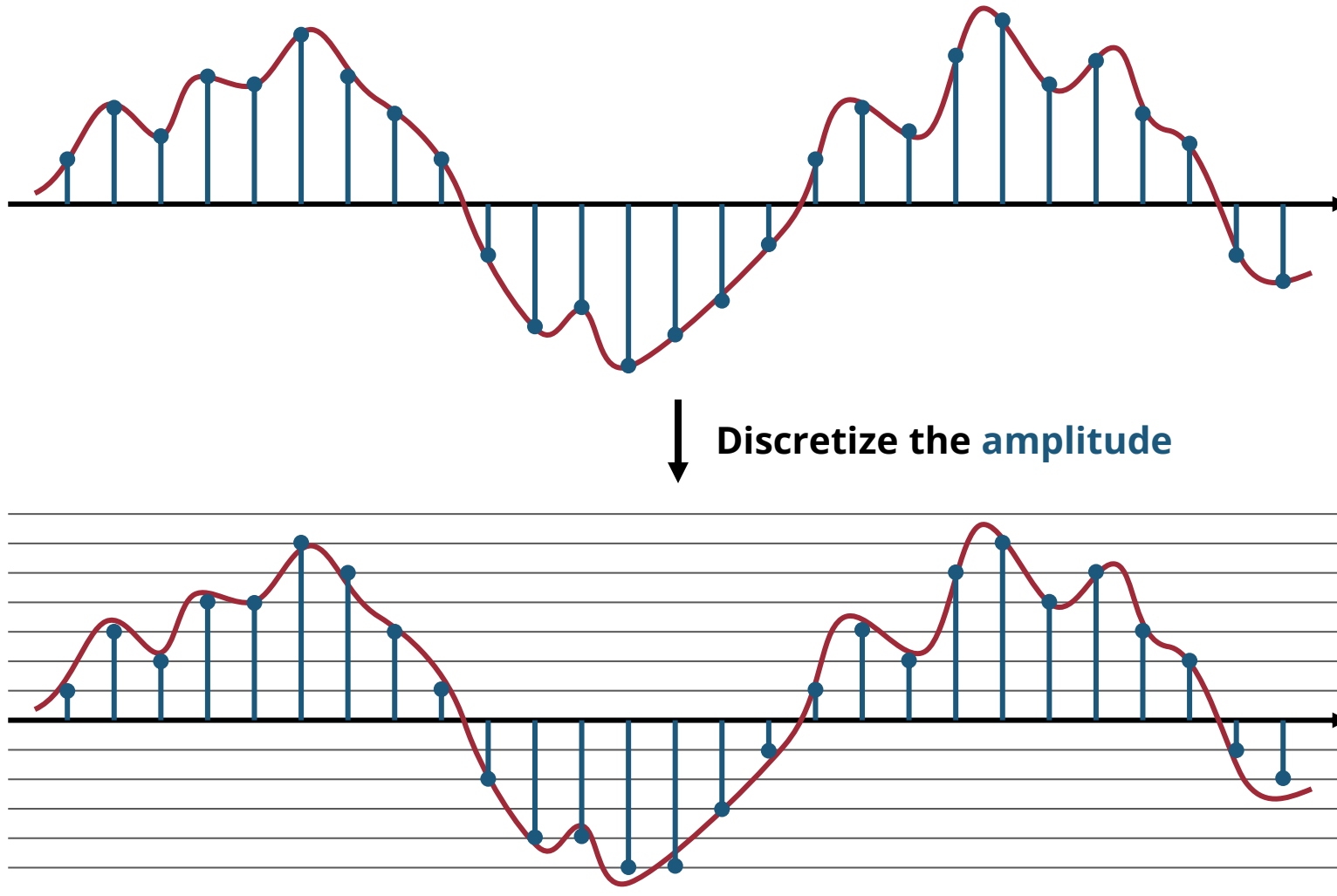
Waveform



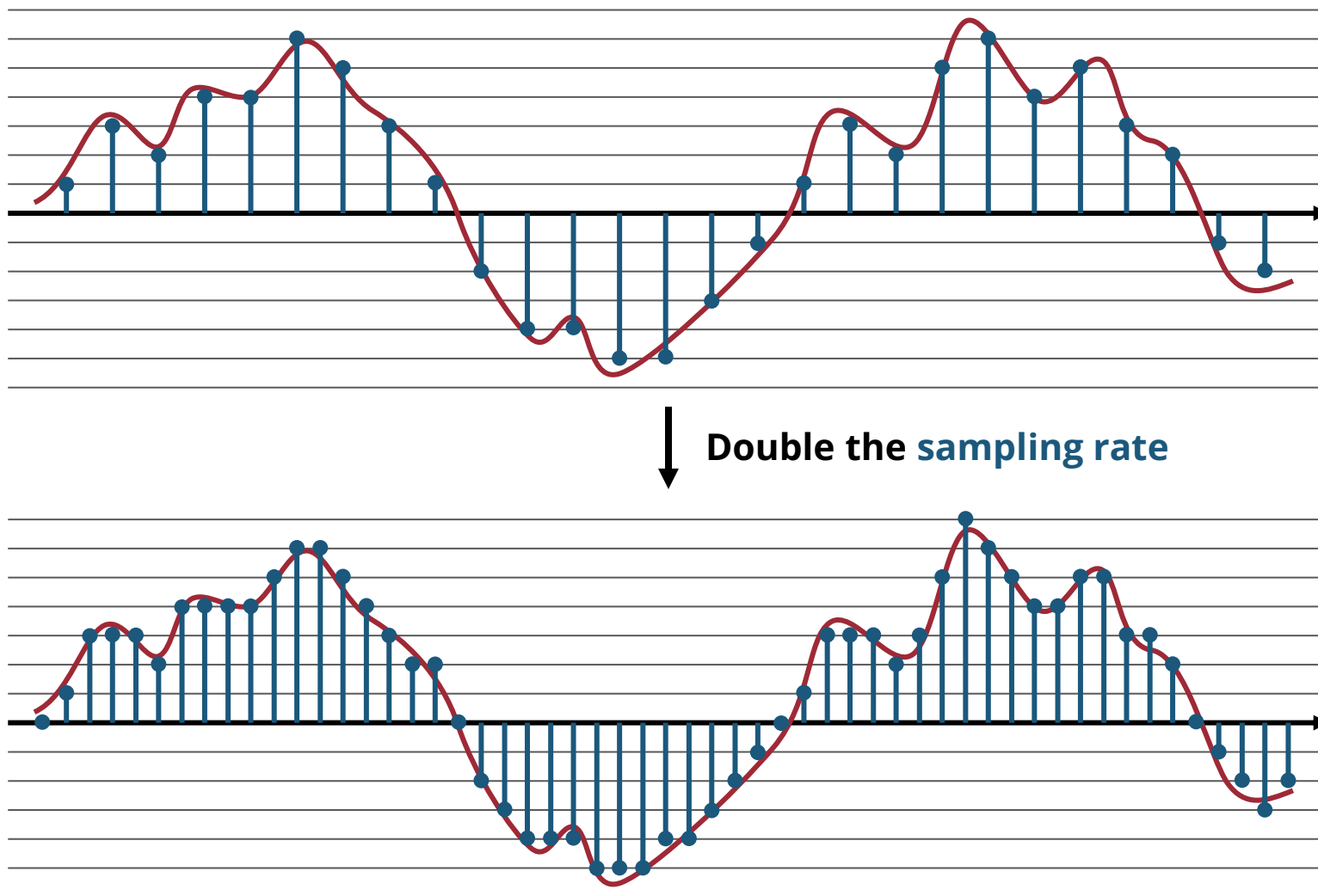
Digitalizing Audio: Timing



Digitalizing Audio: Amplitude



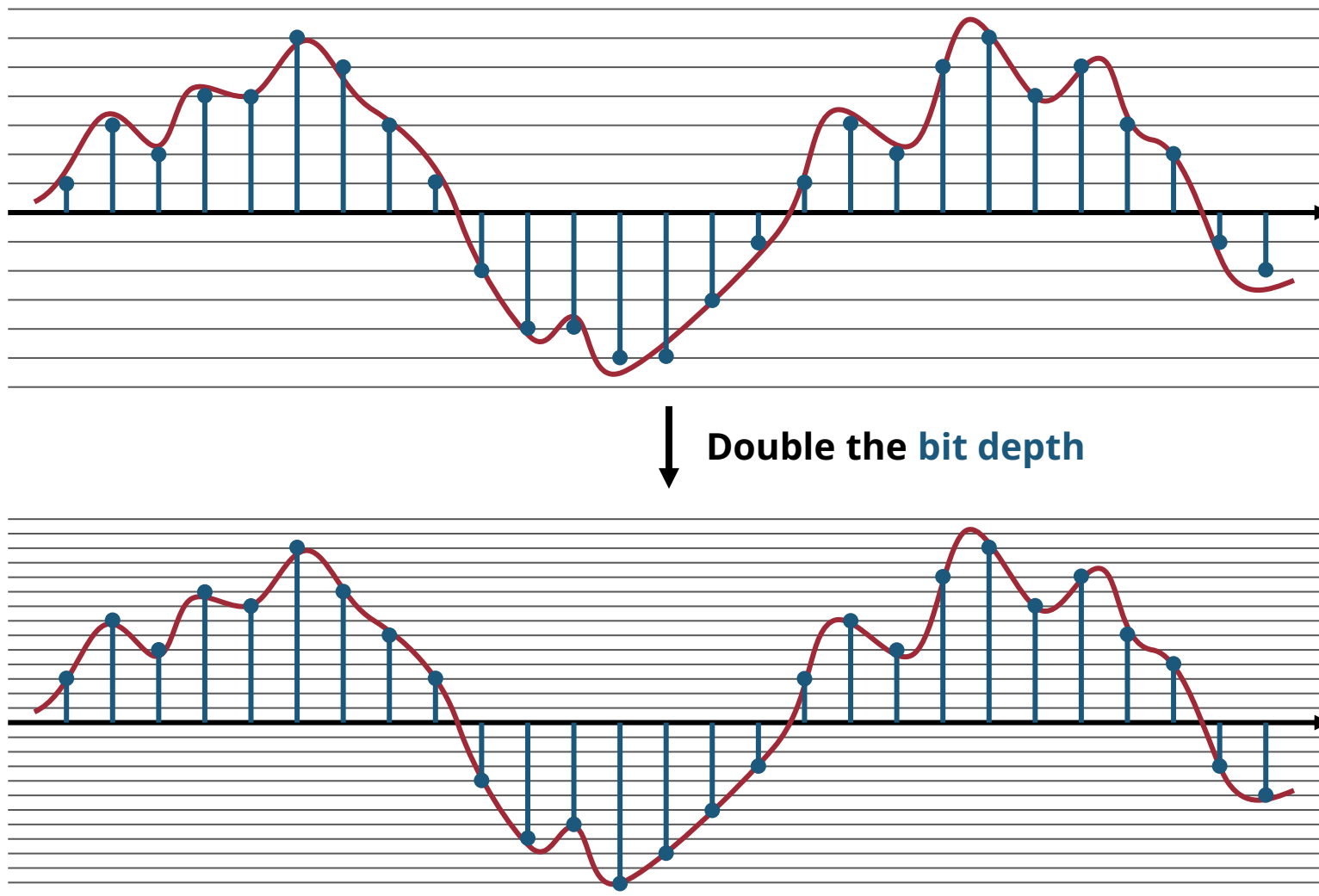
Resolution: Sampling Rate



| 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

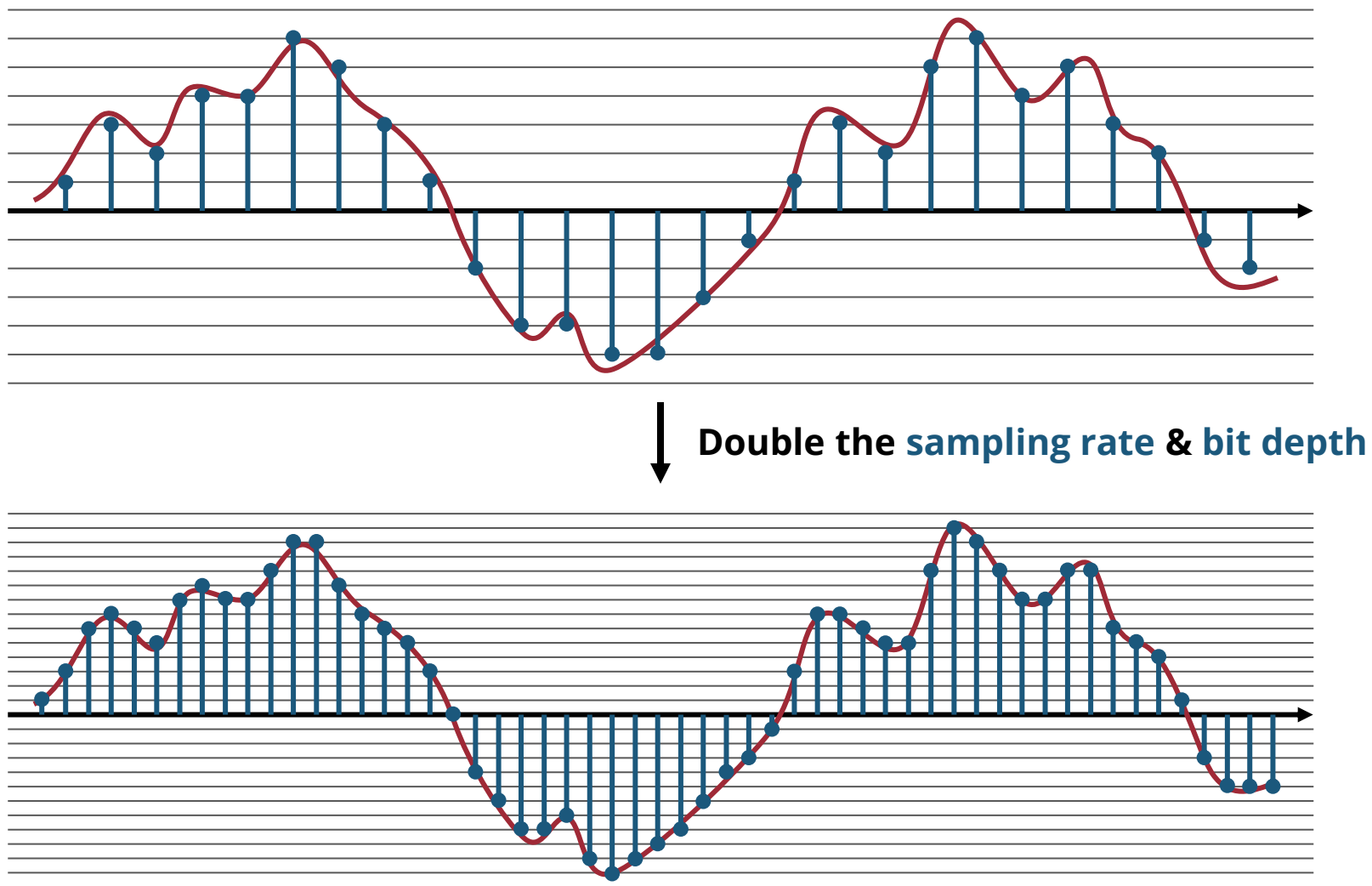
- **Definition: Number of bits used to store each sample**
 - How many bits used to store the amplitude
 - The higher the bit depth, the lower the amplitude distortion
- **Common bit depth**
 - **Chiptunes:** 8 bit
 - **CD:** 16 bit
 - **Modern audio interfaces & DAWs:** 24 bit, 32 bit



| 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 \neq Bit Rate

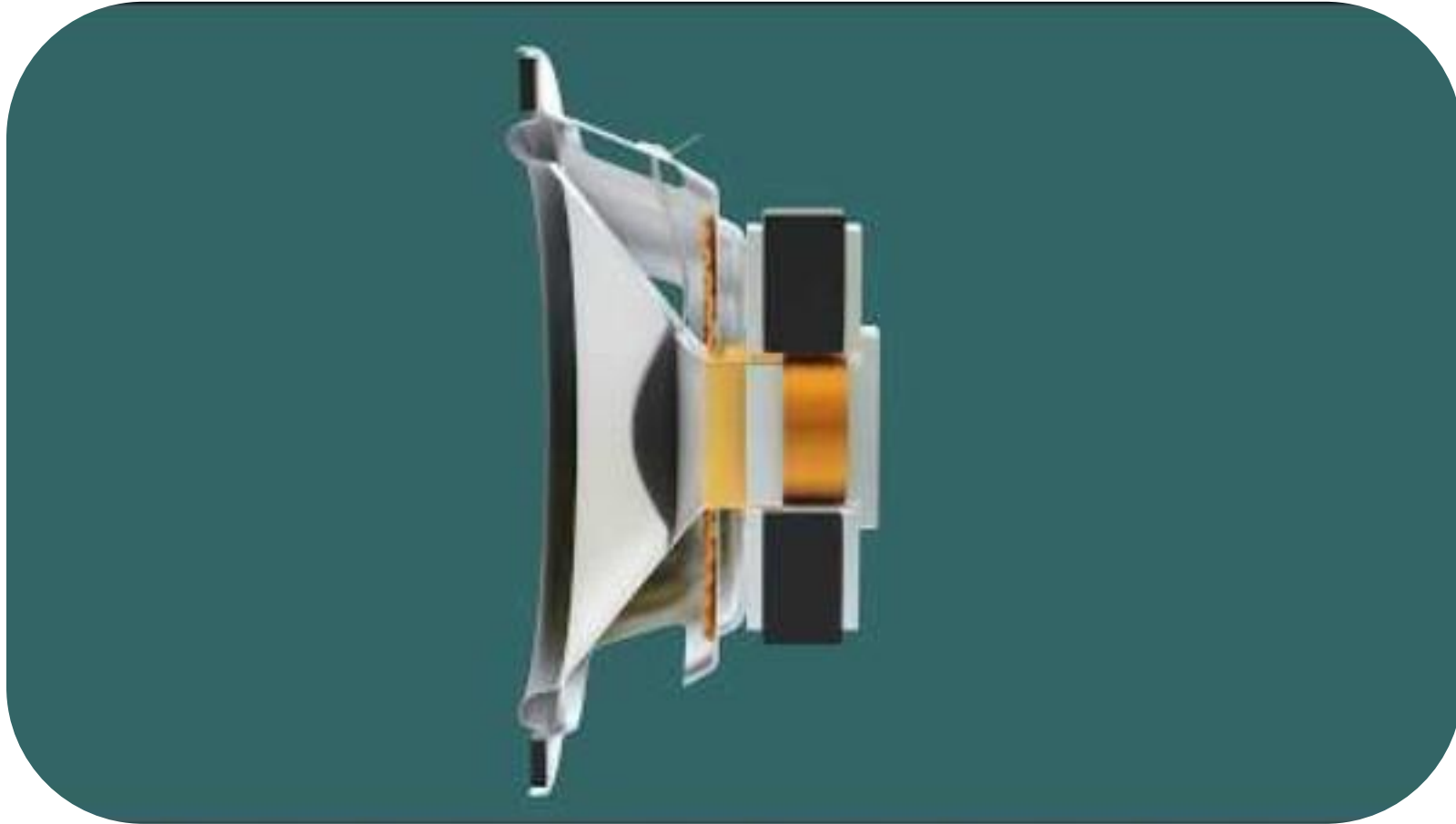
- **Bit Depth:** Number of bits used to store each sample
 - Example: **CD quality** is **16bit/44.1kHz**
- **Bit Rate:** Amount of data transferred per second (unit: bits/sec)
 - Example: **320K MP3** files \rightarrow **320kbps** (320,000 bits per second)
 - Example: **YouTube** recommendation \rightarrow **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

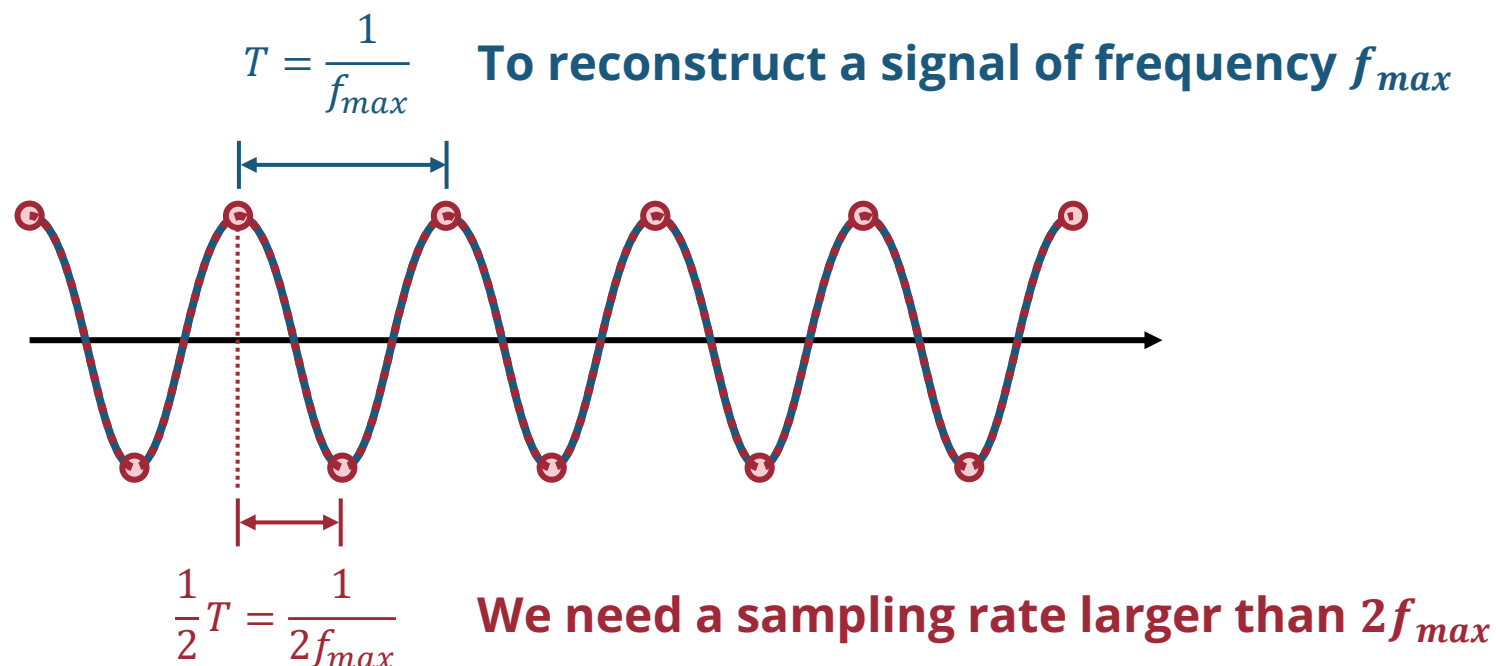


youtu.be/RxdFP31QYAg

Sampling Theorem

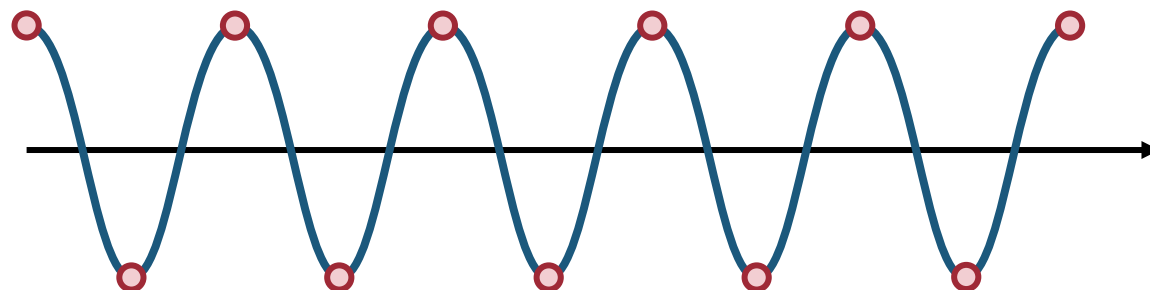
Nyquist–Shannon Sampling Theorem

- **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**

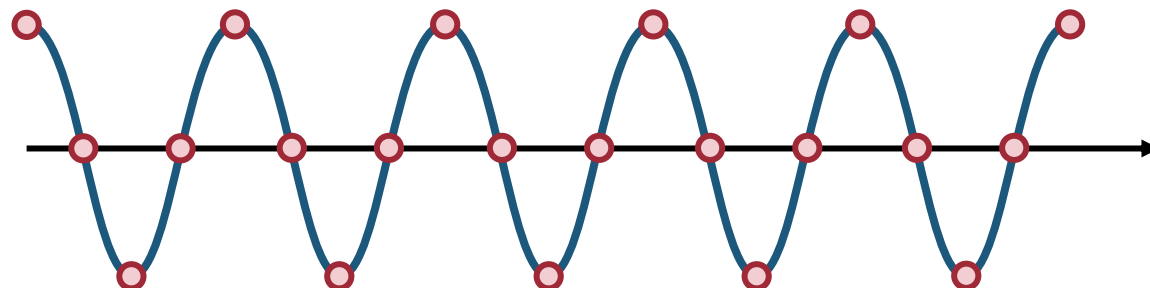


Sampling Theorem: Oversampling

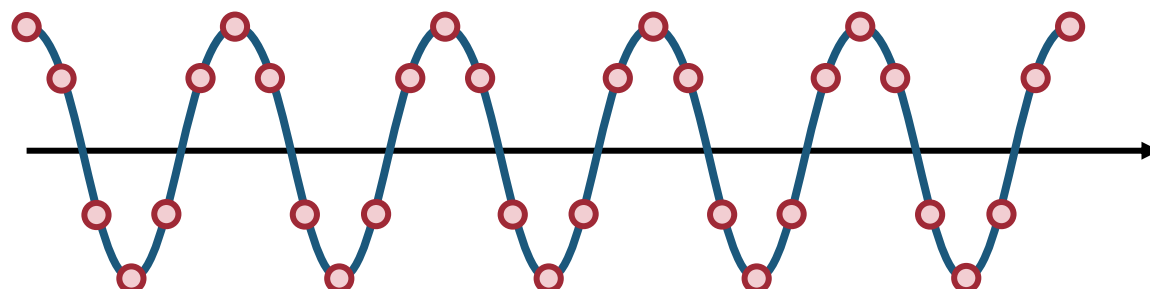
Critically sampled
($f_s = 2f_{max}$)



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



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

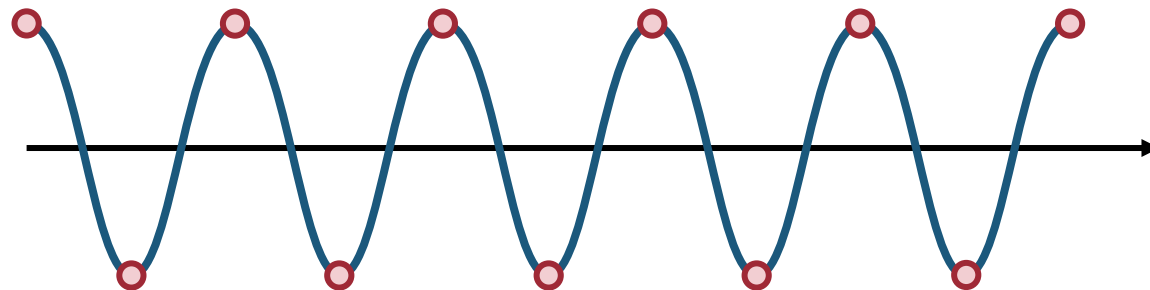


Reconstruction
is possible!

Sampling Theorem: Undersampling

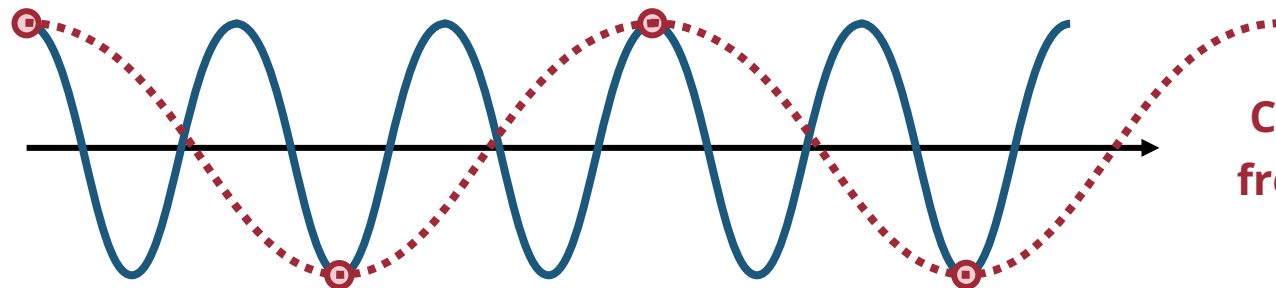
Critically sampled

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



Undersampled

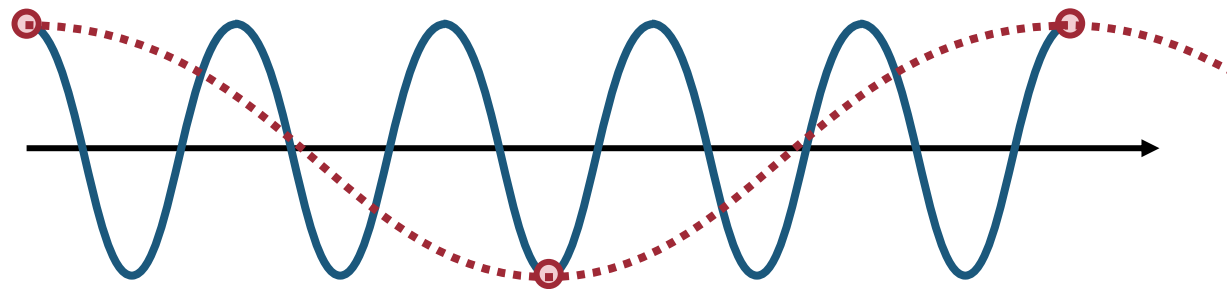
$$(f_s = \frac{2}{3}f_{max})$$



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

Undersampled

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

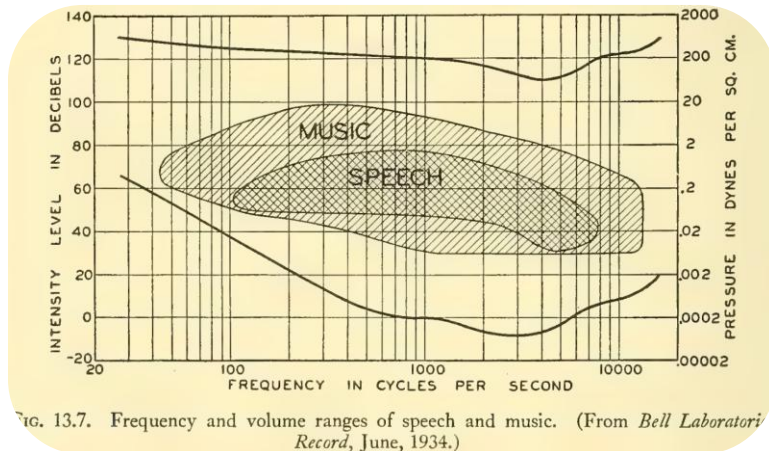


Can only reconstruct
frequency up to $\frac{1}{5}f_{max}$

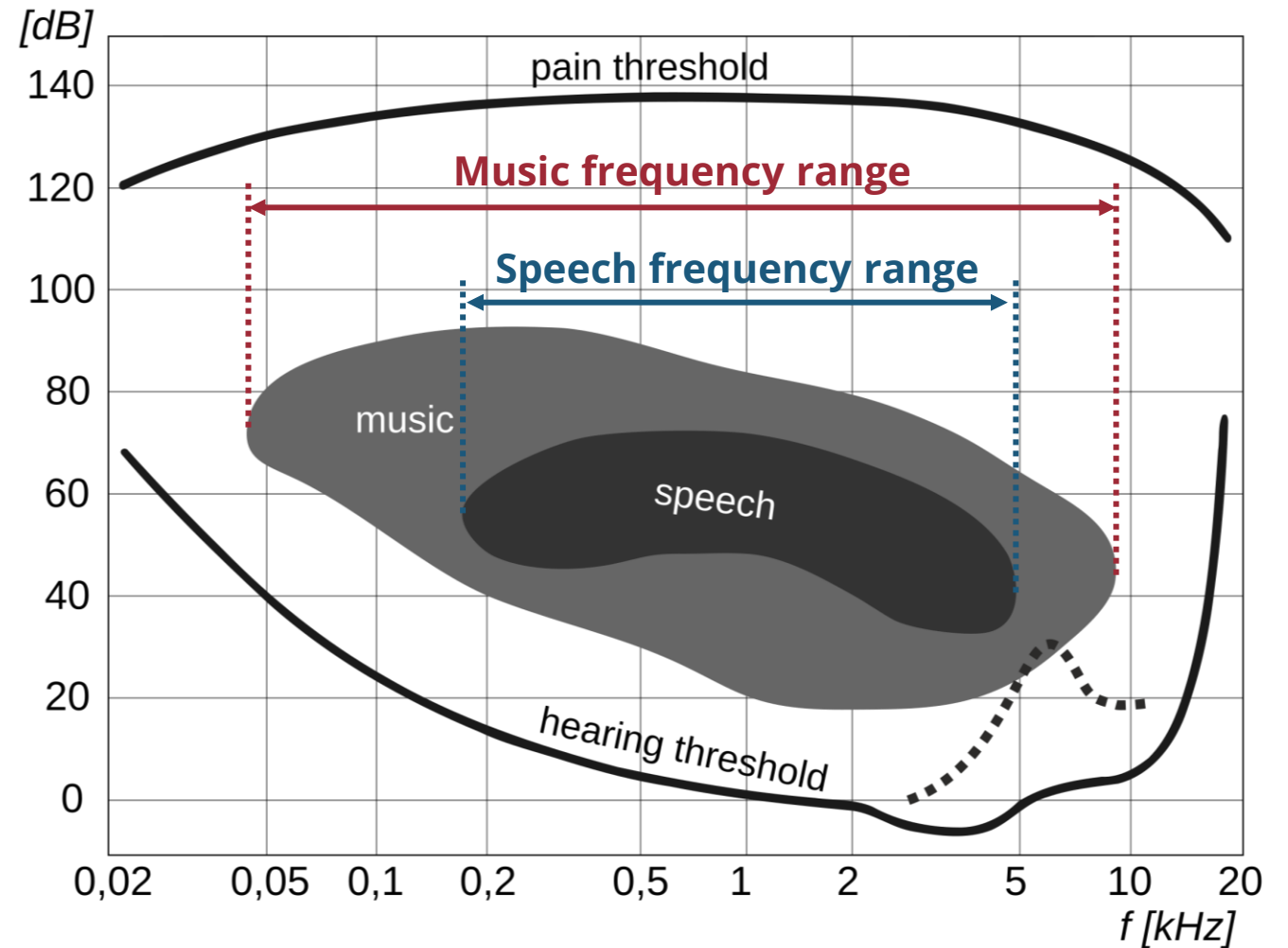
🤔 Sampling Theorem

- **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)



(Source: Wikipedia)

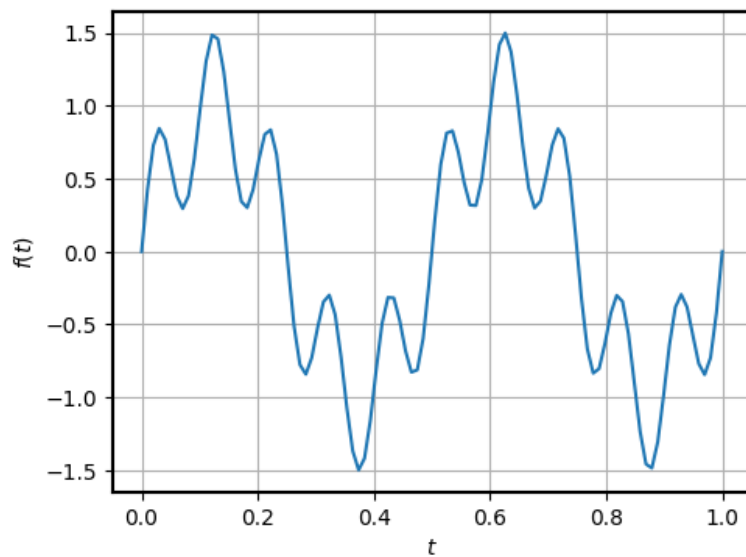
Bell Laboratories Record, 12(6):314, 1934.

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

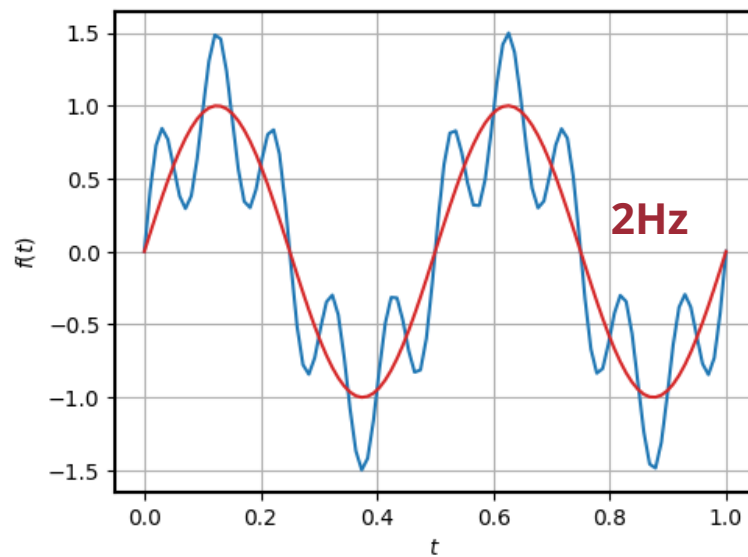
Spectral Analysis

Spectral Analysis

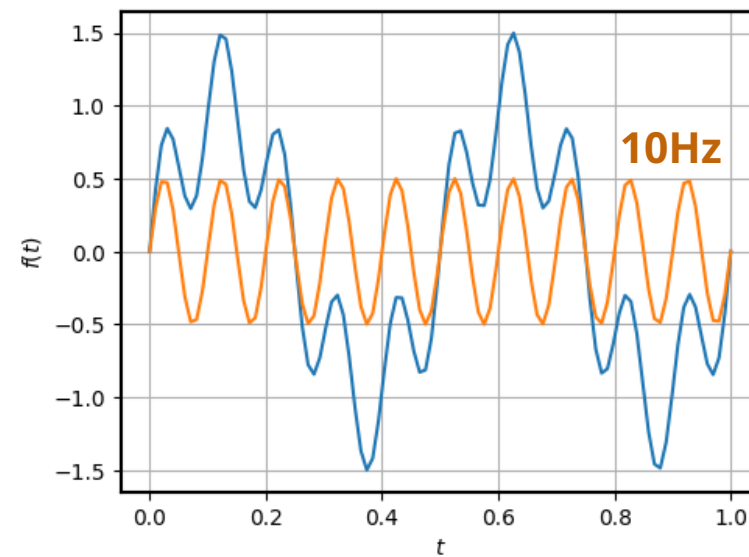
- **Goal:** Analyze the **frequency components** of a signal



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

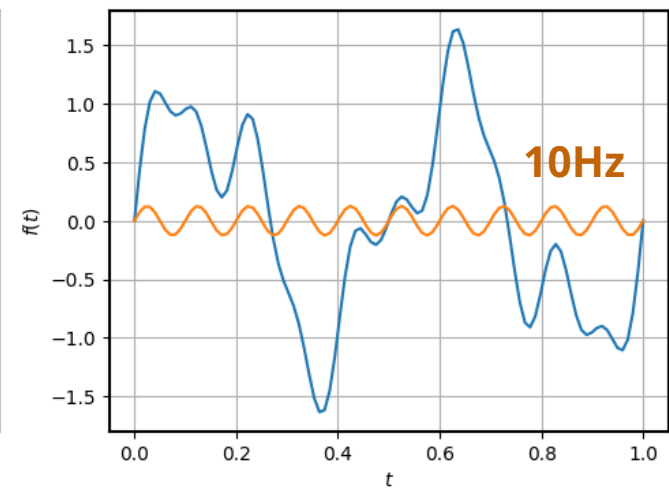
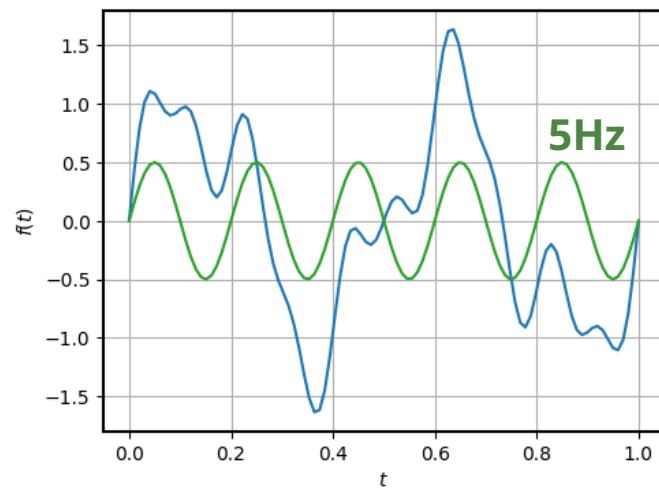
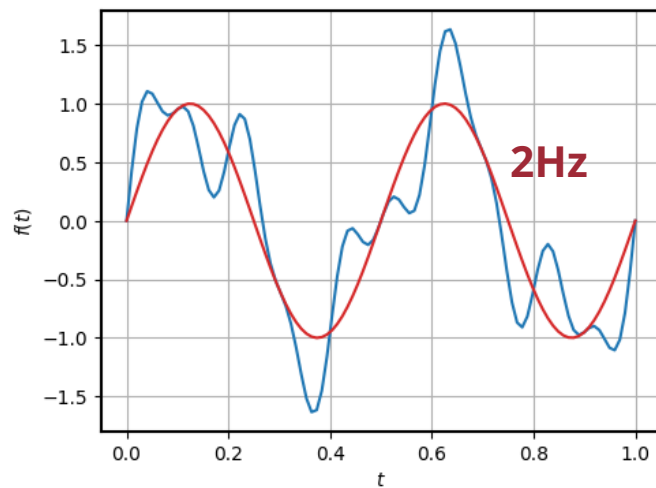
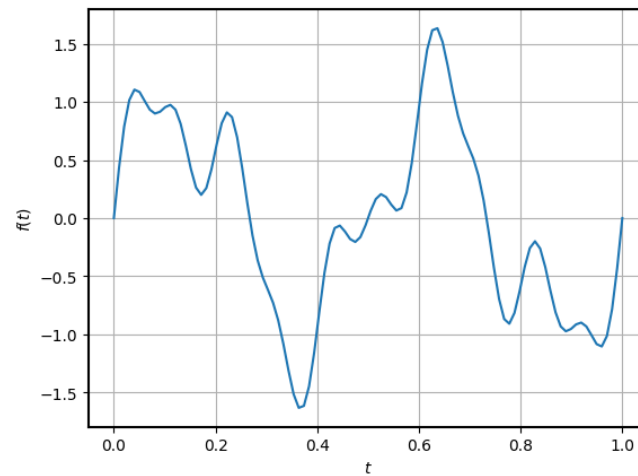


$$\sin(2 \cdot 2\pi t)$$



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

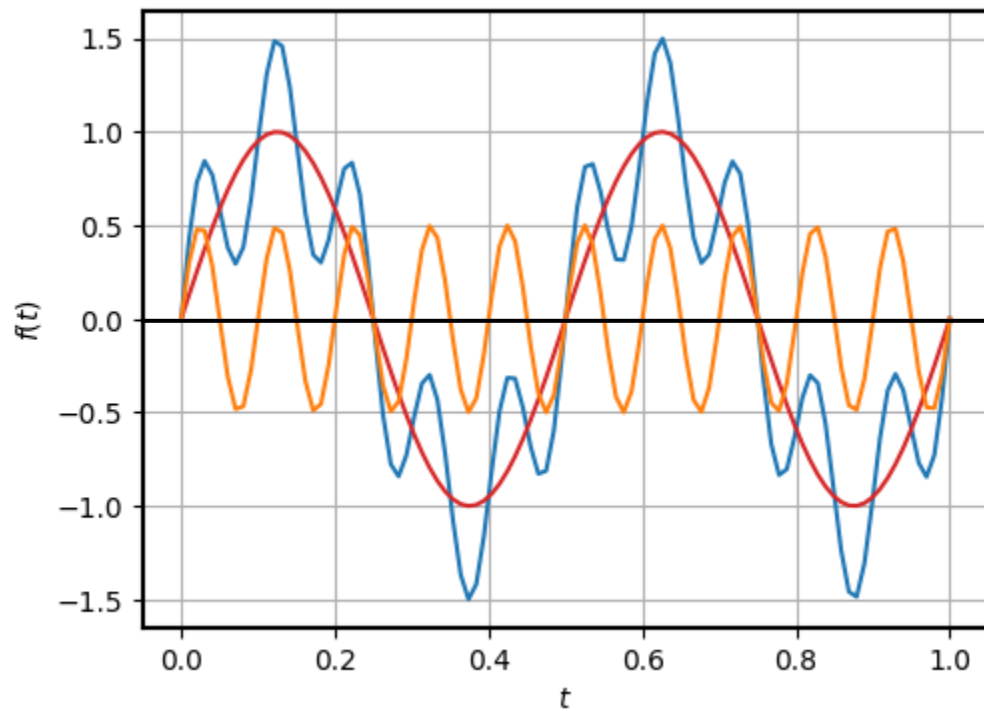
Spectral Analysis



Fourier Transform

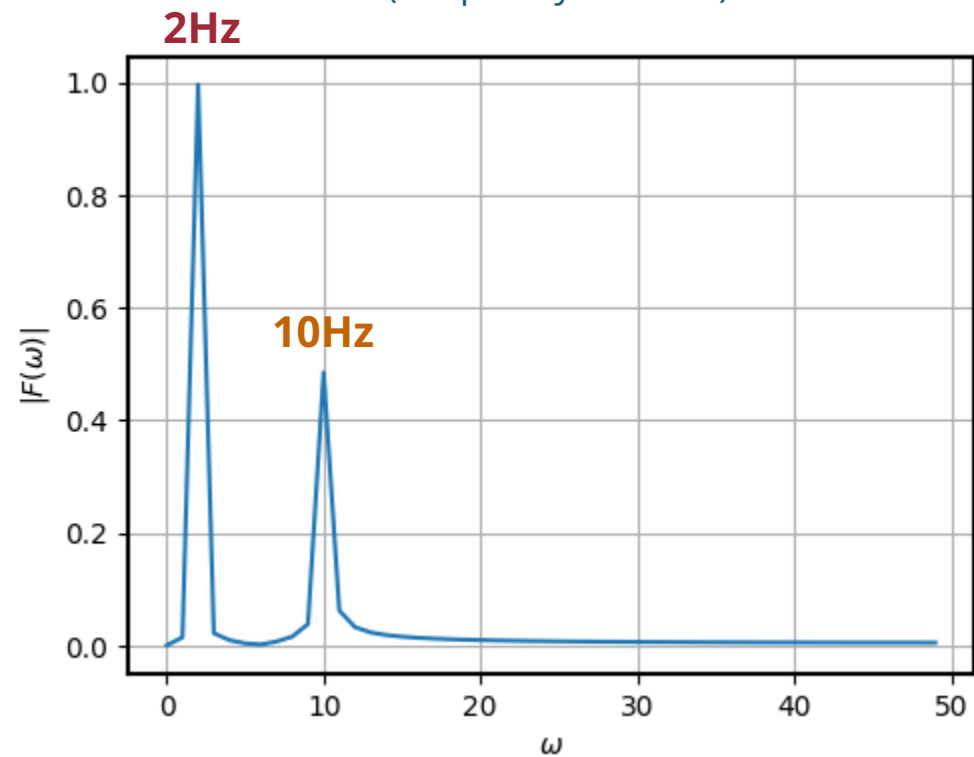
Signal

(time-domain)



Spectrum

(frequency-domain)



Fourier Transform


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

The diagram illustrates the Fourier Transform equation with various annotations:

- Output spectrum:** A green arrow points to the $F(\omega)$ term, which is enclosed in a green box.
- Frequency:** An orange arrow points to the ω symbol inside the $F(\omega)$ box.
- Input signal:** A blue arrow points to the $f(t)$ term, which is enclosed in a blue box.
- Sum over all t :** A purple arrow points to the integral symbol $\int_{-\infty}^{\infty}$, which is enclosed in a purple box.
- Complex exponential:** The term $e^{-j\omega t}$ is enclosed in a red box, with a thinking face emoji above it.
- Differential element:** The term dt is enclosed in a purple box, with a purple arrow pointing to it.

The equation is presented as:
$$F(\omega) = \int_{-\infty}^{\infty} f(t) e^{-j\omega t} dt$$

Demystifying Fourier Transform

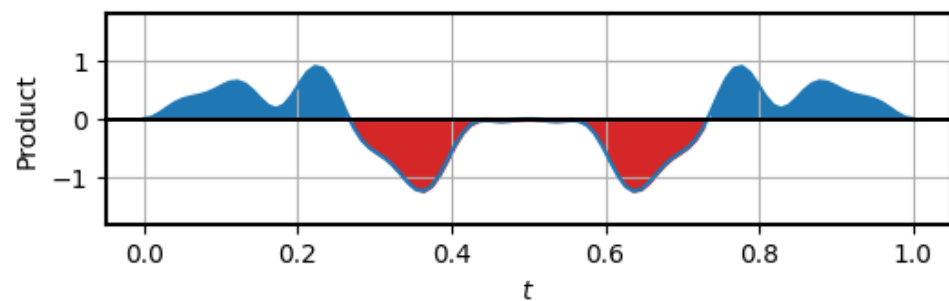
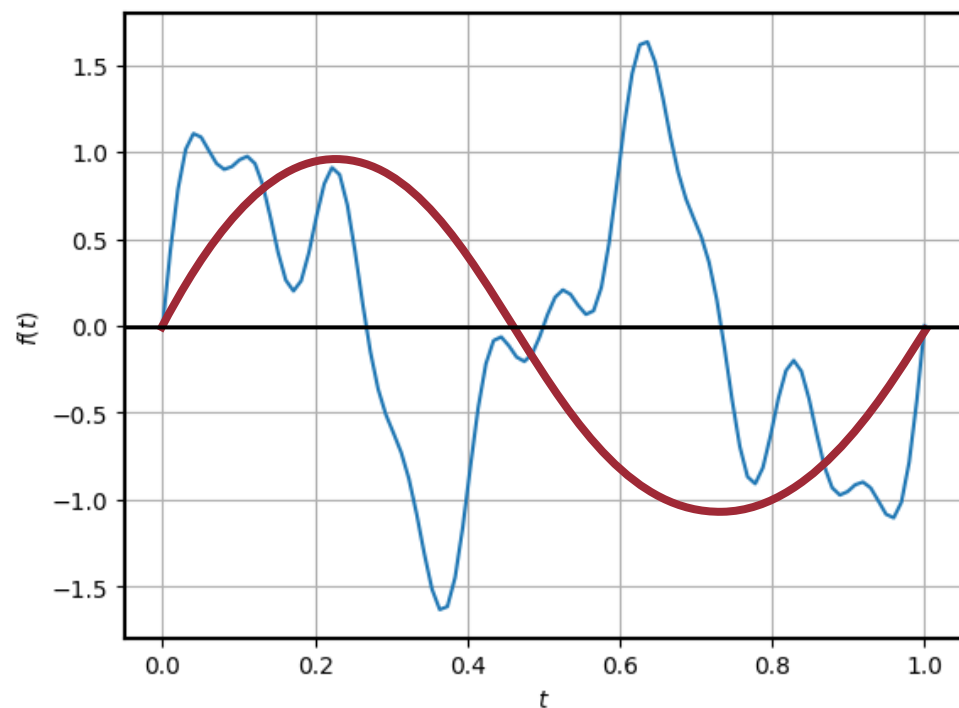
$$F(\omega) = \int_{-\infty}^{\infty} f(t) \boxed{e^{-j\omega t}} dt$$


Euler's formula

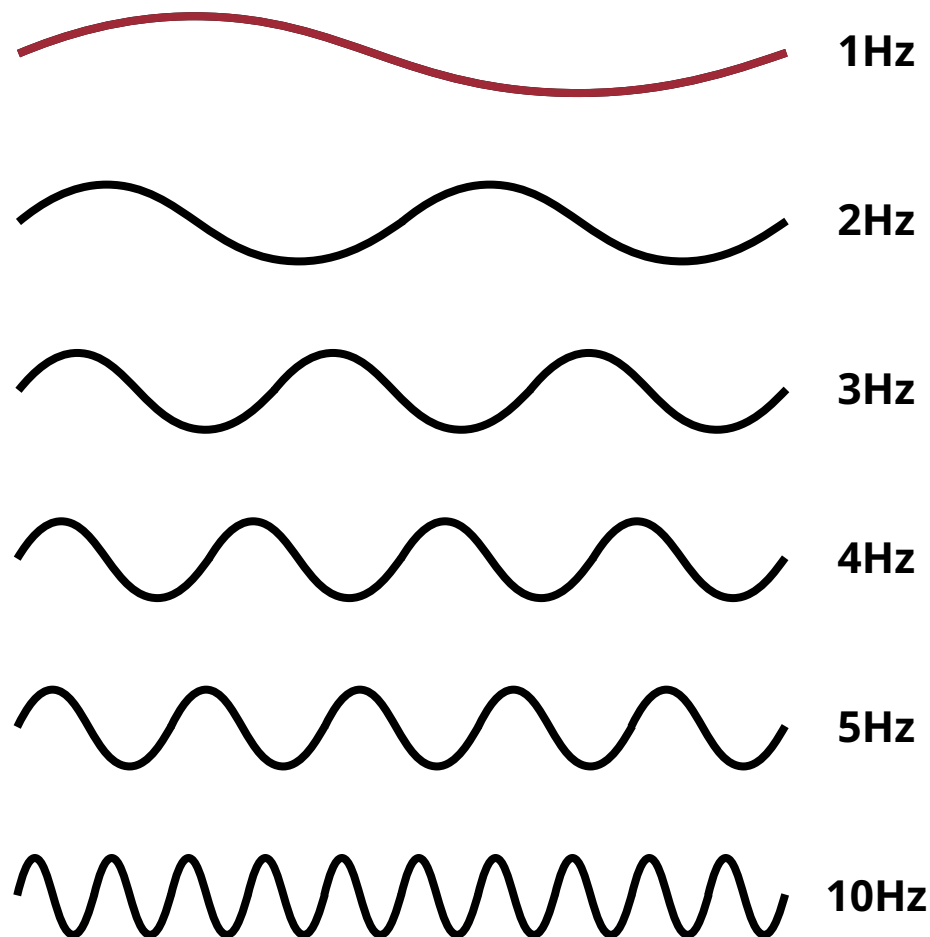
$$e^{-j\theta} = \cos \theta + j \sin \theta$$

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

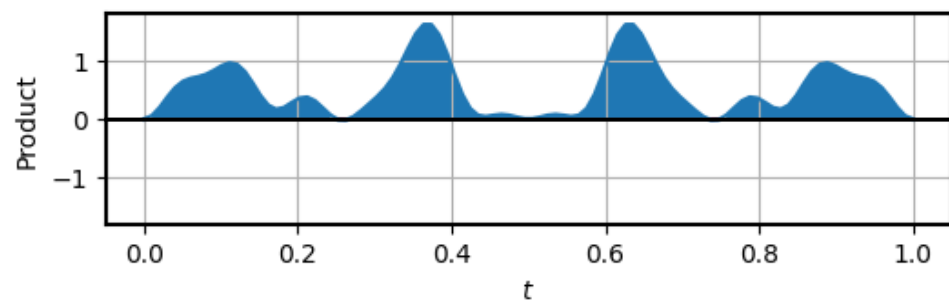
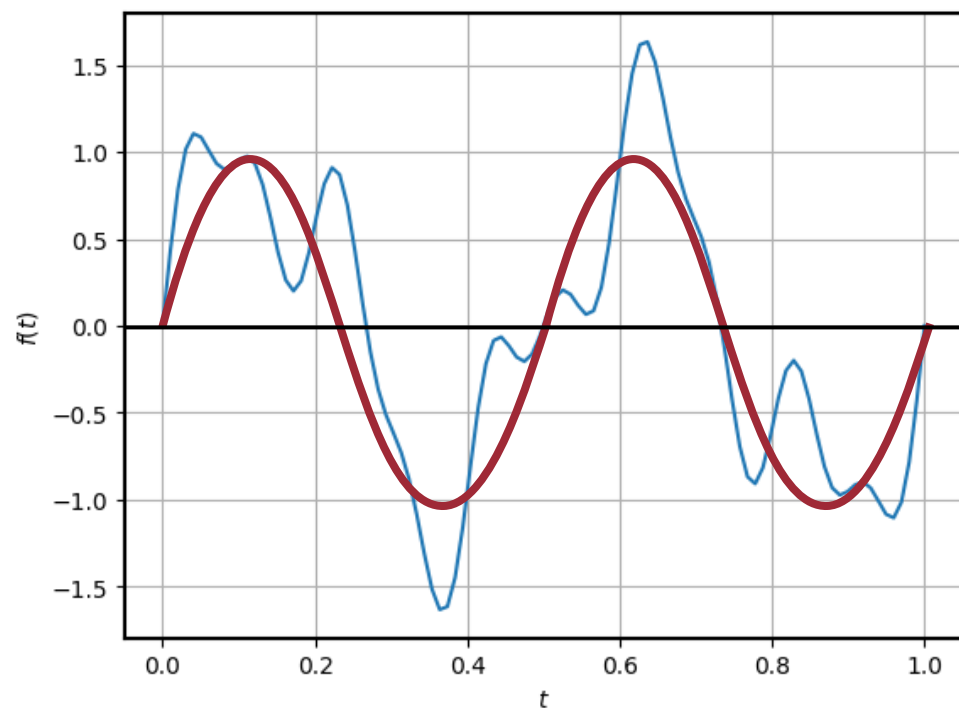
Demystifying Fourier Transform



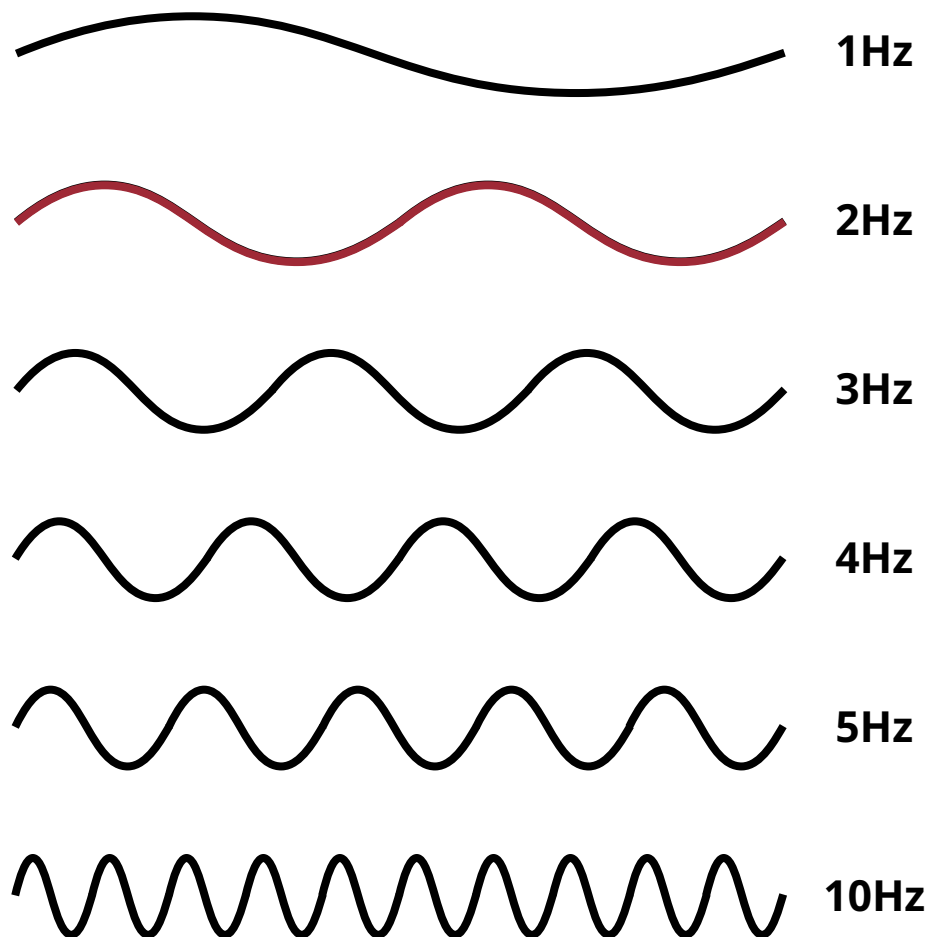
Candidate frequency components



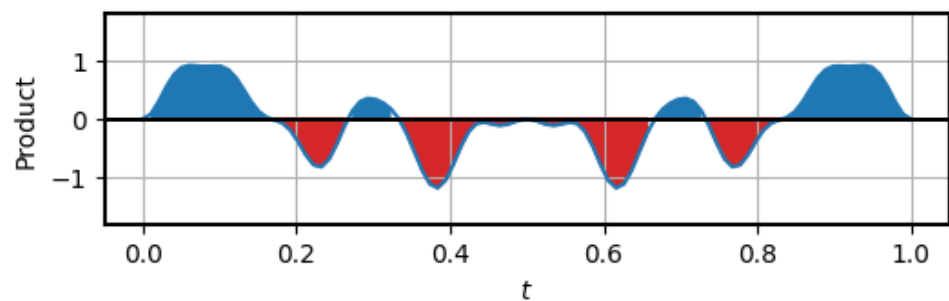
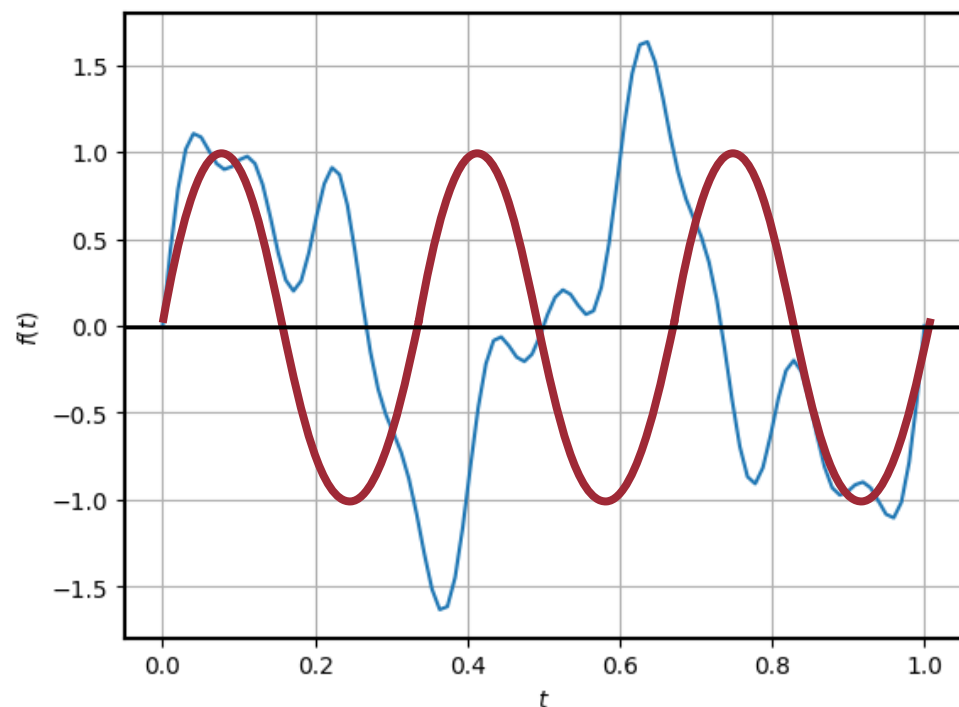
Demystifying Fourier Transform



Candidate frequency components



Demystifying Fourier Transform



Candidate frequency components



1Hz



2Hz



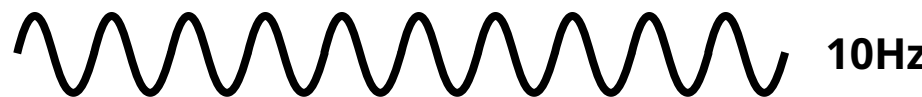
3Hz



4Hz

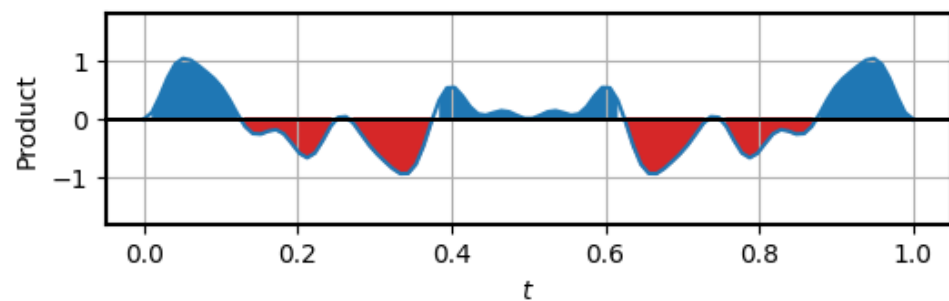
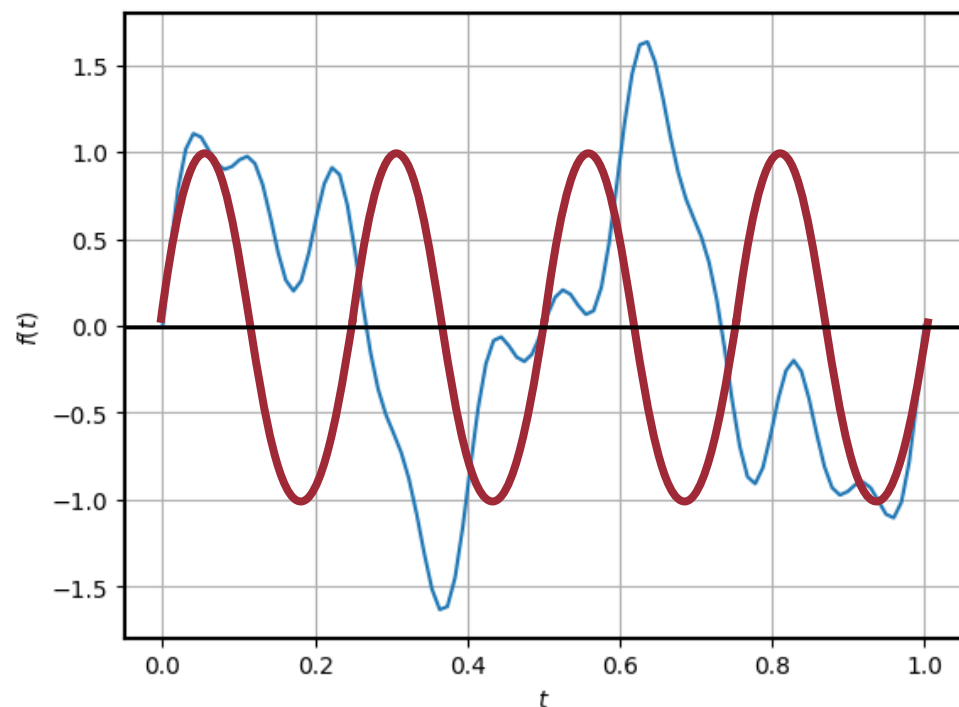


5Hz

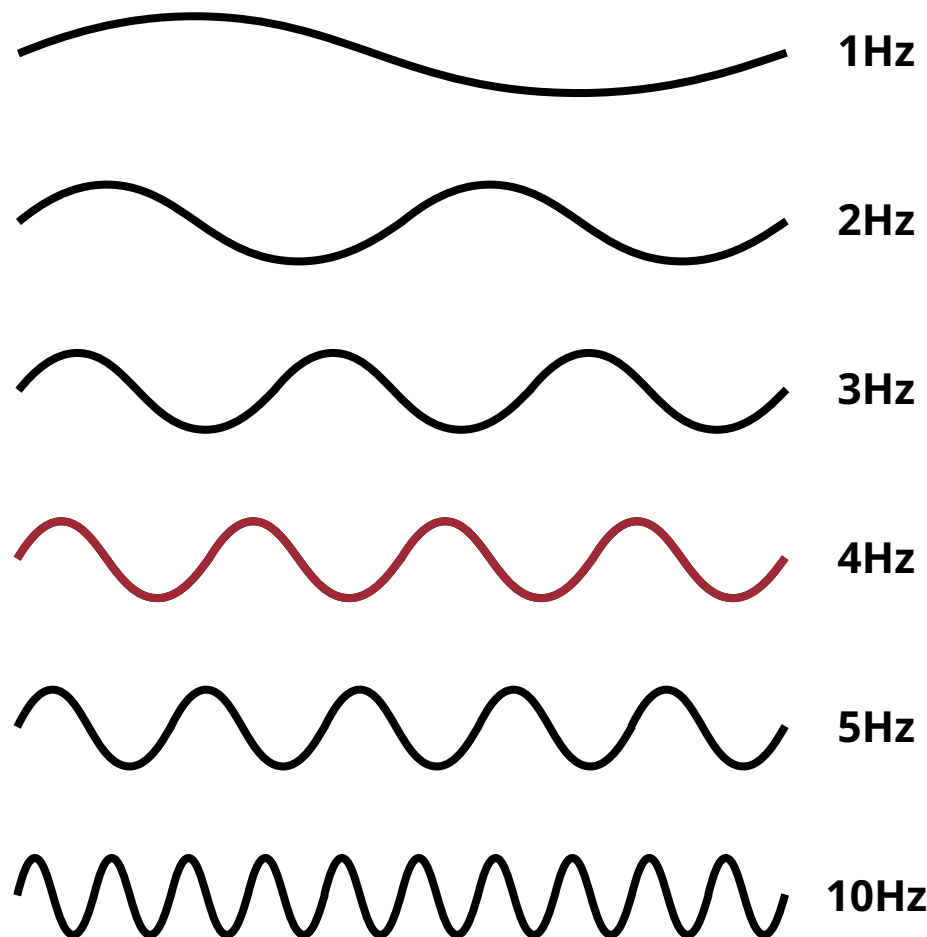


10Hz

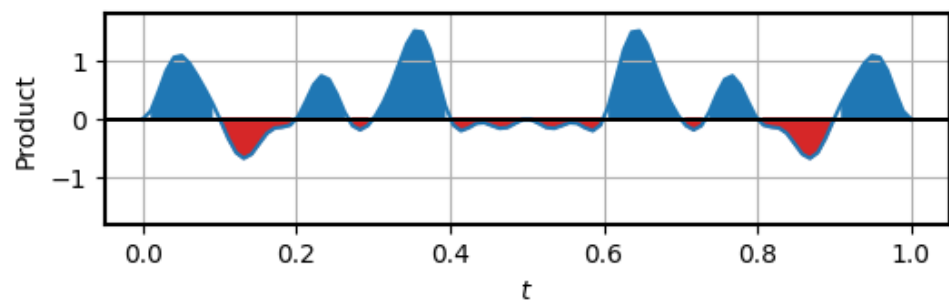
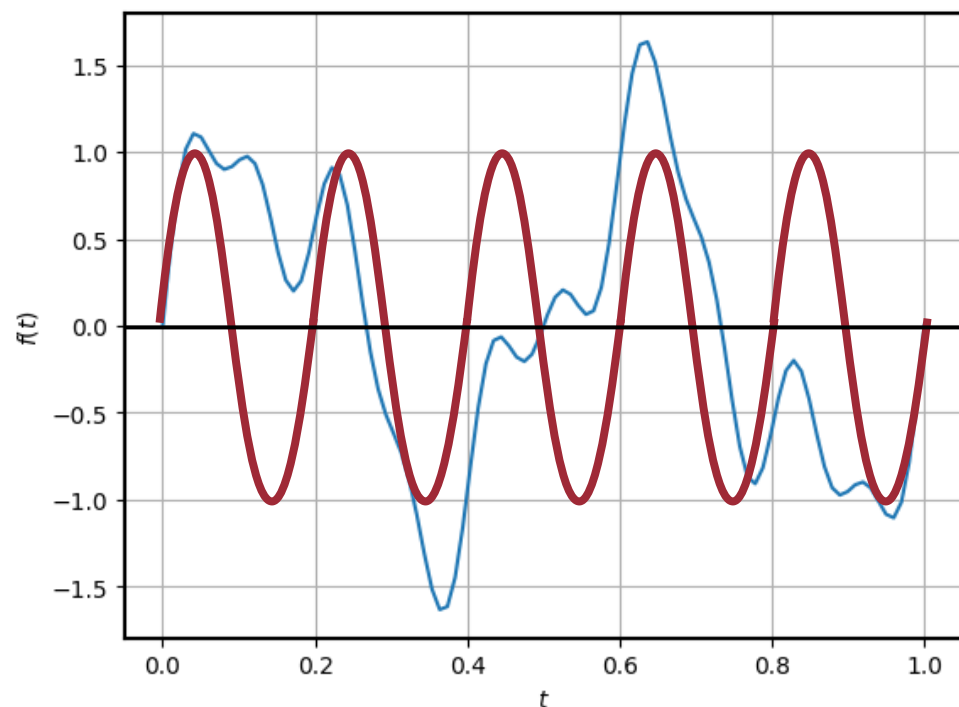
Demystifying Fourier Transform



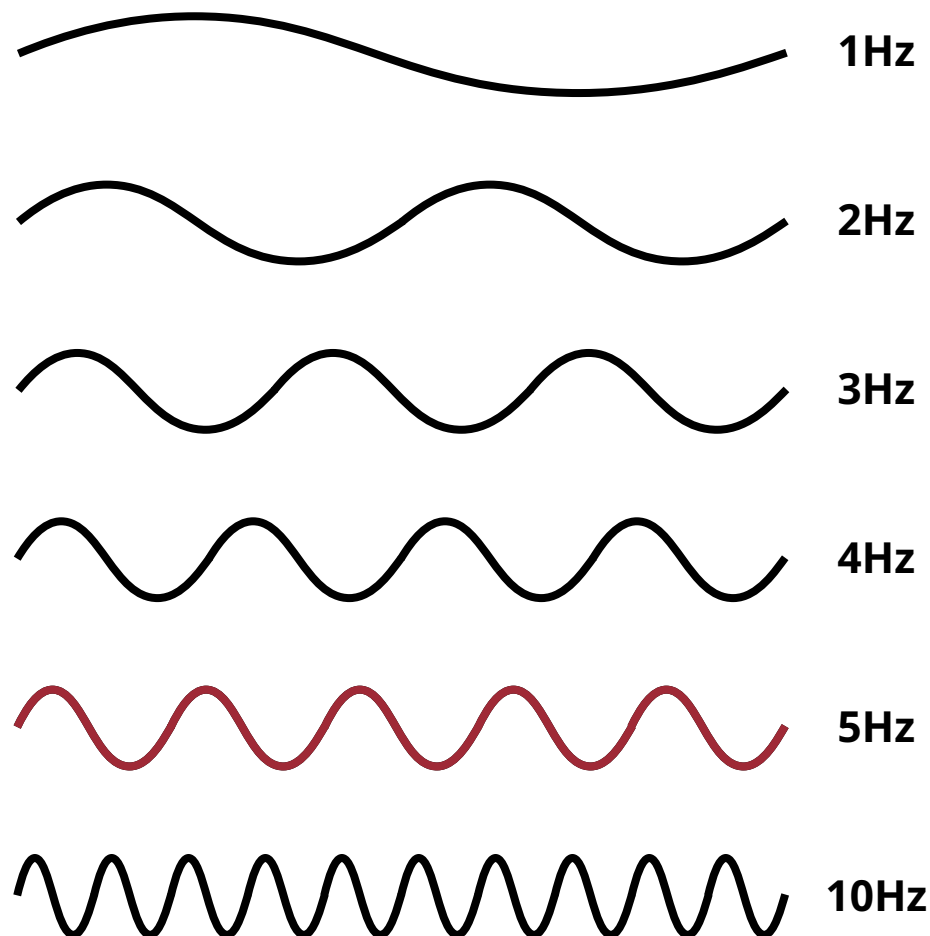
Candidate frequency components



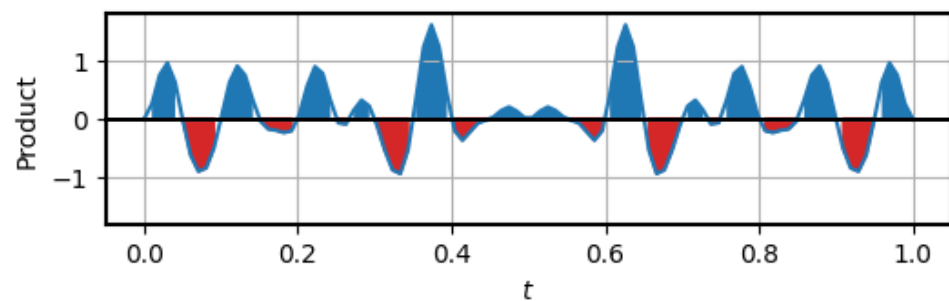
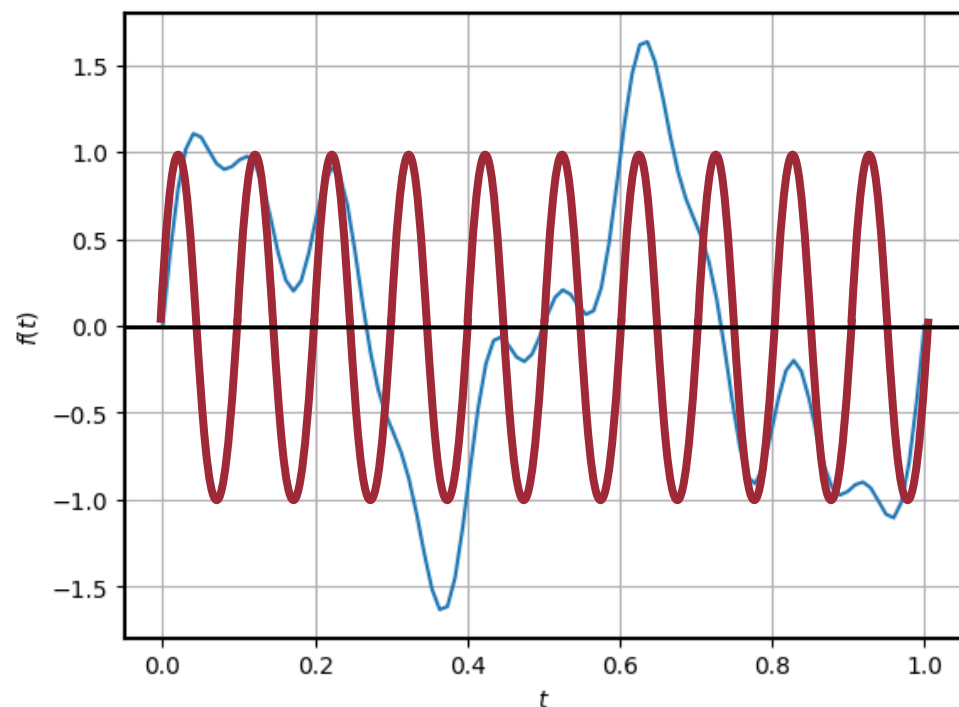
Demystifying Fourier Transform



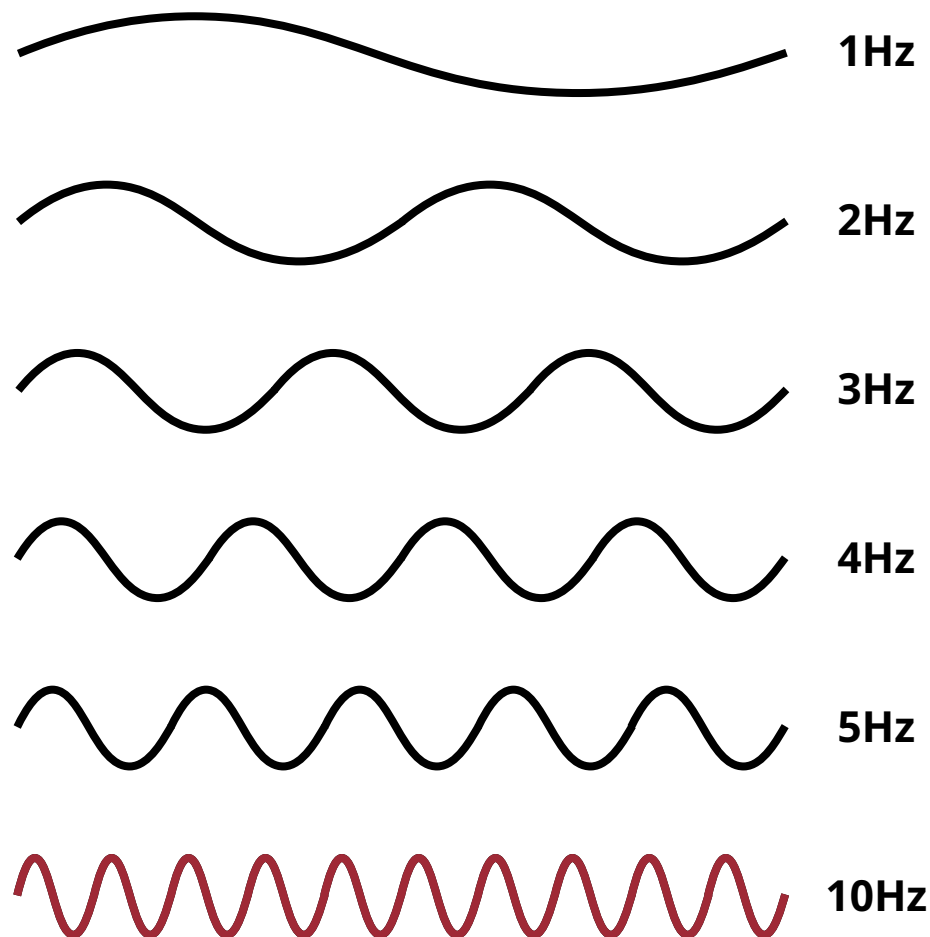
Candidate frequency components



Demystifying Fourier Transform



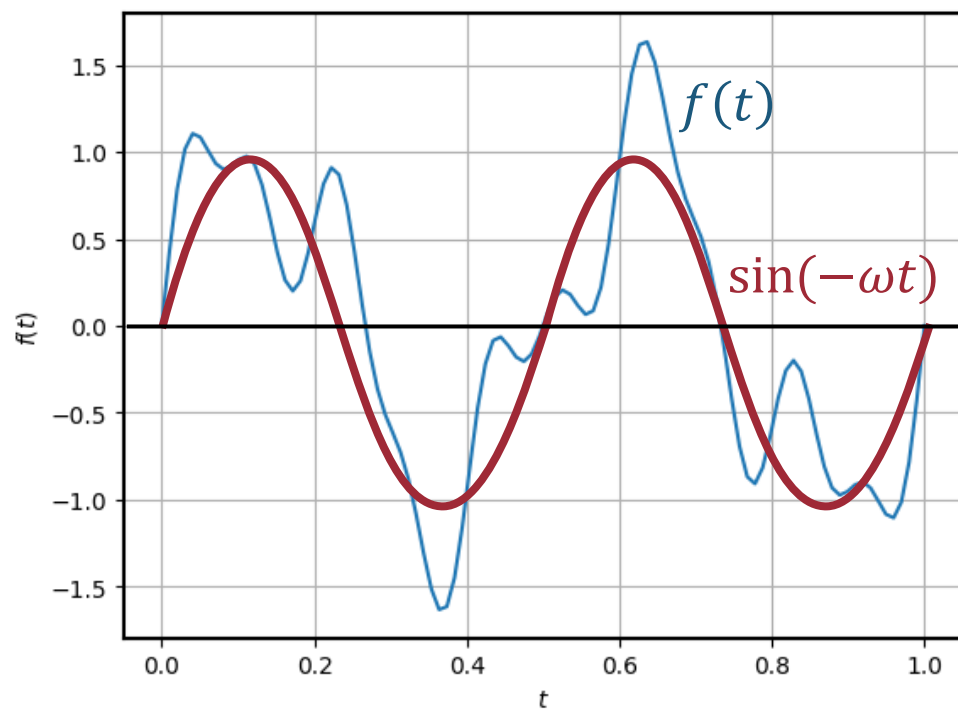
Candidate frequency components



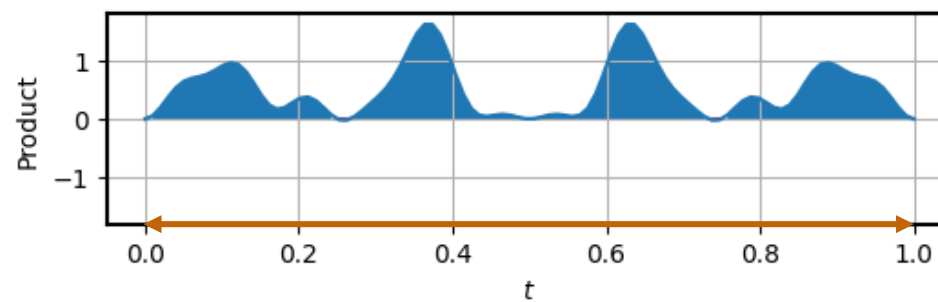
Demystifying Fourier Transform

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

Sum over all t



$f(t) \sin(-\omega t)$

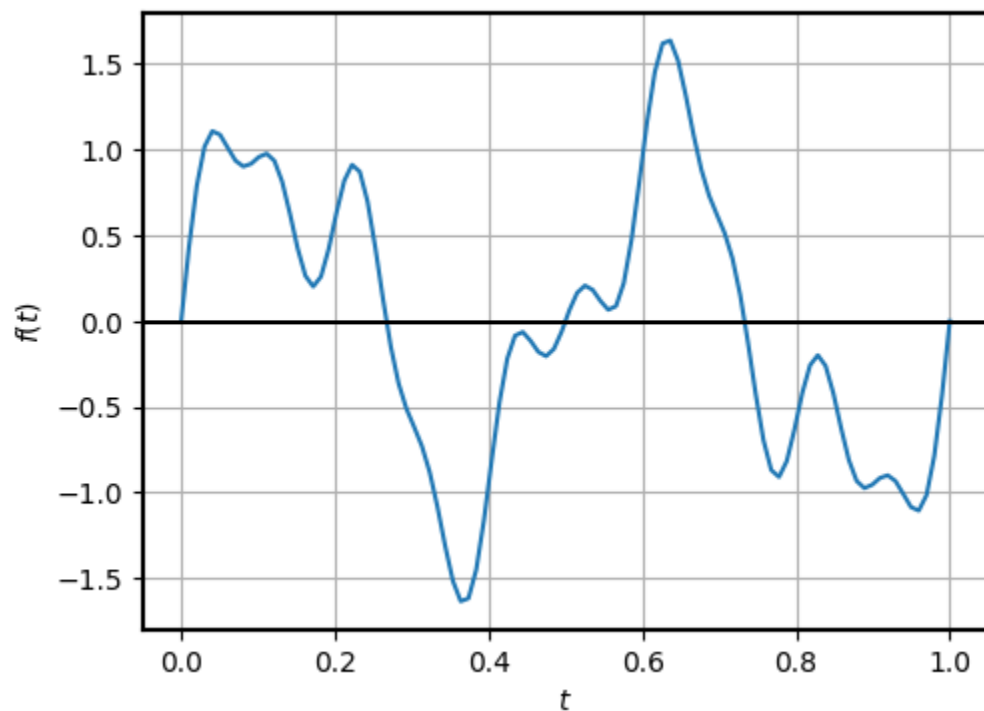


Sum = 0.495

Demystifying Fourier Transform

Signal

(time-domain)

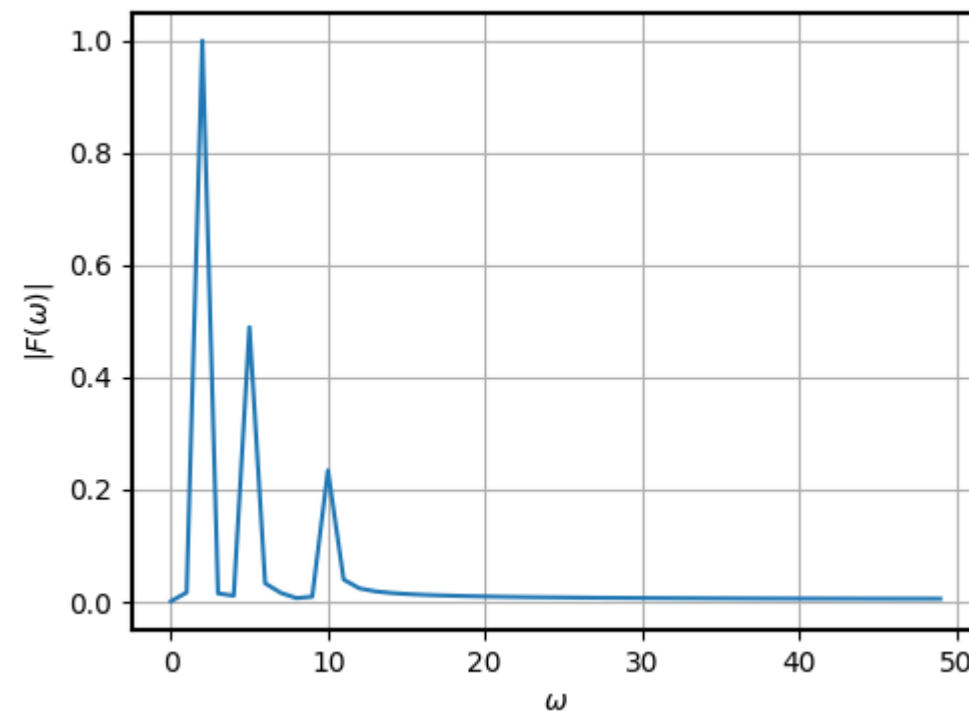


Fourier Transform



Spectrum

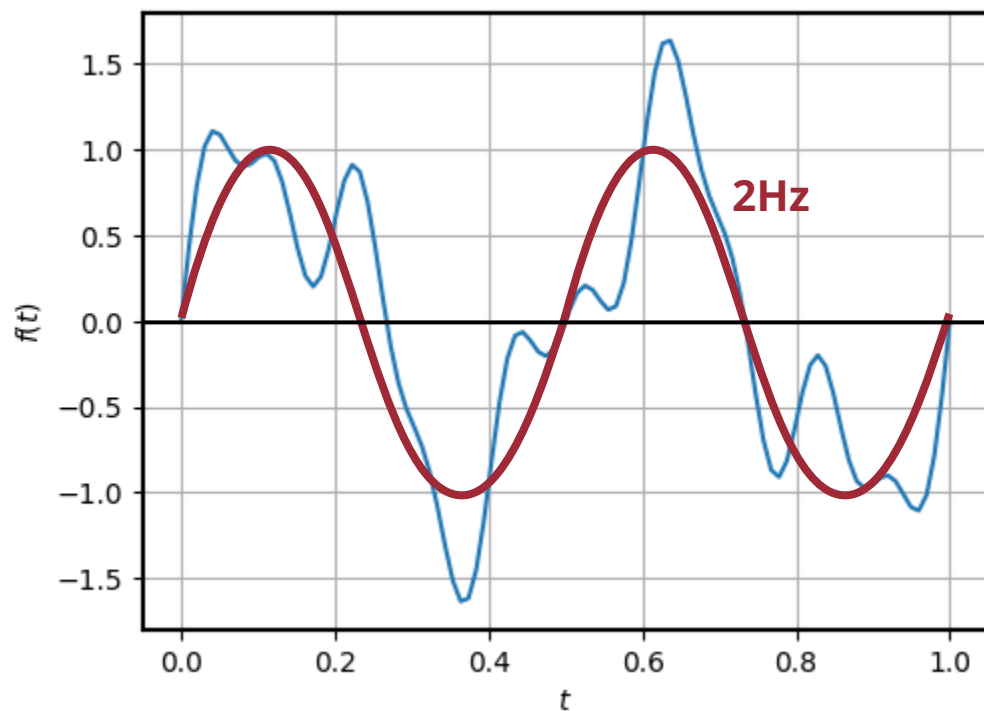
(frequency-domain)



Demystifying Fourier Transform

Signal

(time-domain)

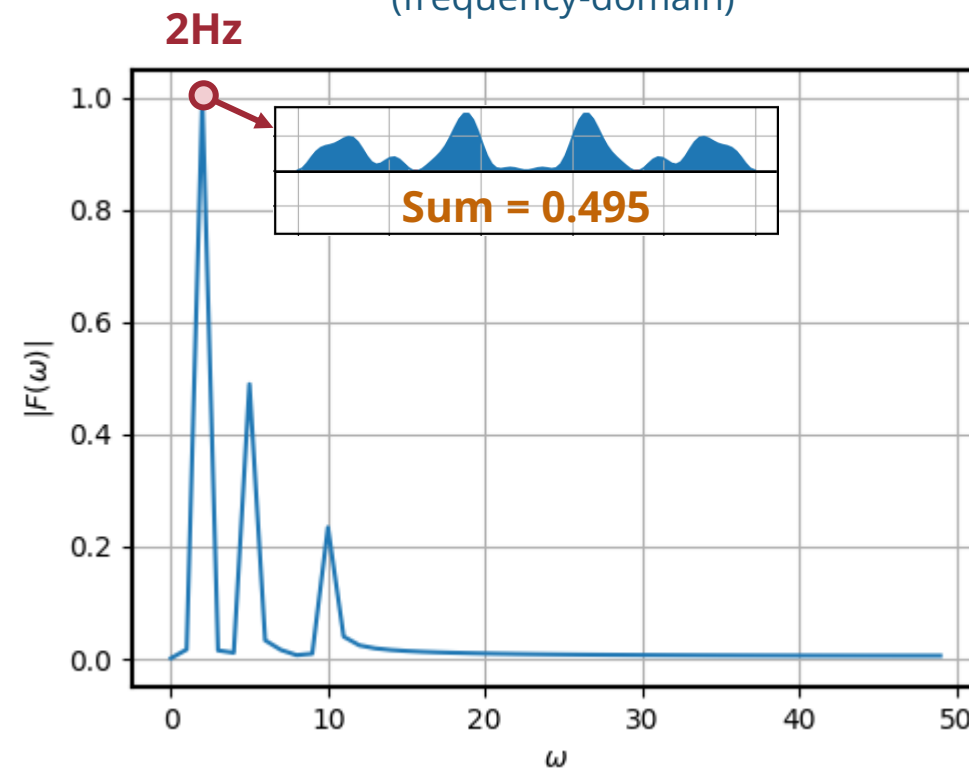


Fourier Transform



Spectrum

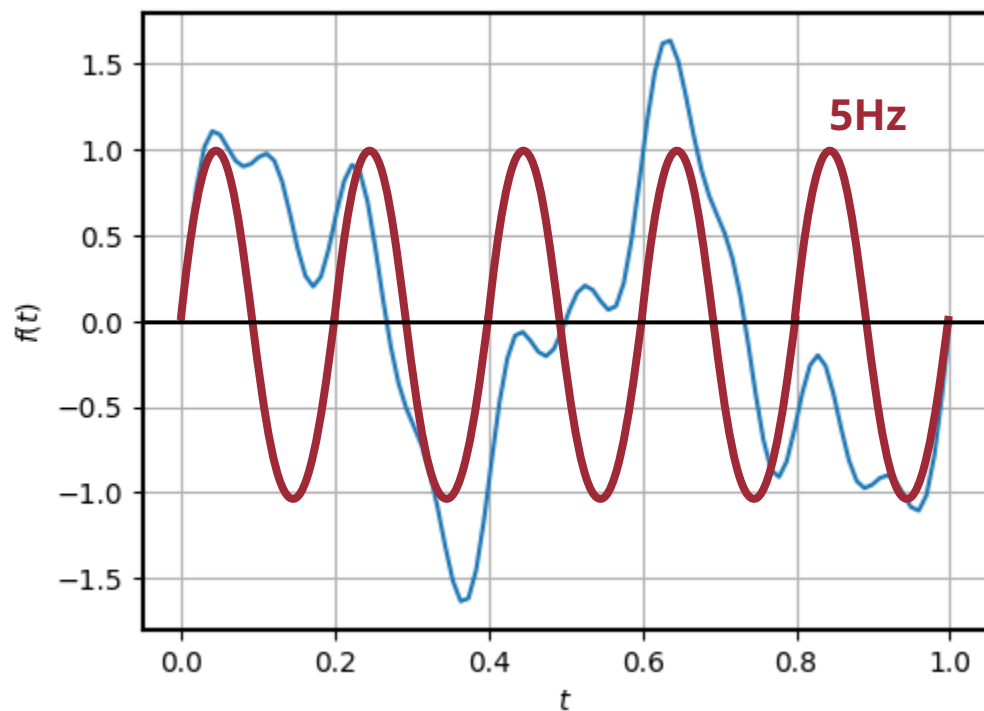
(frequency-domain)



Demystifying Fourier Transform

Signal

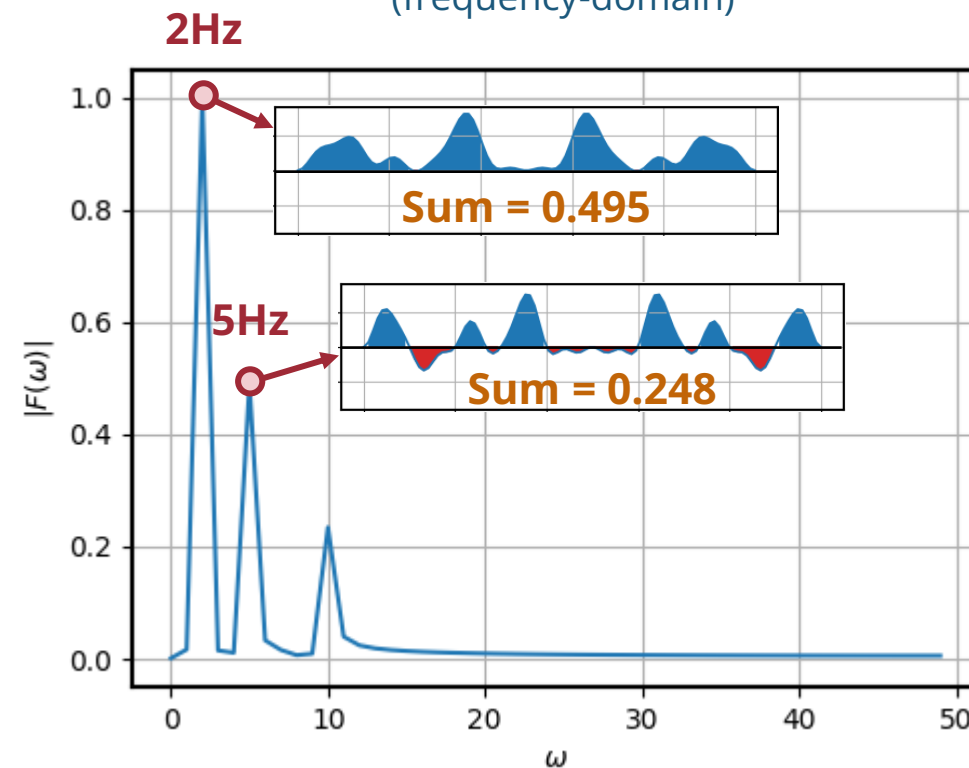
(time-domain)



Fourier Transform
→

Spectrum

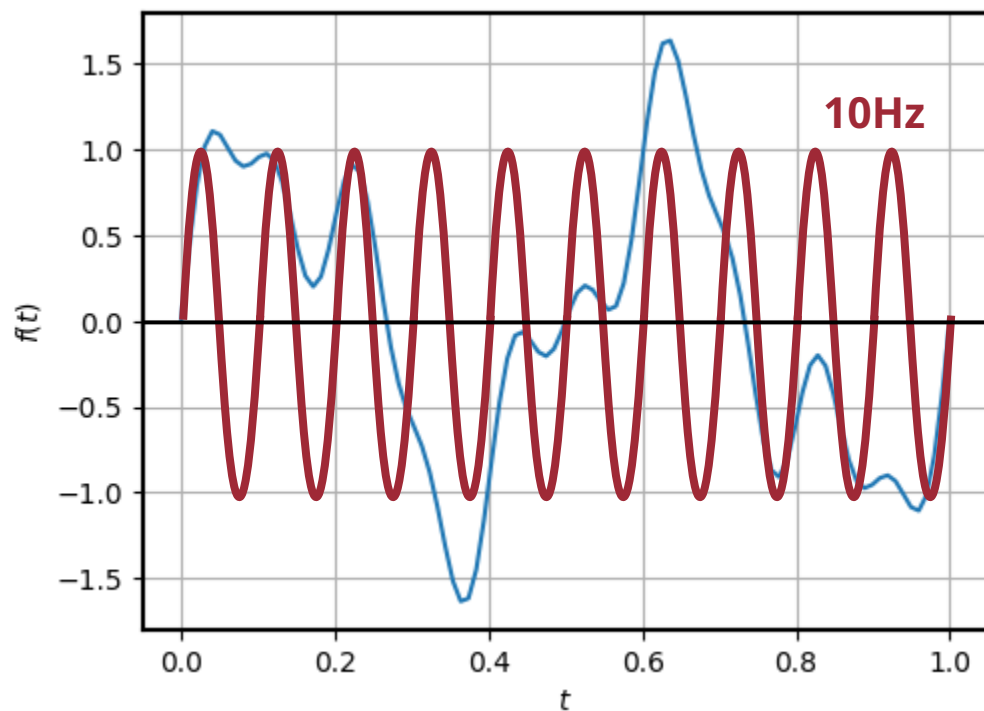
(frequency-domain)



Demystifying Fourier Transform

Signal

(time-domain)

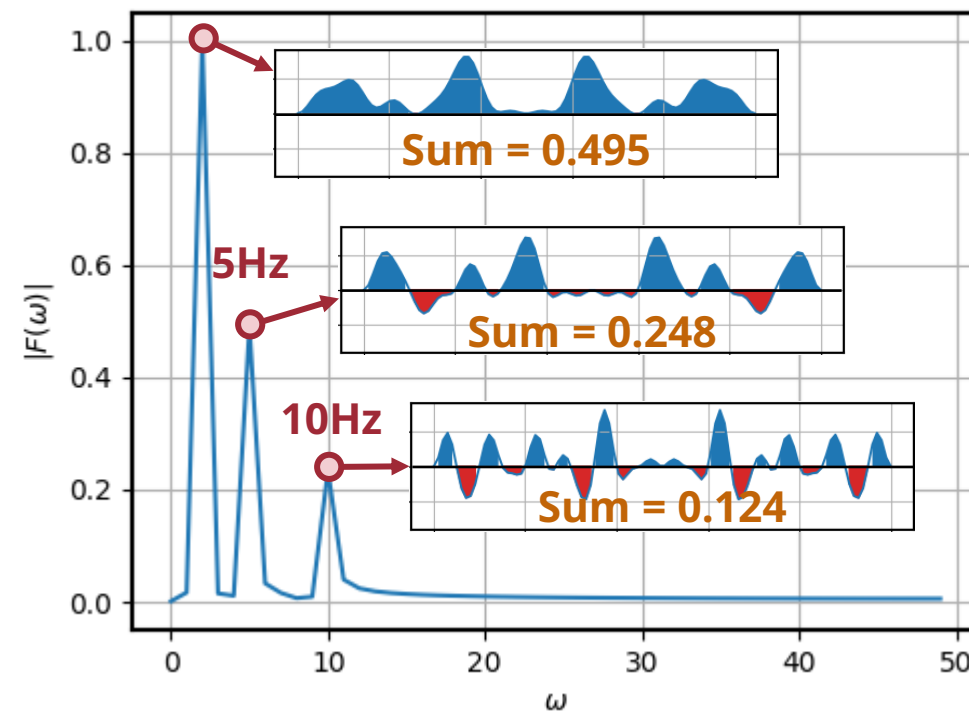


Fourier Transform



Spectrum

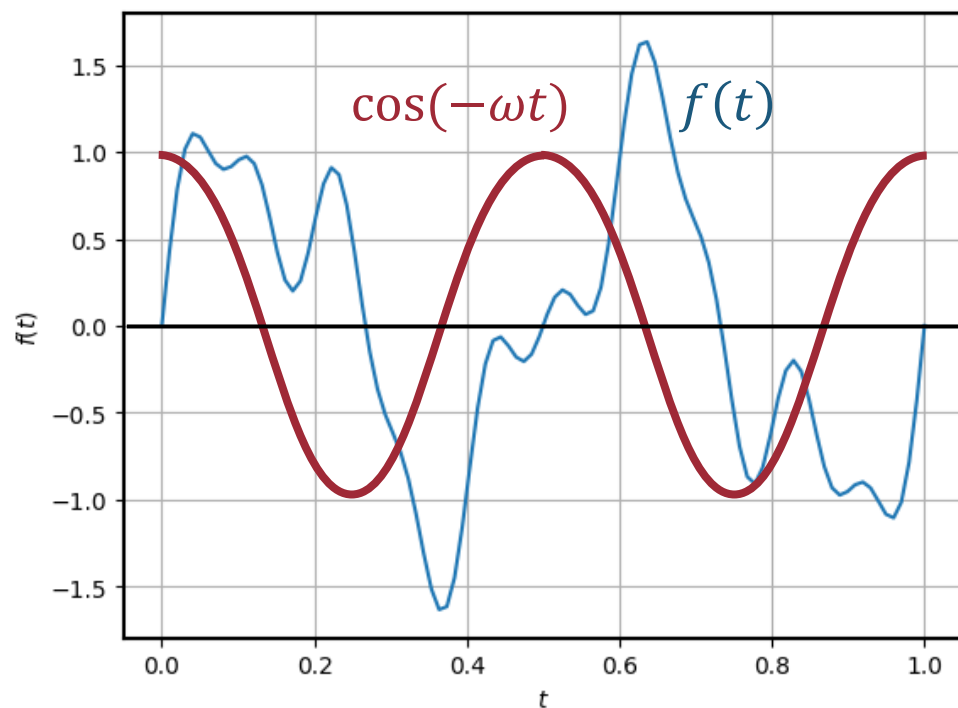
(frequency-domain)



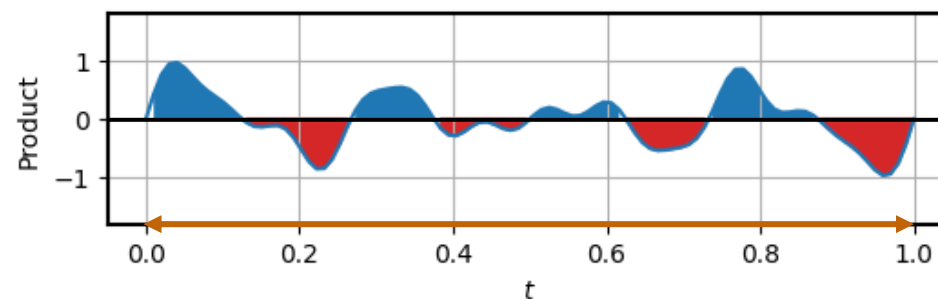
Demystifying Fourier Transform

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

Sum over all t

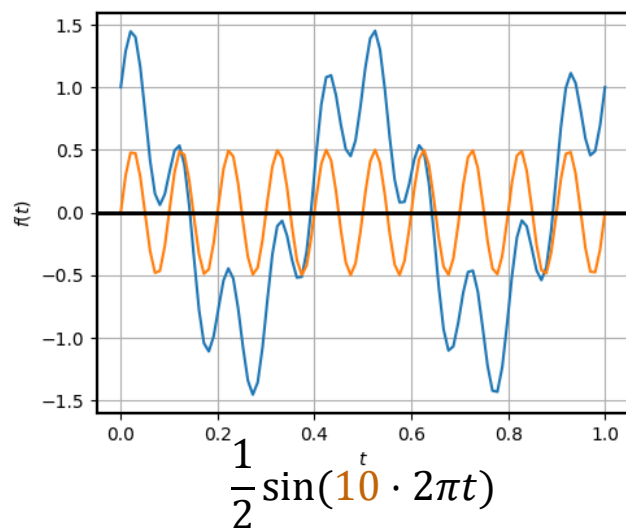
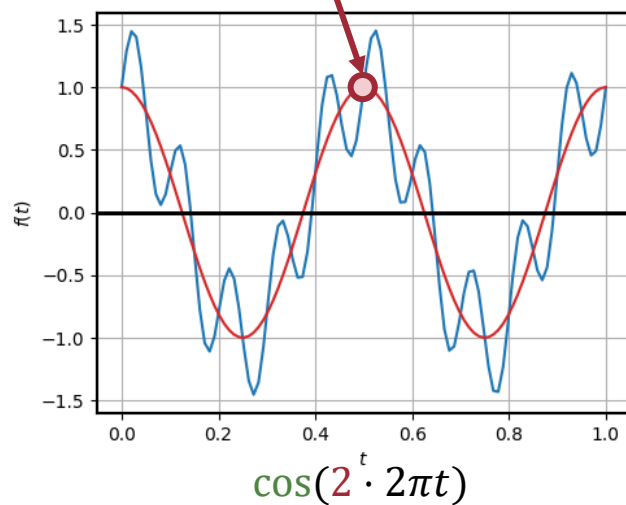
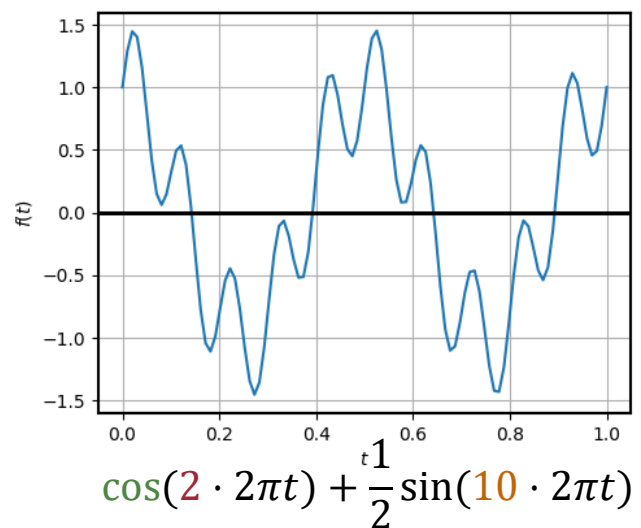
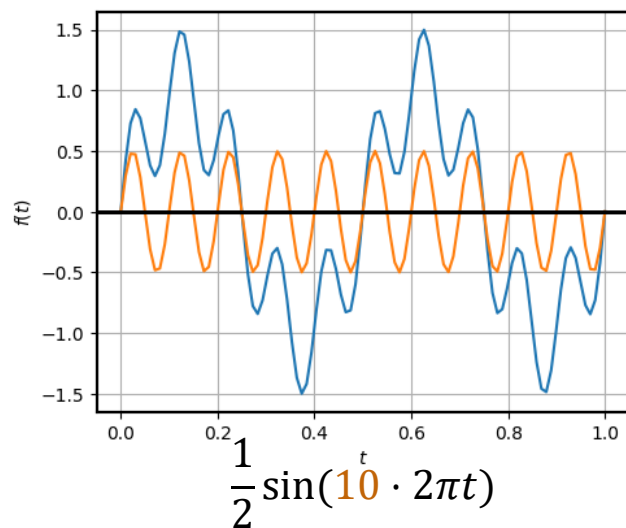
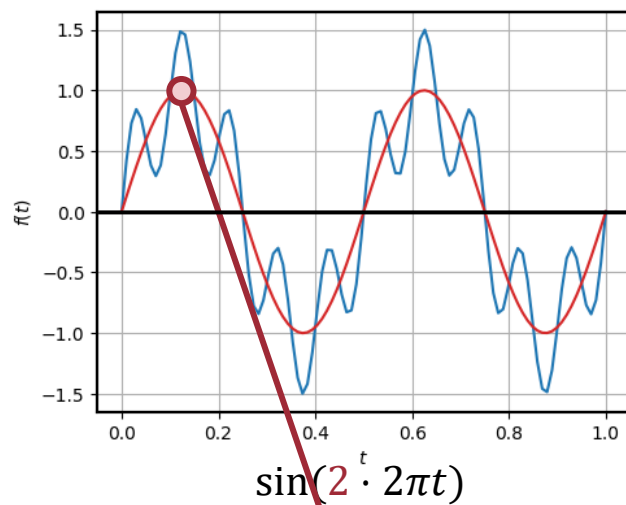
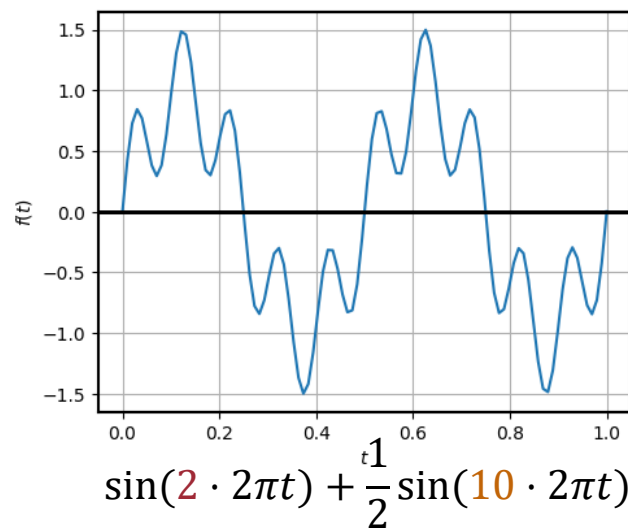


$f(t) \cos(-\omega t)$



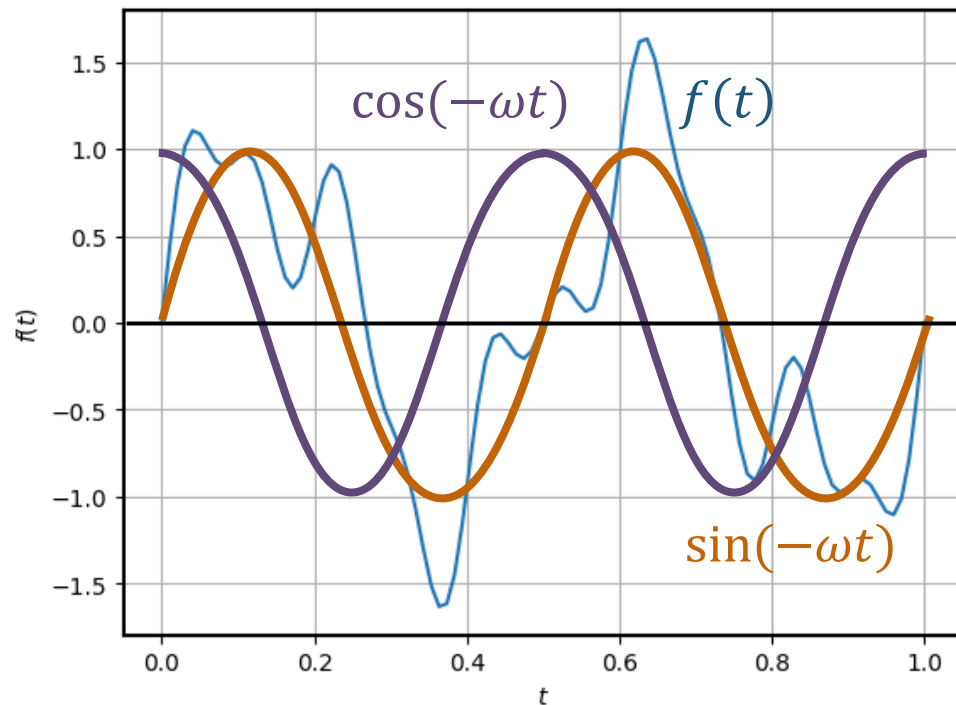
Sum = 0

Phase



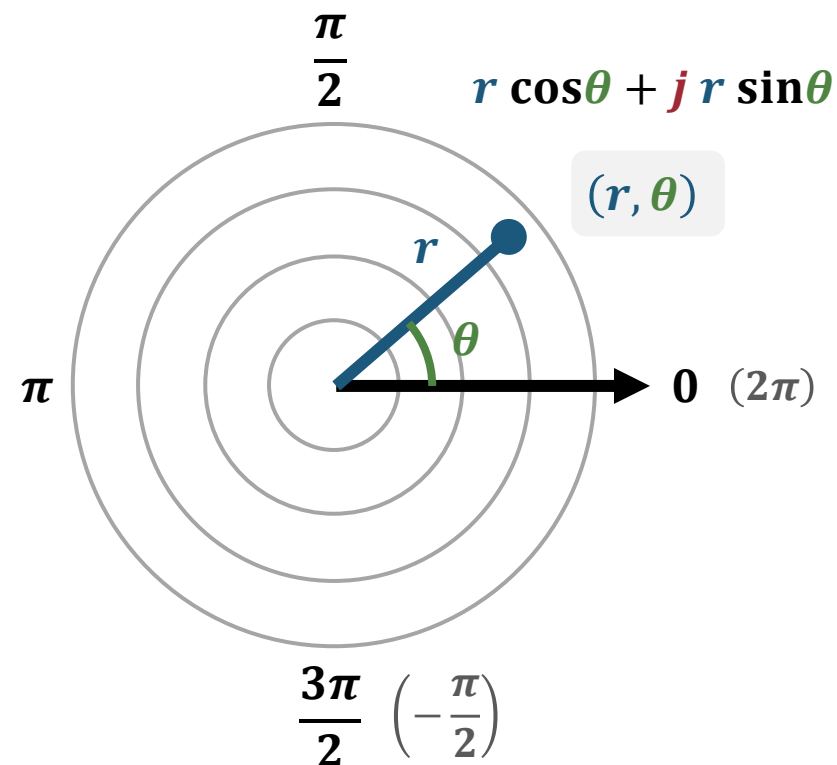
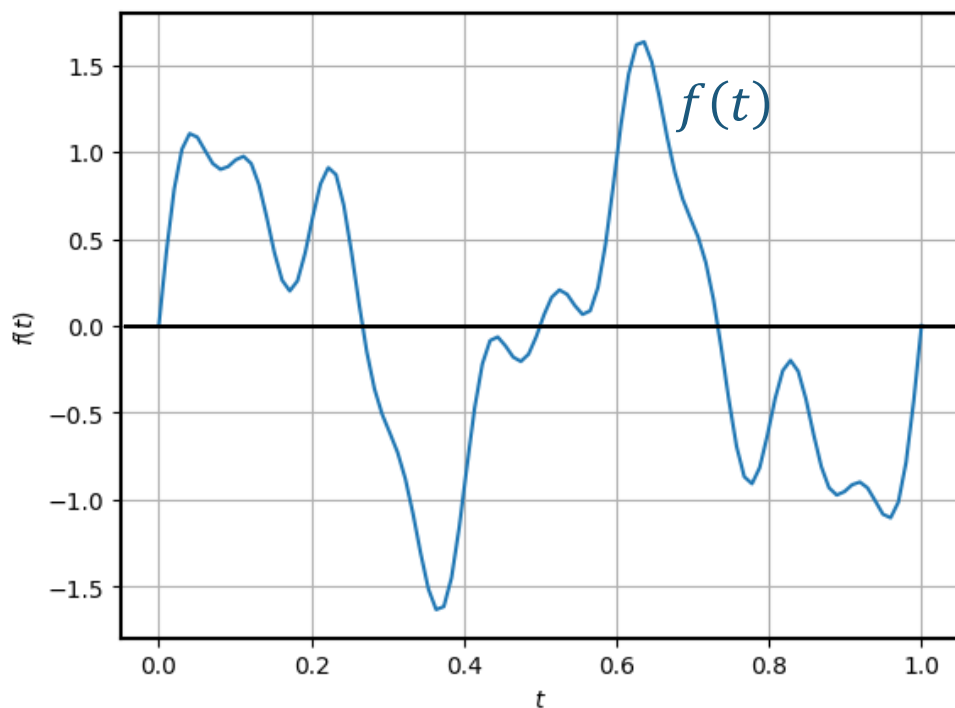
Demystifying Fourier Transform

$$F(\omega) = \int_{-\infty}^{\infty} \underbrace{f(t) \cos(-\omega t)}_{\text{Real part}} + \underbrace{j f(t) \sin(-\omega t)}_{\text{Imaginary part}} dt$$



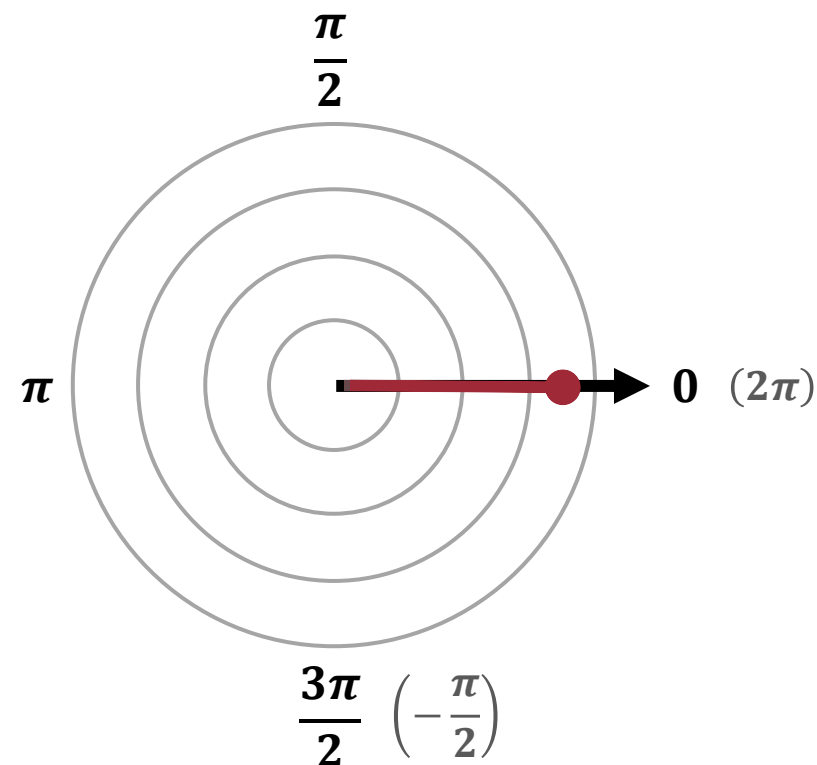
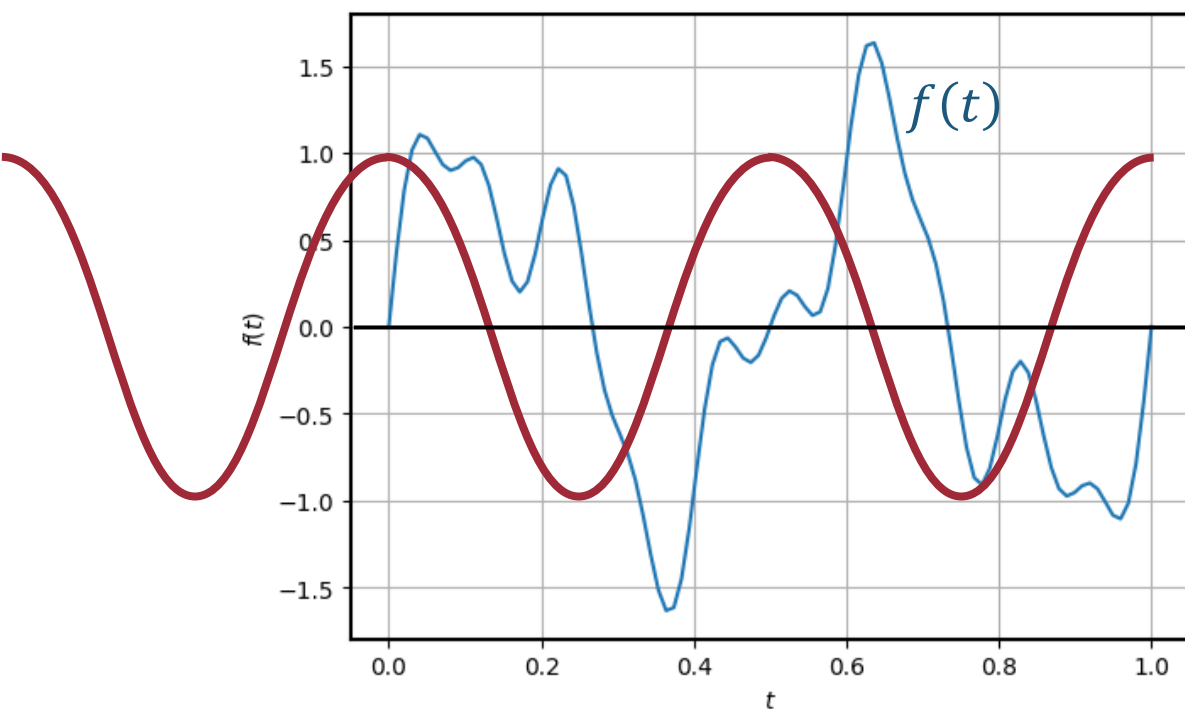
Demystifying Fourier Transform

$$F(\omega) = \int_{-\infty}^{\infty} \underbrace{f(t) \cos(-\omega t)}_{\text{Real part}} + \underbrace{j f(t) \sin(-\omega t)}_{\text{Imaginary part}} dt$$



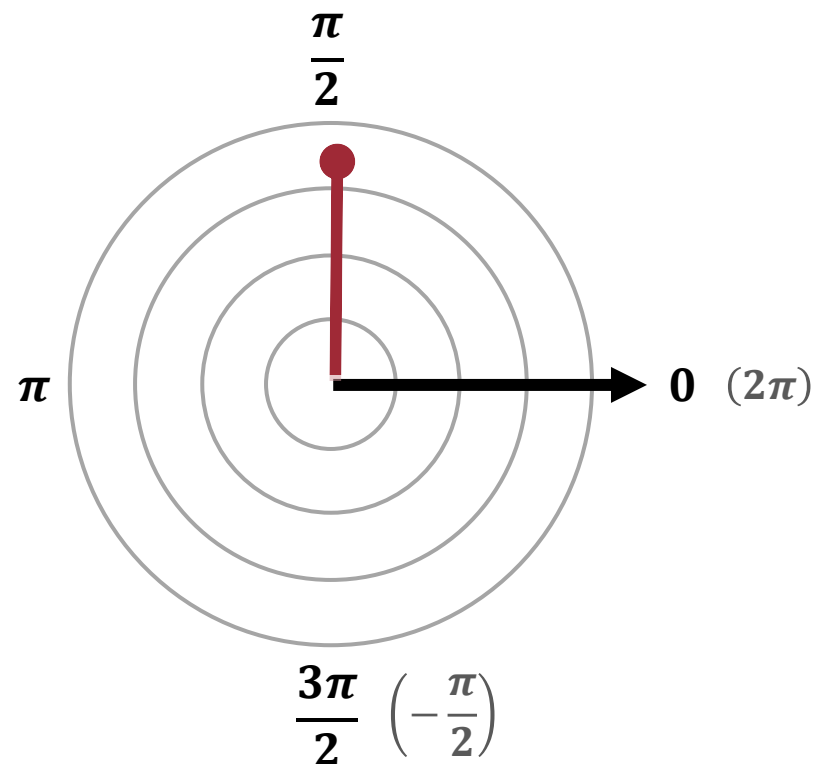
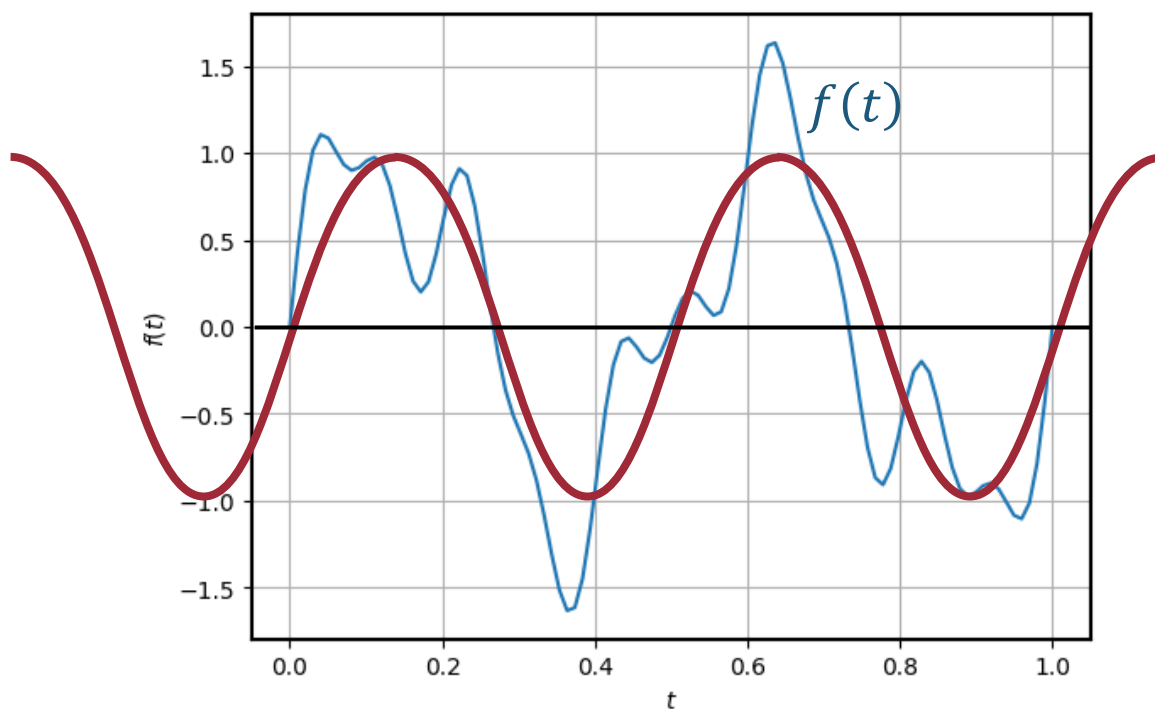
Demystifying Fourier Transform

$$F(\omega) = \int_{-\infty}^{\infty} \underbrace{f(t) \cos(-\omega t)}_{\text{Real part}} + \underbrace{j f(t) \sin(-\omega t)}_{\text{Imaginary part}} dt$$



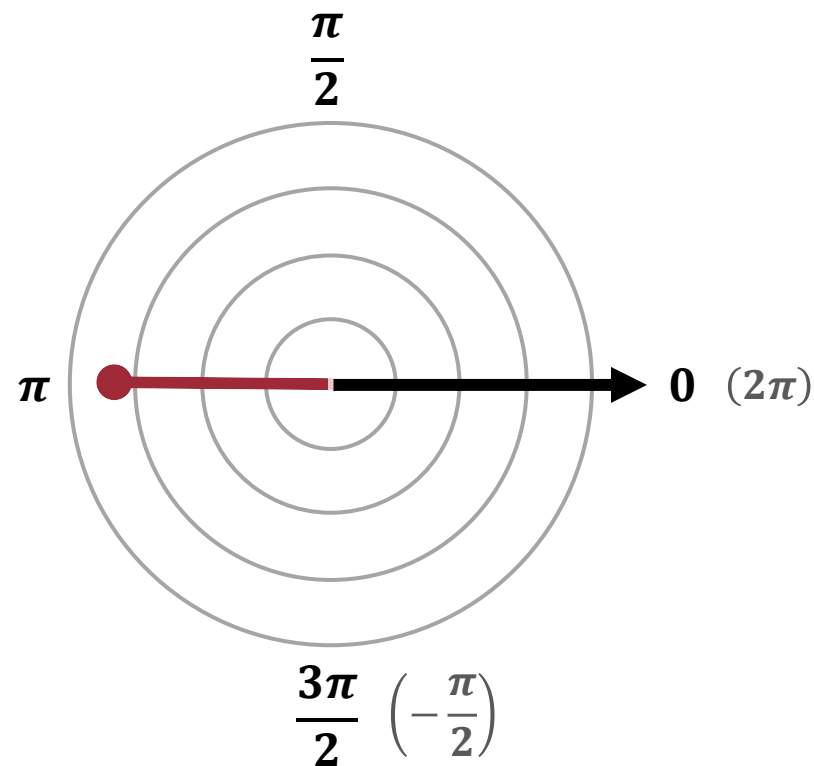
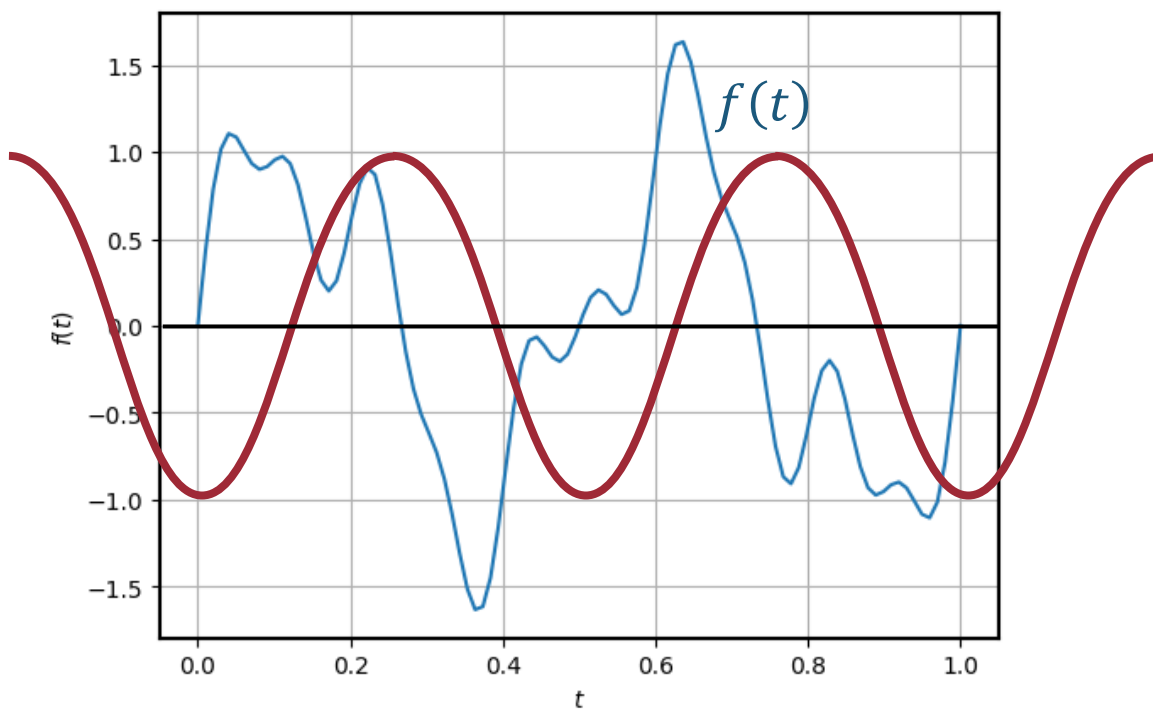
Demystifying Fourier Transform

$$F(\omega) = \int_{-\infty}^{\infty} \underbrace{f(t) \cos(-\omega t)}_{\text{Real part}} + \underbrace{j f(t) \sin(-\omega t)}_{\text{Imaginary part}} dt$$



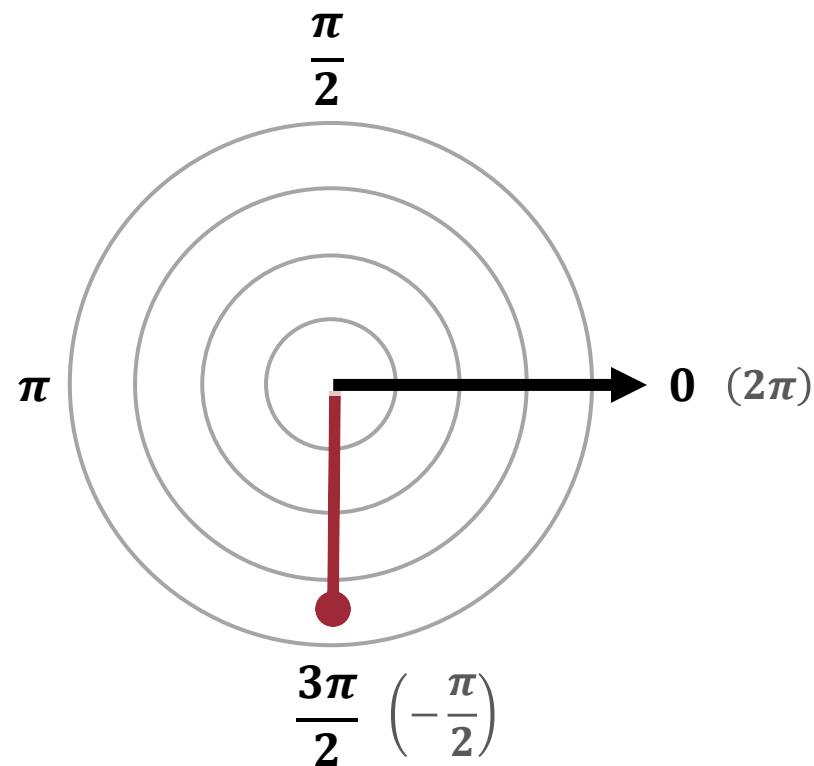
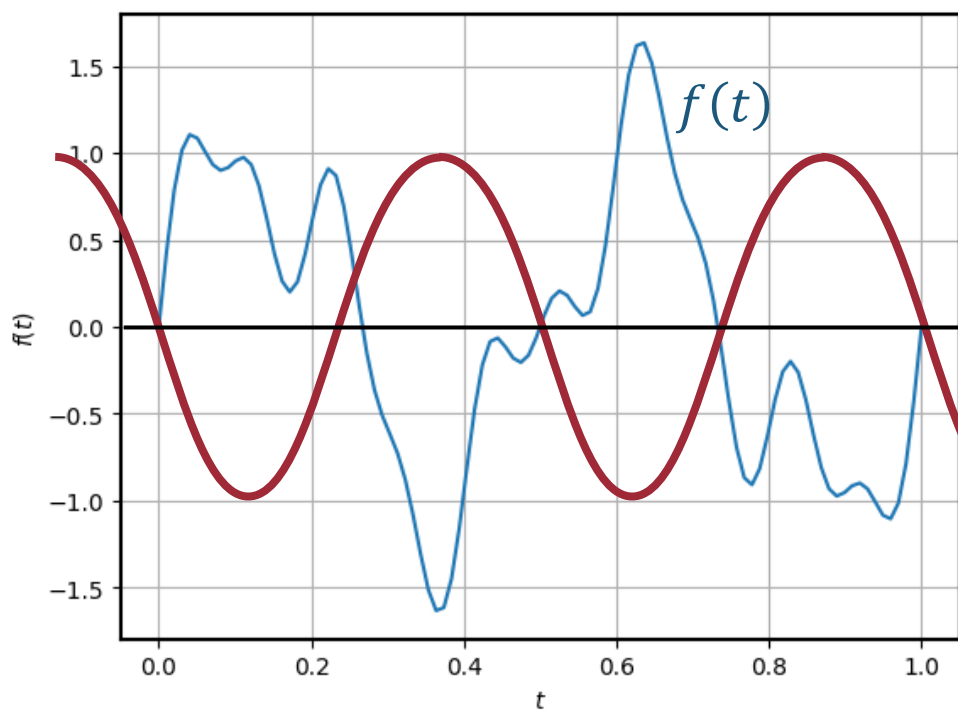
Demystifying Fourier Transform

$$F(\omega) = \int_{-\infty}^{\infty} \underbrace{f(t) \cos(-\omega t)}_{\text{Real part}} + \underbrace{j f(t) \sin(-\omega t)}_{\text{Imaginary part}} dt$$

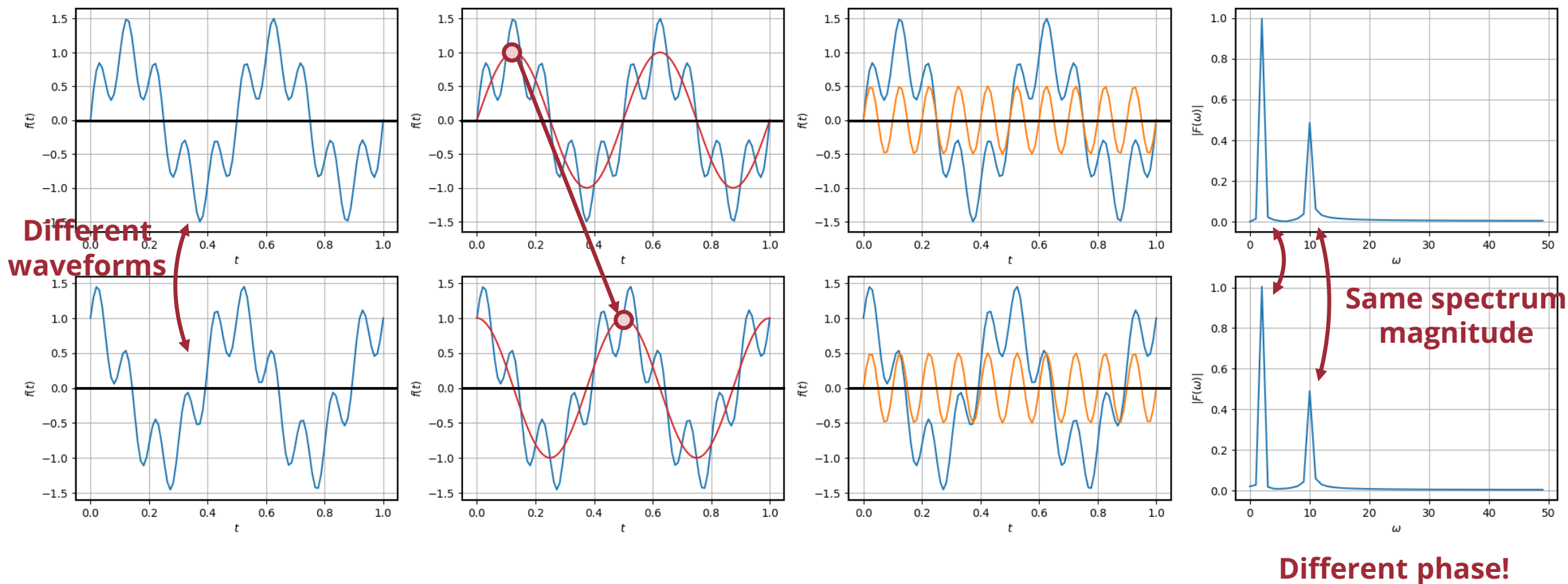


Demystifying Fourier Transform

$$F(\omega) = \int_{-\infty}^{\infty} \underbrace{f(t) \cos(-\omega t)}_{\text{Real part}} + \underbrace{j f(t) \sin(-\omega t)}_{\text{Imaginary part}} dt$$



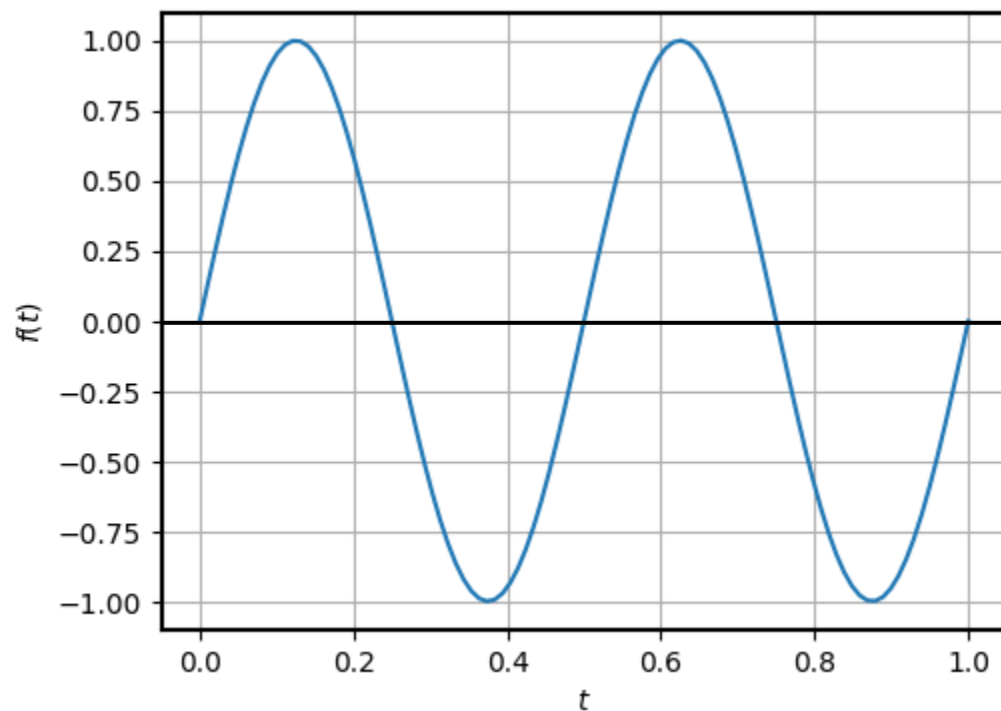
Magnitude & Phase



Example: A 2Hz Sine Wave

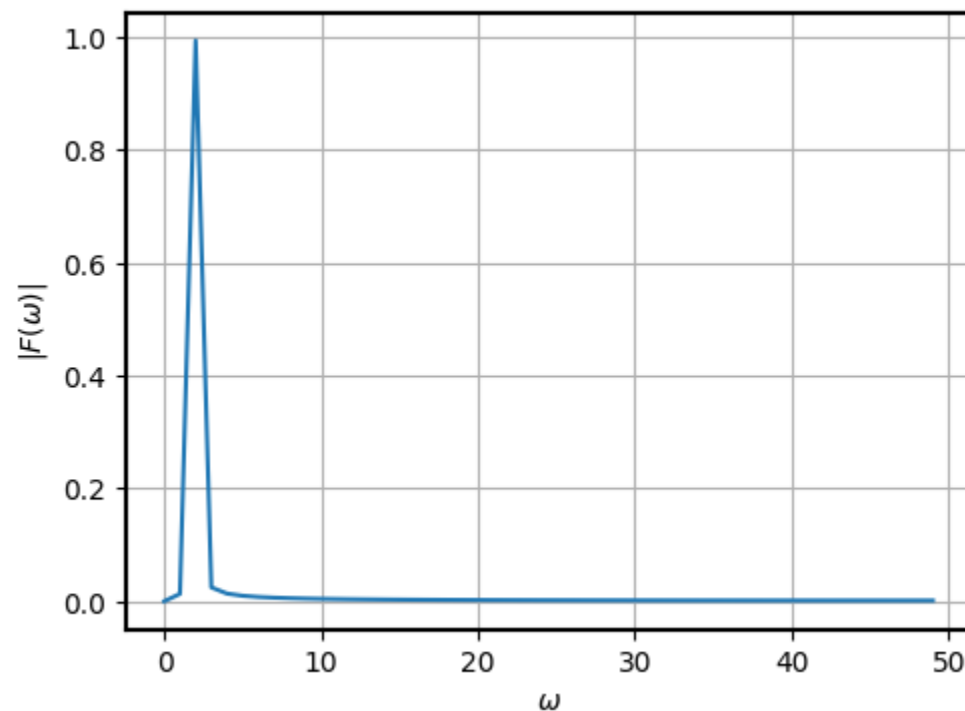
Signal

(time-domain)



Spectrum

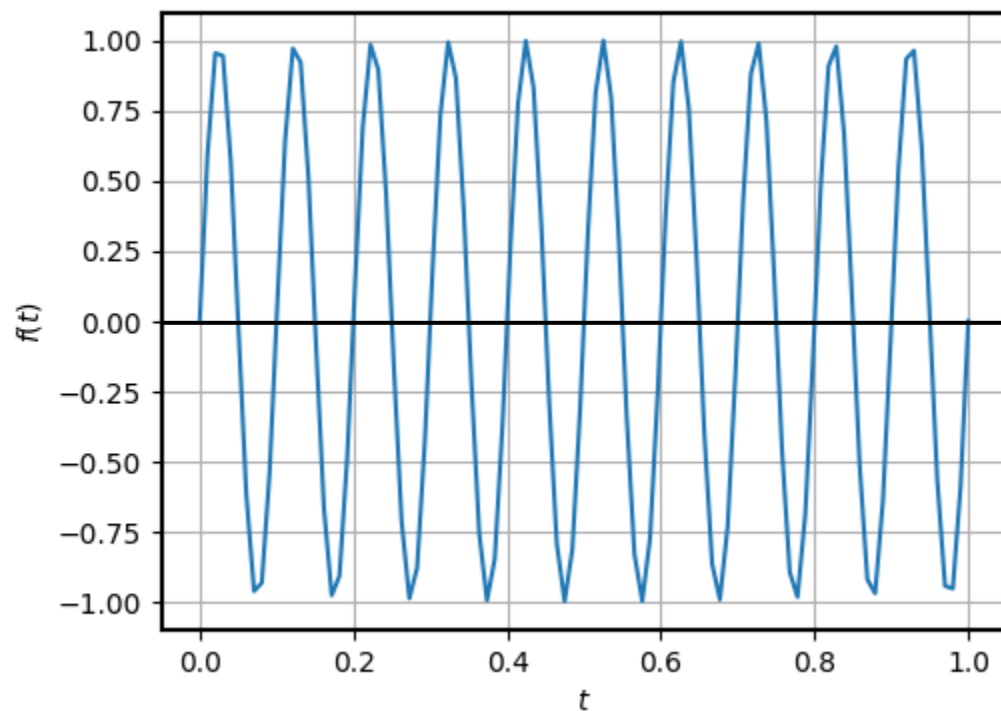
(frequency-domain)



Example: A 10Hz Sine Wave

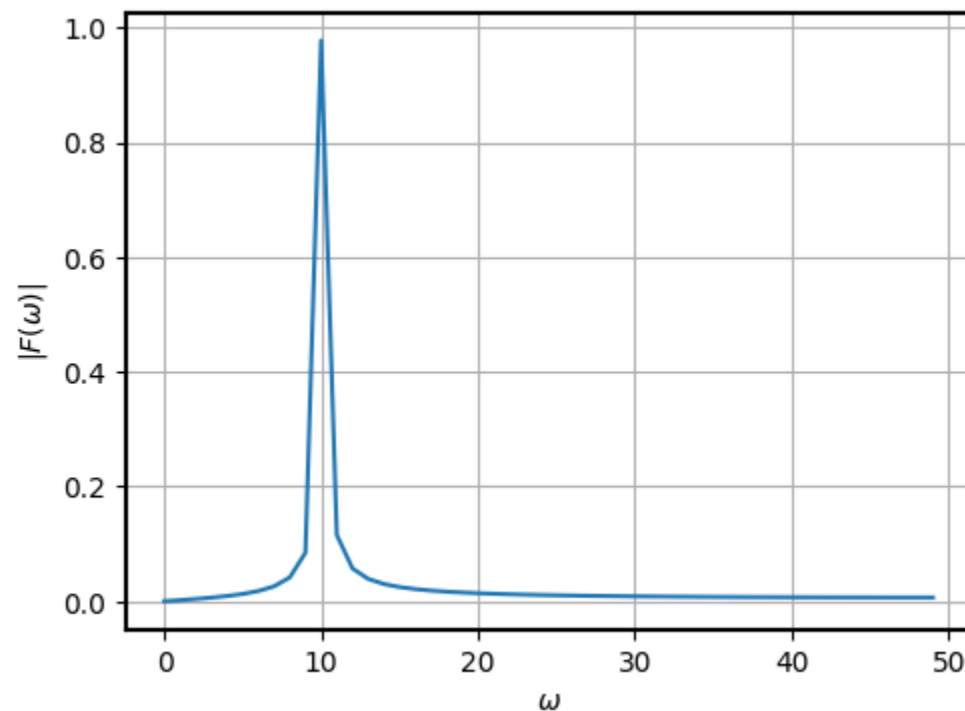
Signal

(time-domain)



Spectrum

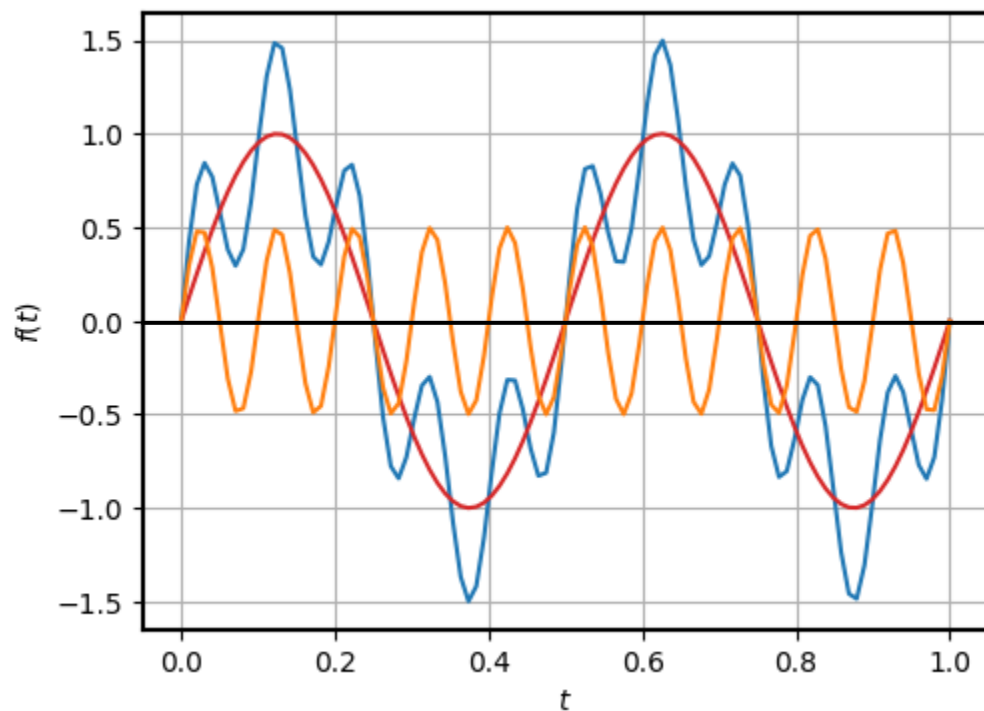
(frequency-domain)



Example: Sum of 2Hz & 10Hz Sine Waves

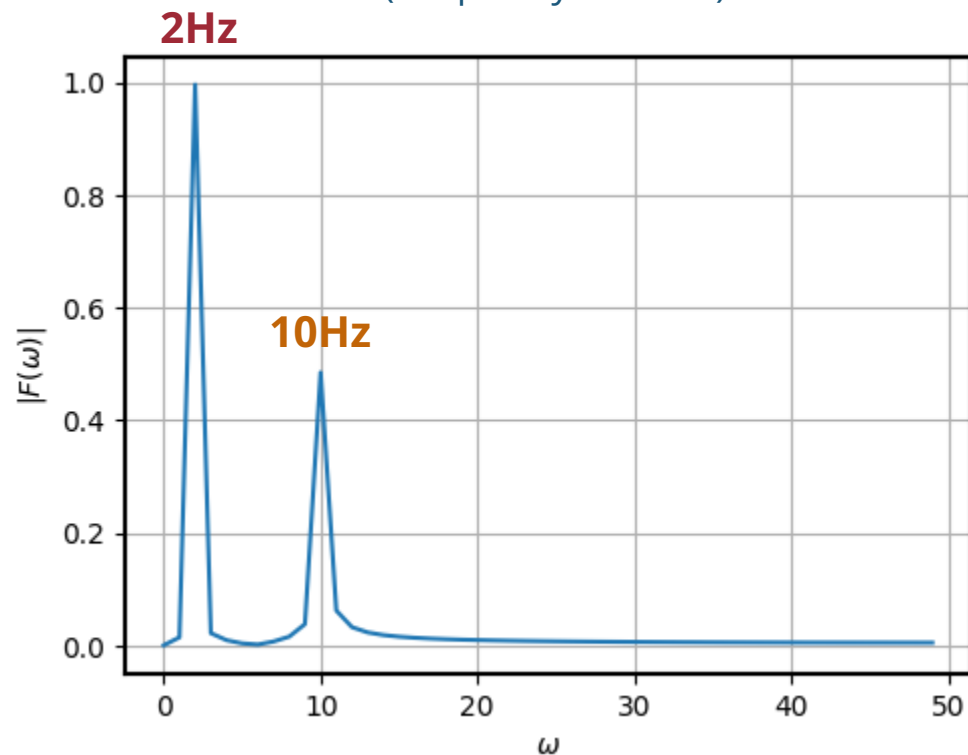
Signal

(time-domain)



Spectrum

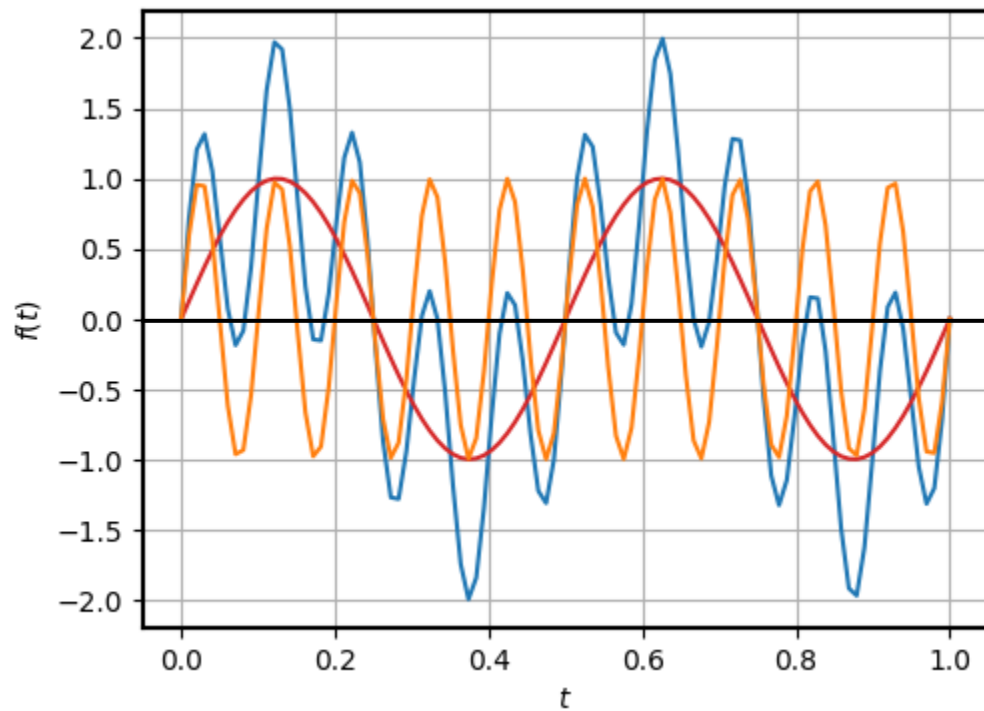
(frequency-domain)



How about this?

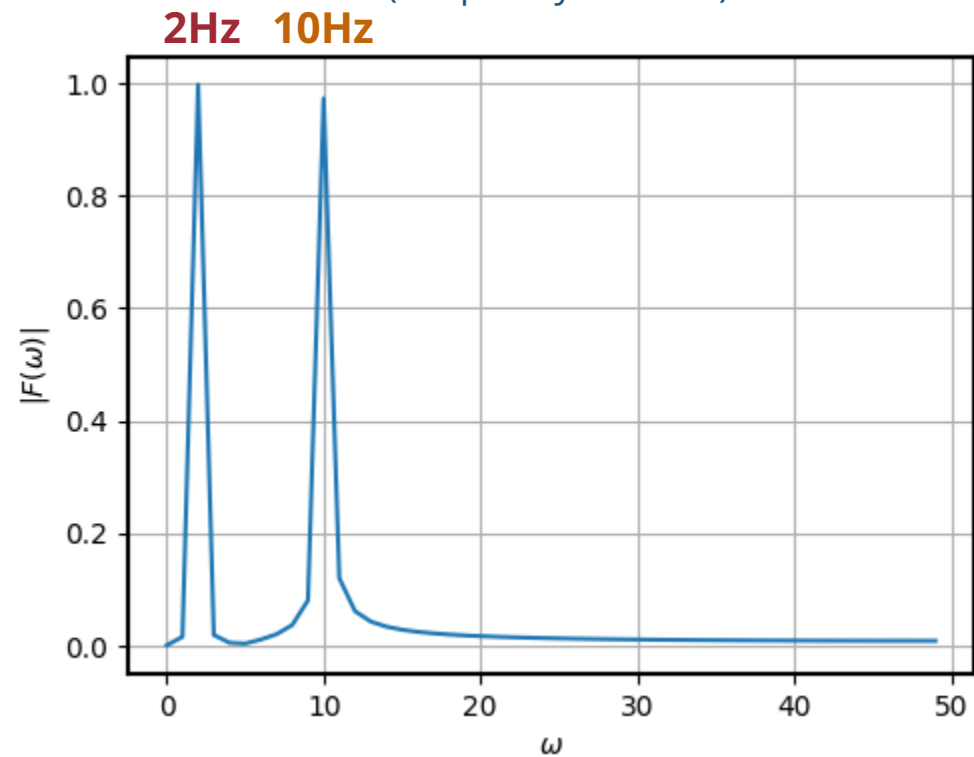
Signal

(time-domain)



Spectrum

(frequency-domain)



Fourier Transform

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

The diagram illustrates the Fourier Transform equation with various components and annotations:

- Output spectrum:** A green box containing $F(\omega)$ with an upward arrow pointing to the text "Output spectrum".
- Frequency:** An orange arrow points from the ω in $F(\omega)$ to the text "Frequency".
- Sum over all t :** A purple arrow points from the integral symbol $\int_{-\infty}^{\infty}$ to the text "Sum over all t ".
- Input signal:** A blue box containing $f(t)$ with an upward arrow pointing to the text "Input signal".
- Sine and cosine waves of frequency ω :** A red box containing $e^{-j\omega t}$ with a red arrow pointing to the text "Sine and cosine waves of frequency ω ".
- dt :** A purple box containing dt with a purple arrow pointing to it.

The equation is represented as:
$$F(\omega) = \int_{-\infty}^{\infty} f(t) e^{-j\omega t} dt$$

Fourier Transform

- **Intuition:** **Analysis** through **resynthesis**!

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

The diagram illustrates the Fourier Transform equation $F(\omega) = \int_{-\infty}^{\infty} f(t) e^{-j\omega t} dt$ with arrows indicating the direction of the process. A red arrow points from the term $e^{-j\omega t}$ to the word "Analysis" above it. A purple arrow points from the term dt to the word "Synthesis" below it. Another purple arrow points from the integral symbol $\int_{-\infty}^{\infty}$ to the word "Synthesis" below it.

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:

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

Fourier Transform vs. Discrete Fourier Transform

Fourier Transform

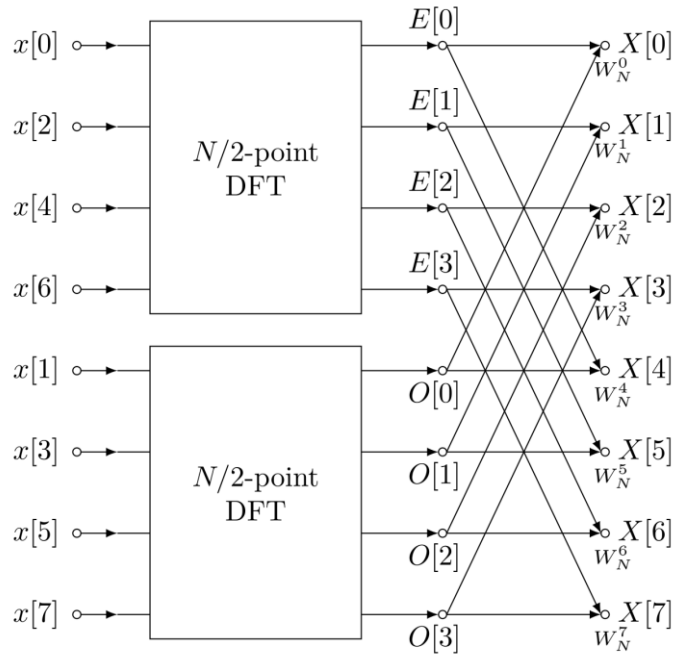
$$F(\omega) = \int_{-\infty}^{\infty} f(t) e^{-j\omega t} dt$$

Discrete Fourier Transform

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

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

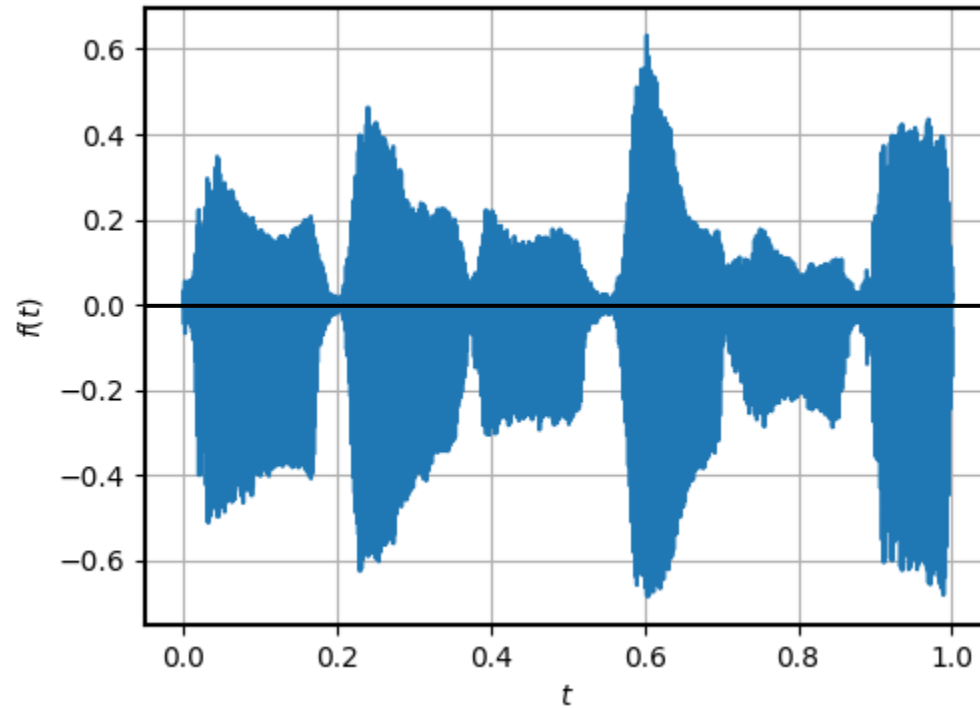
IEEE COMPUTER SOCIETY
www.computer.org/ise

Time-Frequency Analysis

Fourier Transform of a Trumpet Sound

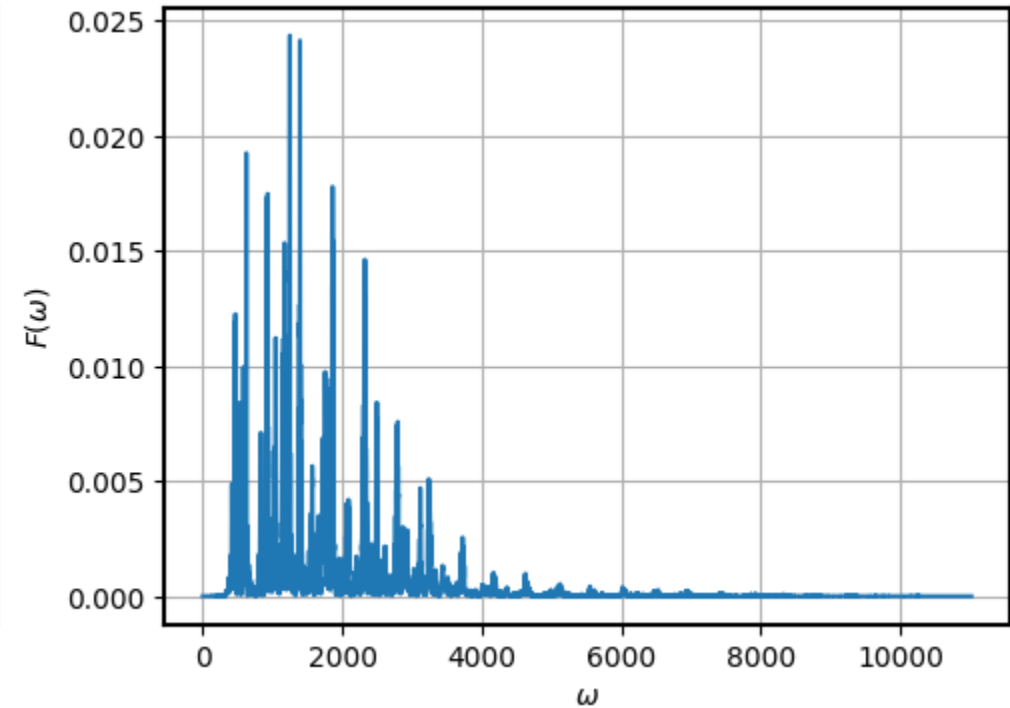
Signal

(time-domain)



Spectrum

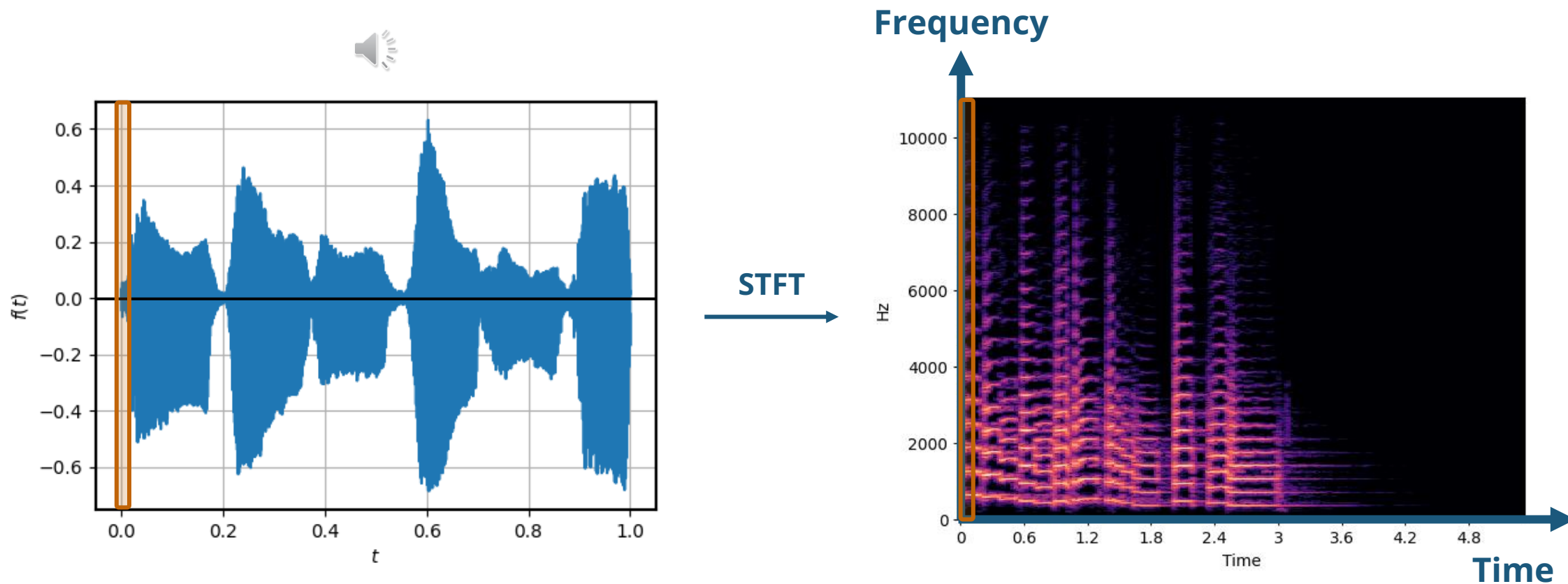
(frequency-domain)



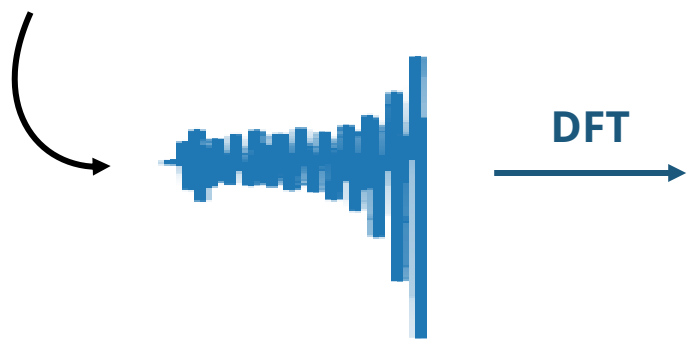
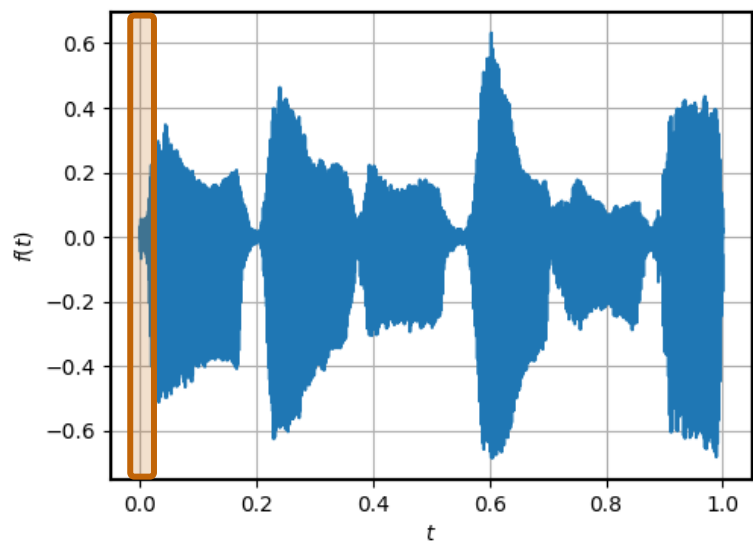
Fourier Transform cannot localize! 🙄

Short-Time Fourier Transform (STFT)

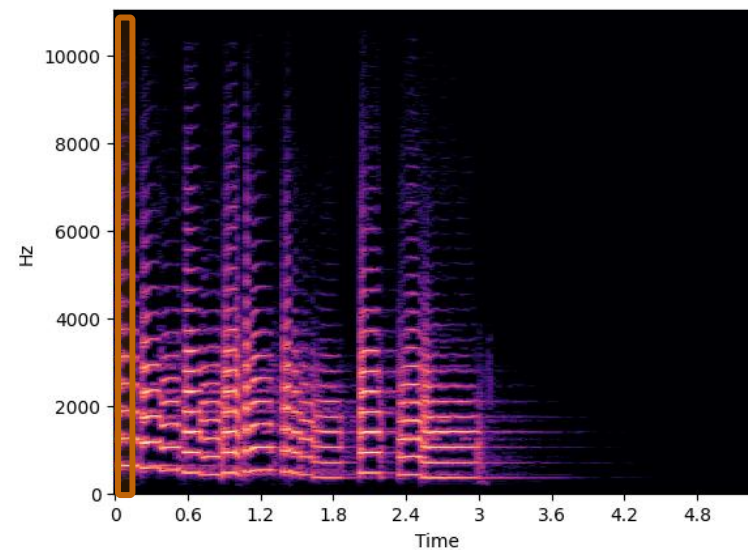
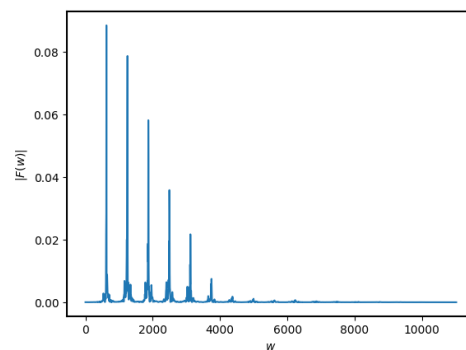
- **Intuition:** Slice the audio into chunks and apply Fourier transform



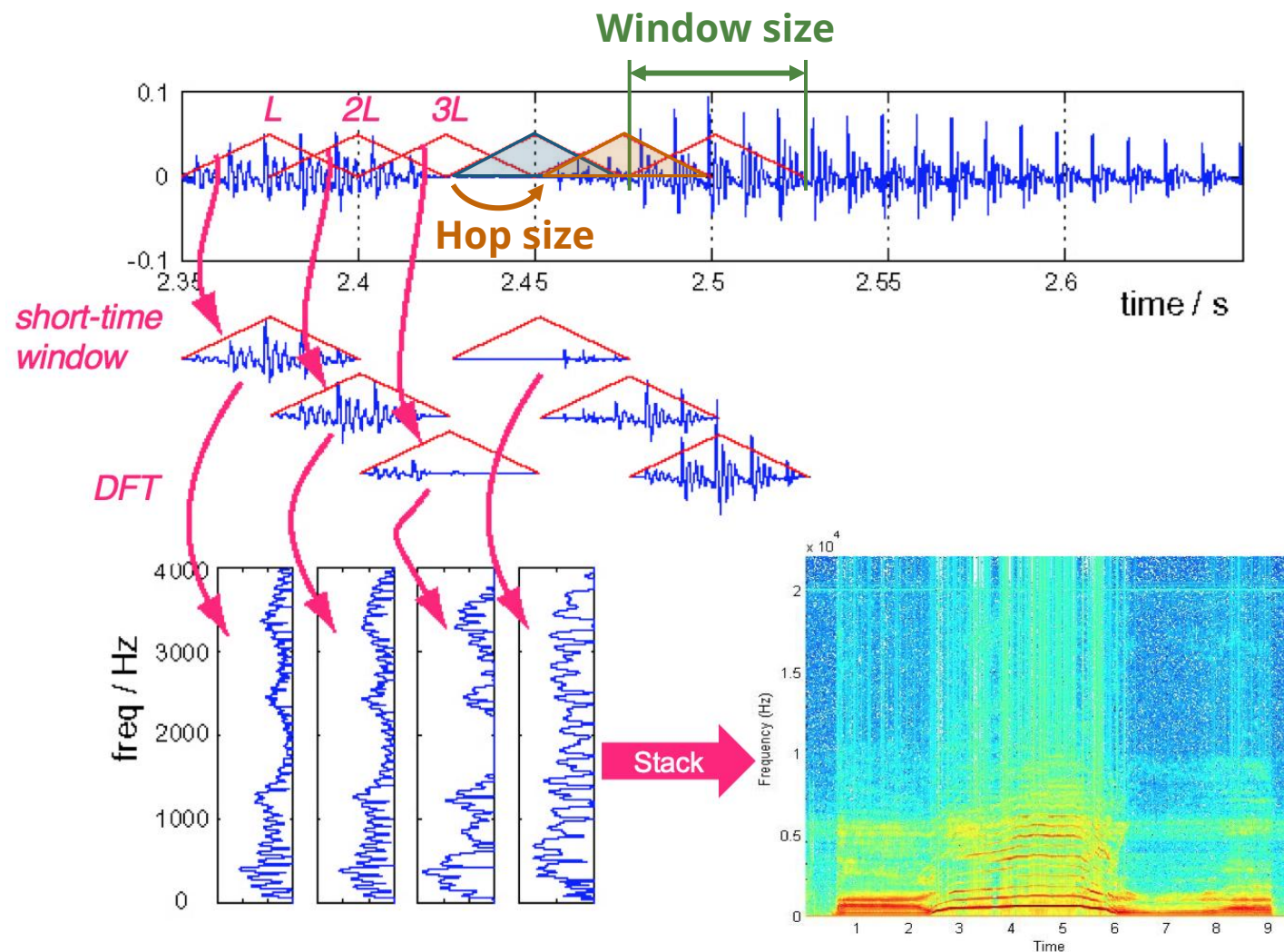
Short-Time Fourier Transform (STFT)



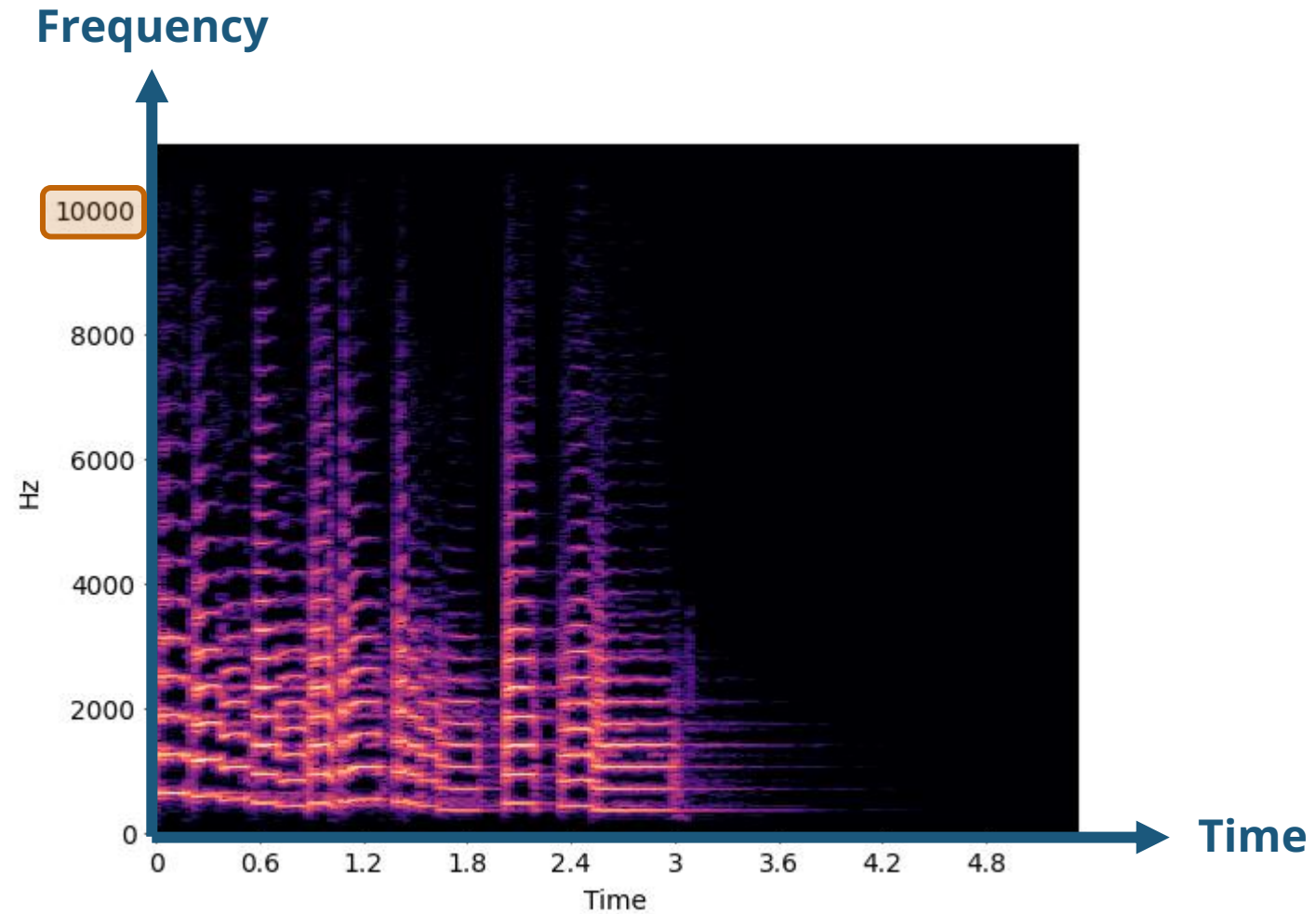
DFT



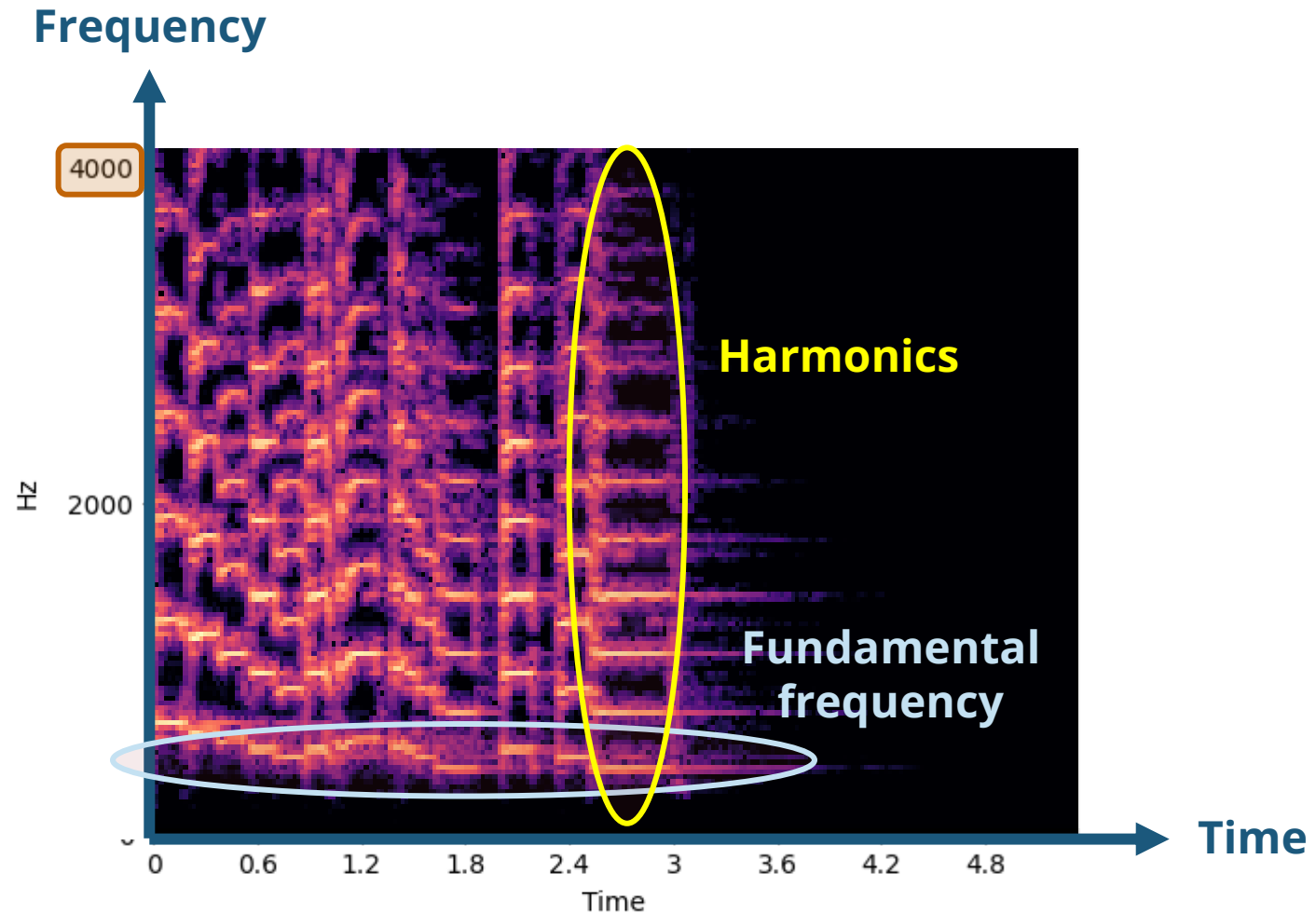
Short-Time Fourier Transform (STFT)



| Spectrogram



| Spectrogram



Timbre

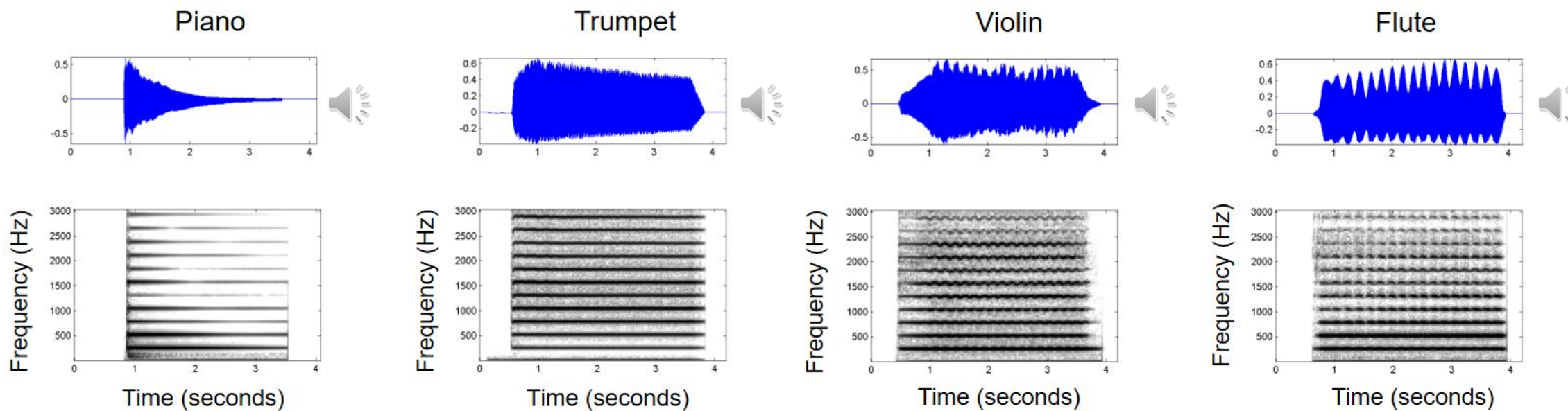
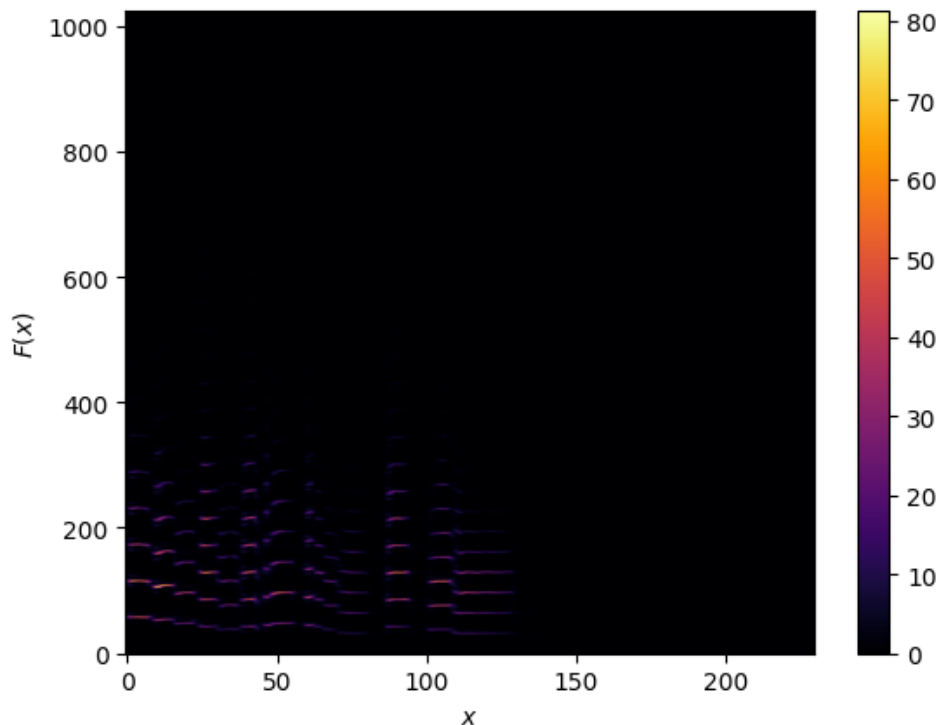


Figure 1.23 from [Müller, FMP, Springer 2015]

(Source: Müller et al., 2021)

Example: `librosa.stft`



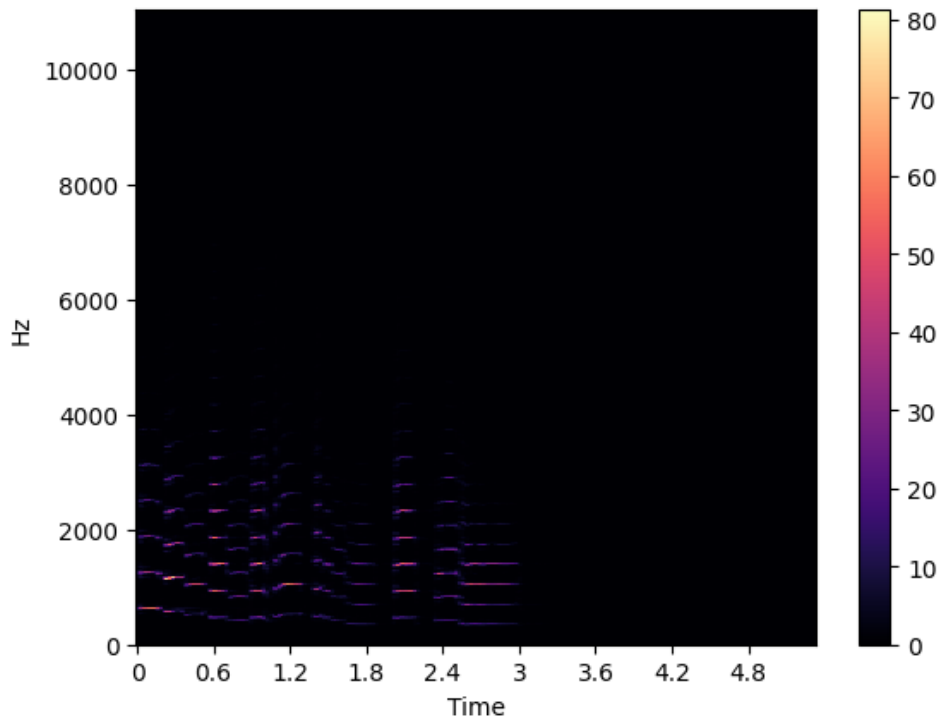
```
# Load the example audio in librosa
y, sr = librosa.load(librosa.example("trumpet"))

# 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`



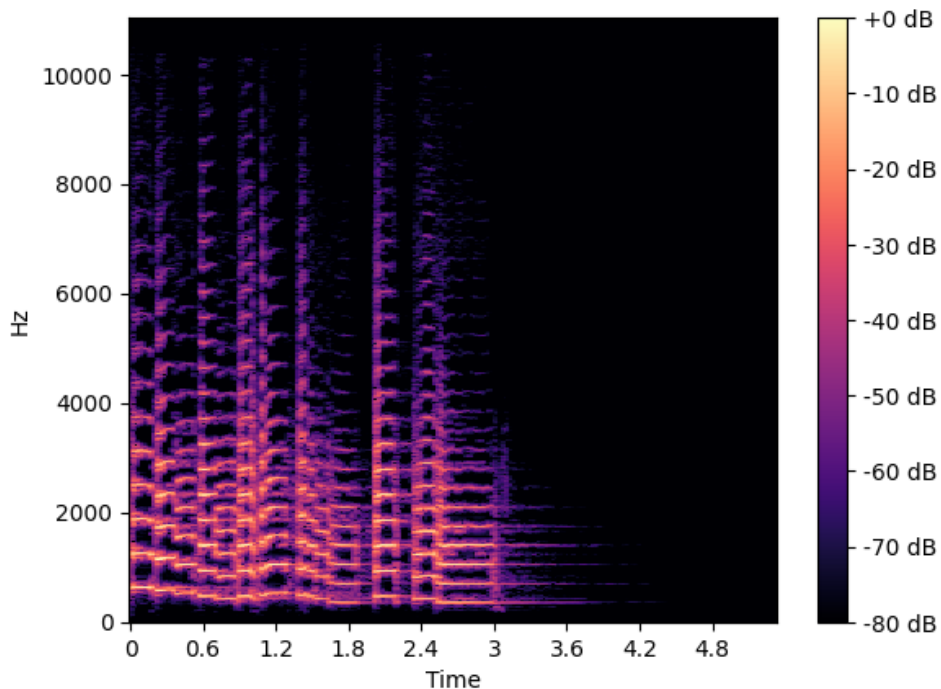
```
# Load the example audio in librosa
y, sr = librosa.load(librosa.example("trumpet"))

# 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"))

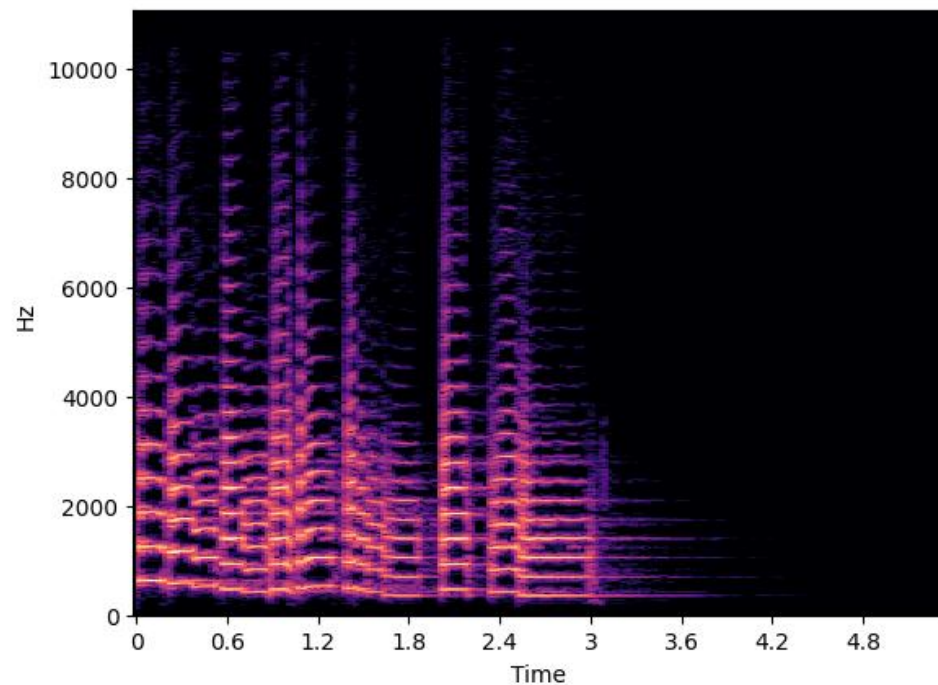
# Compute the spectrogram
S = np.abs(librosa.stft(y))
S_db = librosa.amplitude_to_db(S, ref=np.max)

# Plot the spectrogram
im = librosa.display.specshow(S_db, x_axis="time",
                               y_axis="linear")

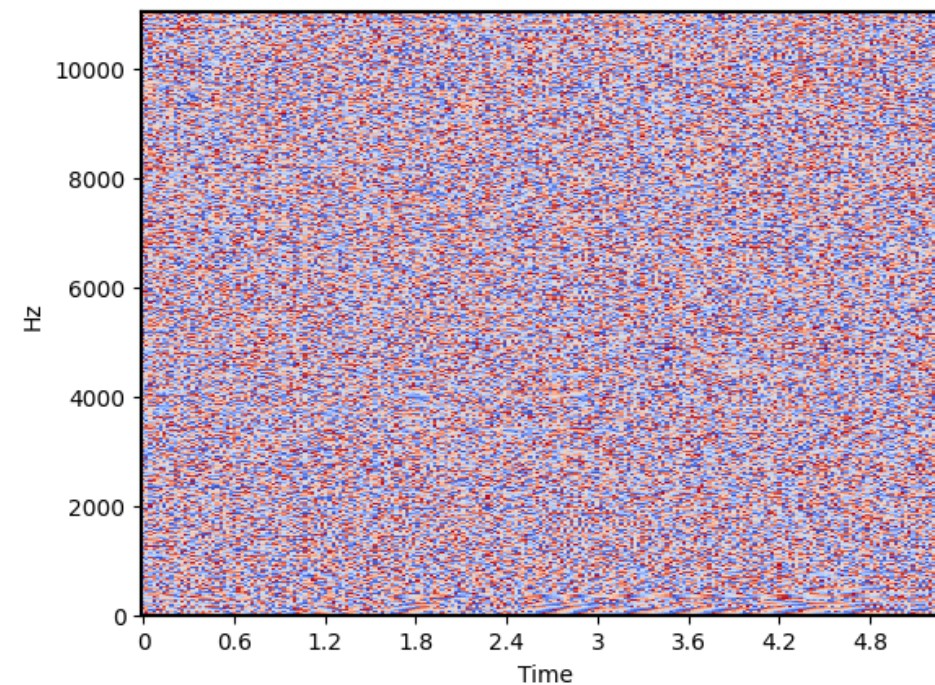
plt.colorbar(im, format="%+2.0f dB")
plt.show()
```

Example: Magnitude & Phase

Magnitude

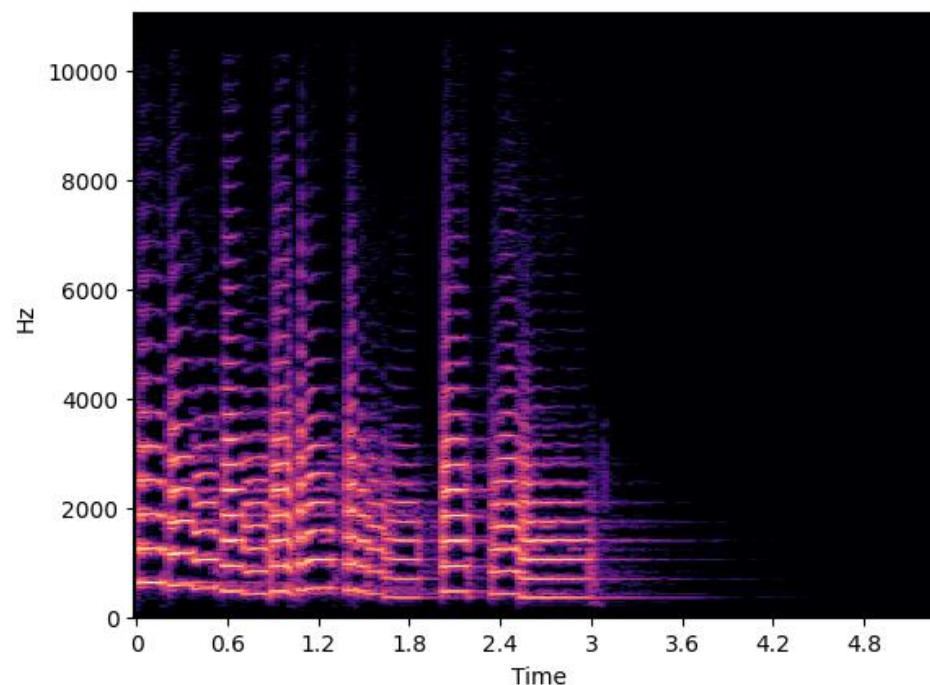
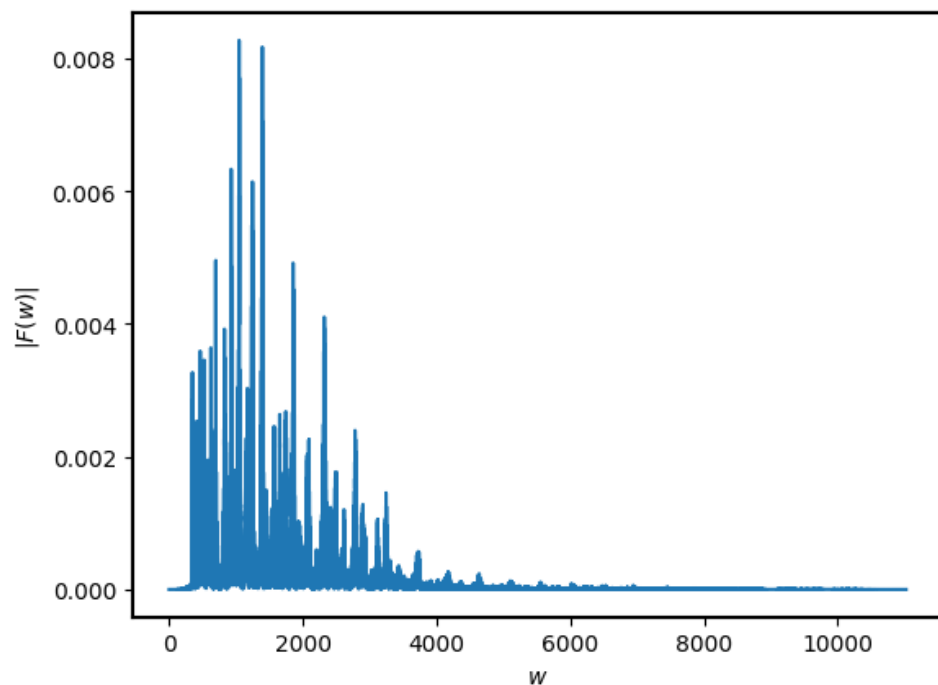


Phase



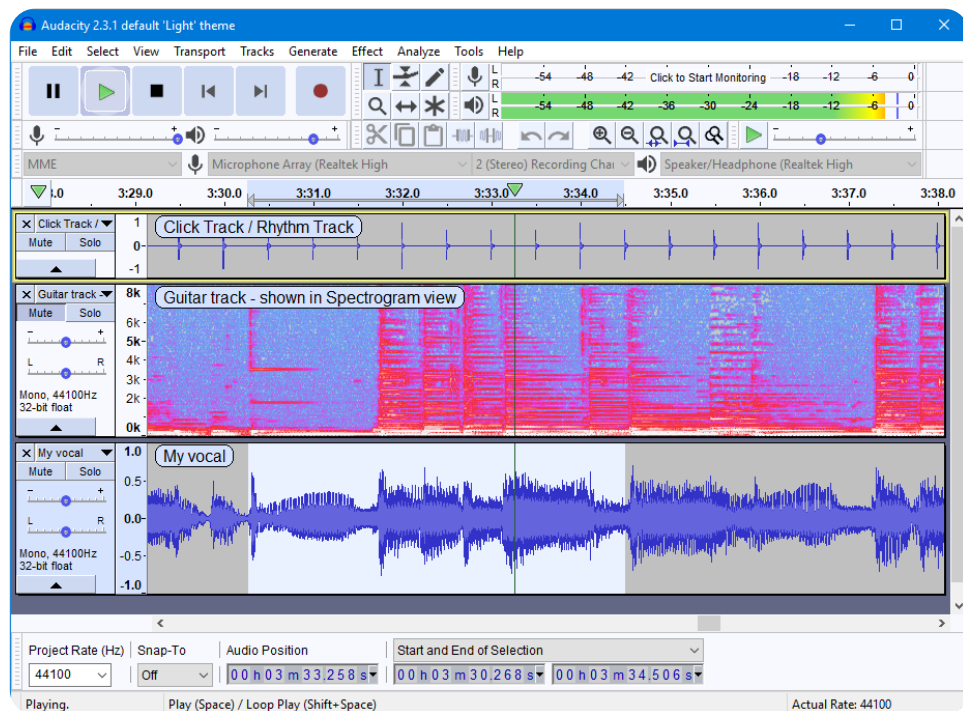
🔥 PA2: Spectral Analysis

- Use [librosa](https://librosa.org/doc/latest/index.html) to process audio files
 - Fast Fourier transform (**FFT**)
 - Short-time Fourier transform (**STFT**)



Software

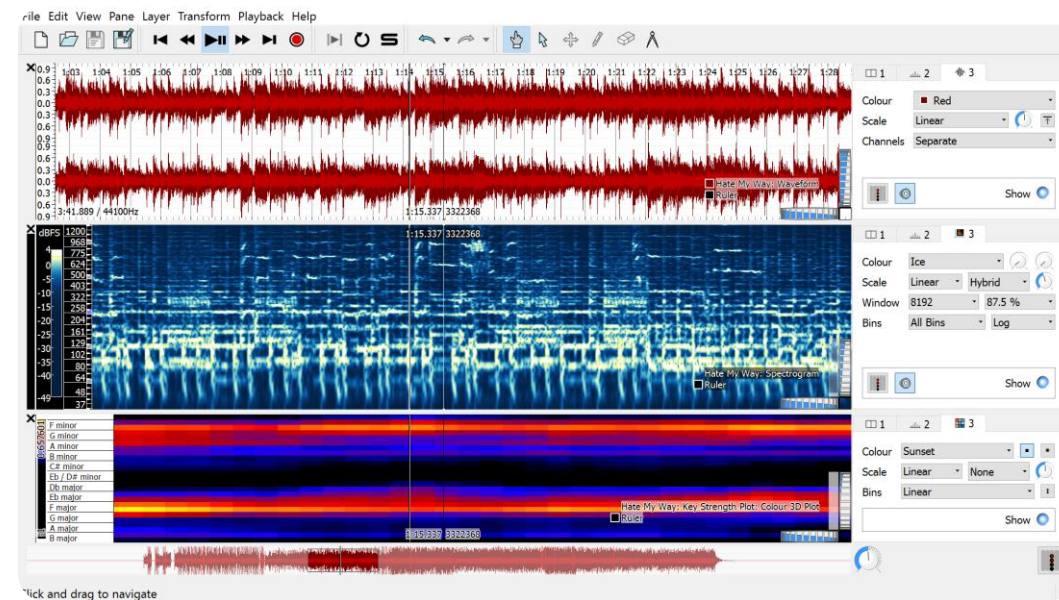
Audacity



(Source: audacity-2.3.1 via Internet Archive)

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

Sonic Visualiser

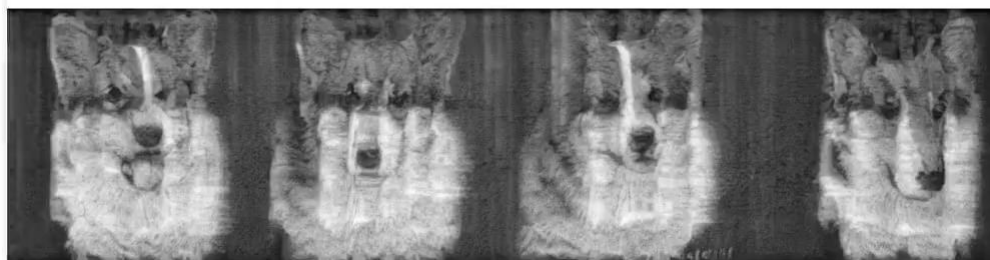
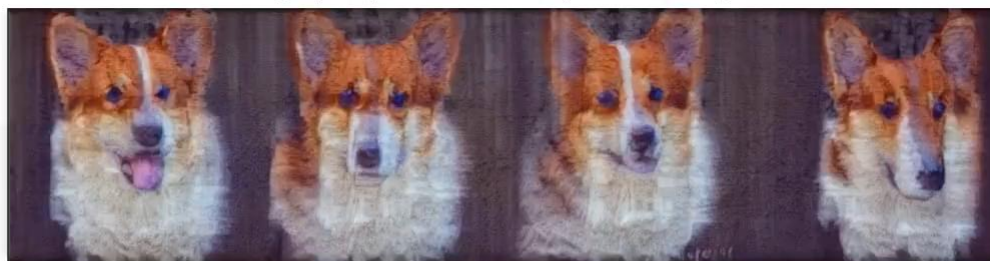


(Source: sonicvisualizer.org)

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



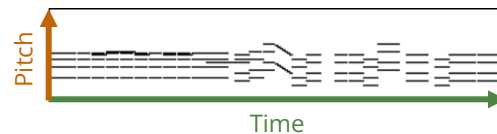
Symbolic music representations

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



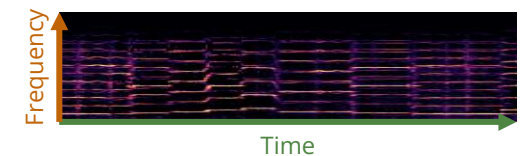
Audio-domain music representations

Time series-based



Waveform

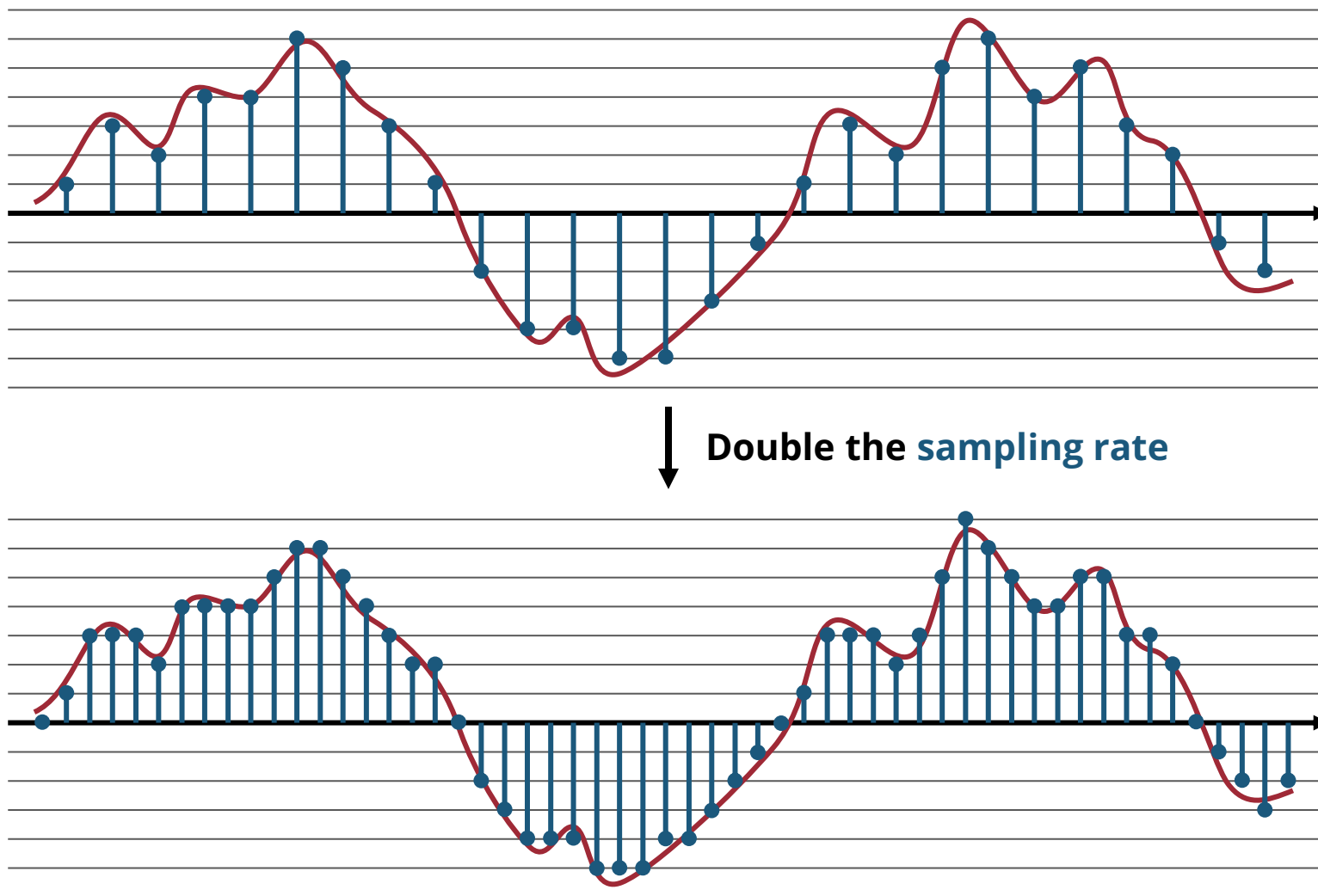
Image-based



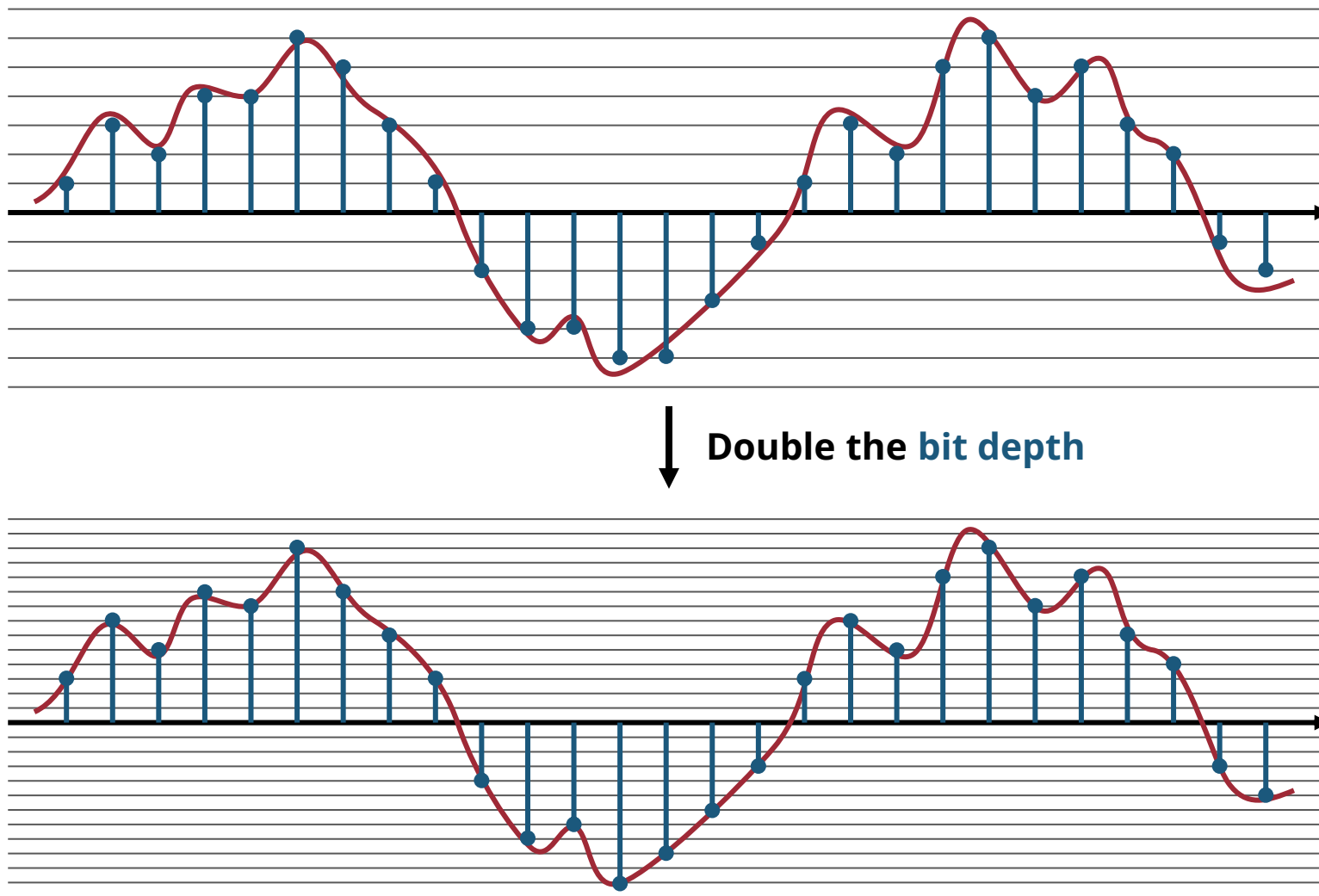
Spectrogram

Today's topic!

Resolution: Sampling Rate



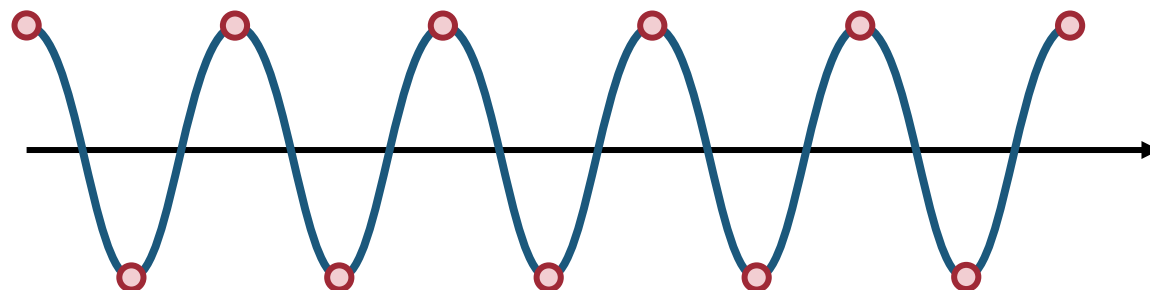
Resolution: Bit Depth



Sampling Theorem: Undersampling

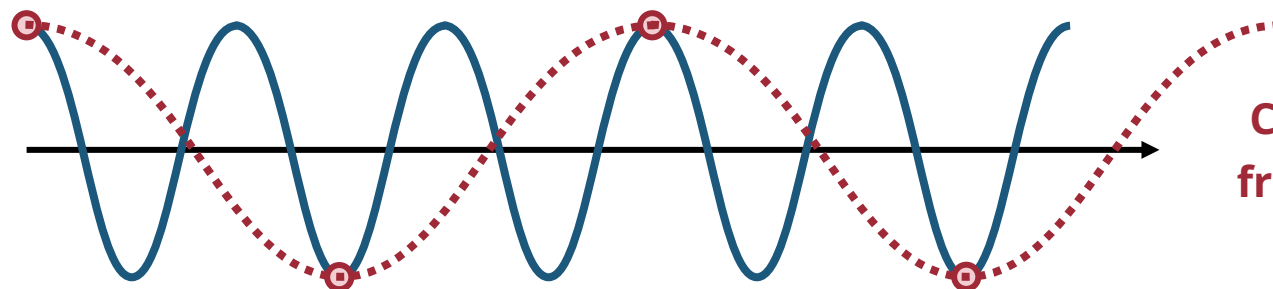
Critically sampled

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



Undersampled

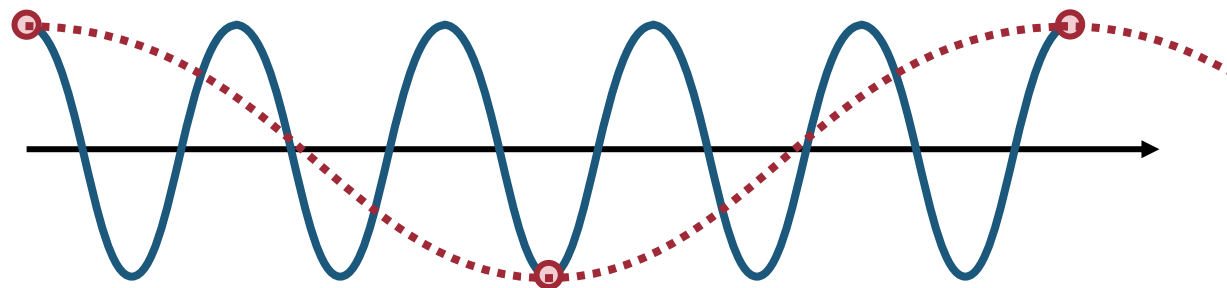
$$(f_s = \frac{2}{3}f_{max})$$



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

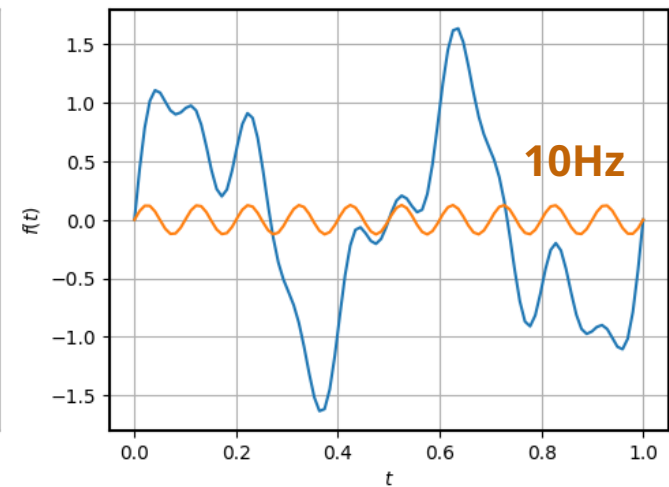
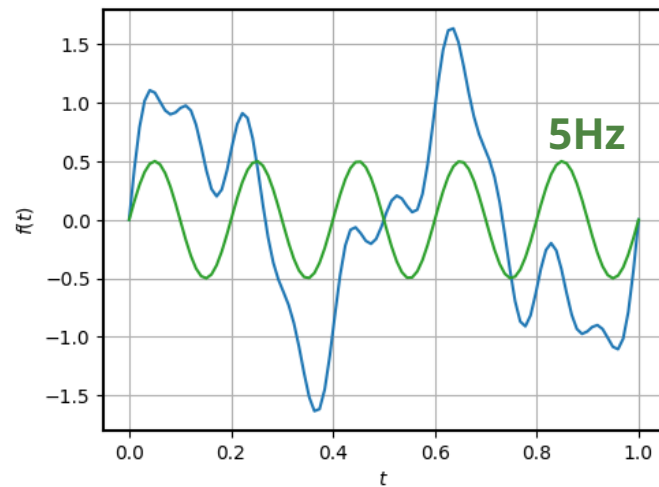
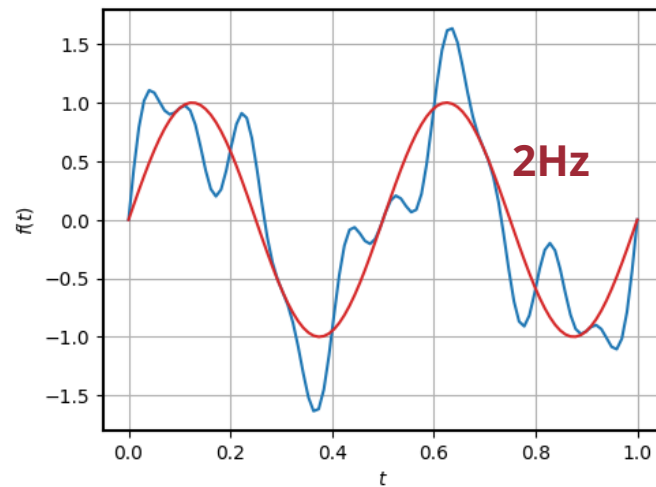
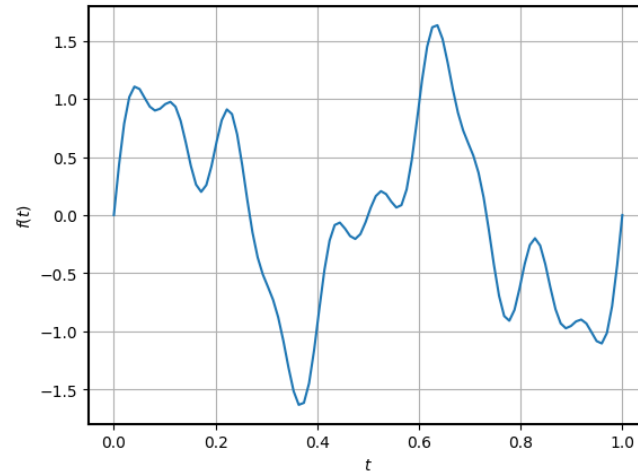
Undersampled

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



Can only reconstruct
frequency up to $\frac{1}{5}f_{max}$

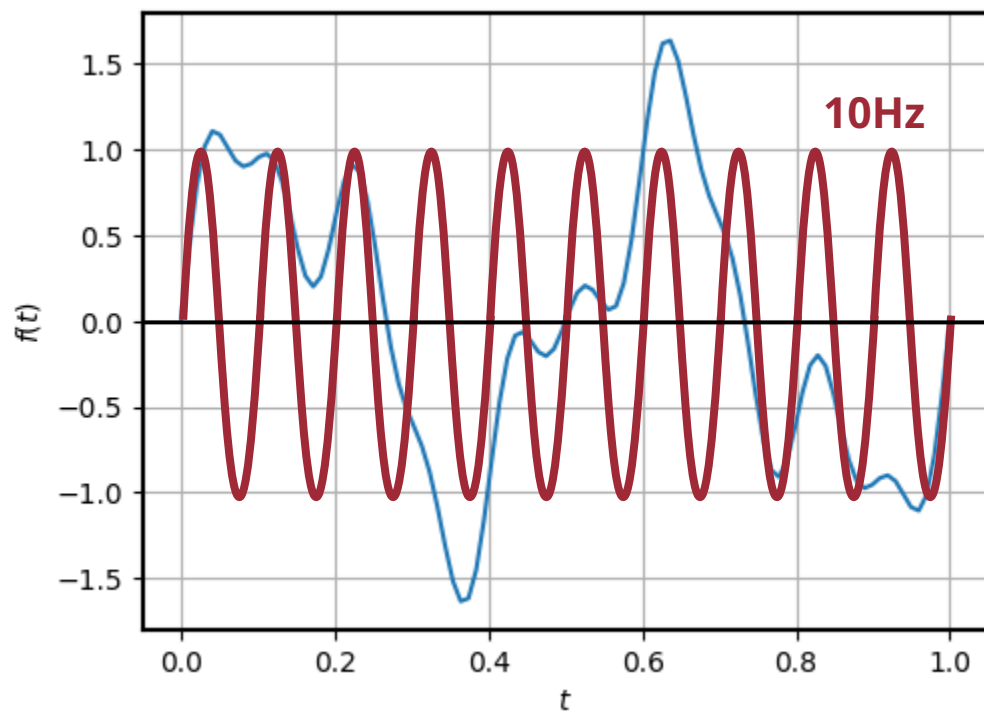
Spectral Analysis



Demystifying Fourier Transform

Signal

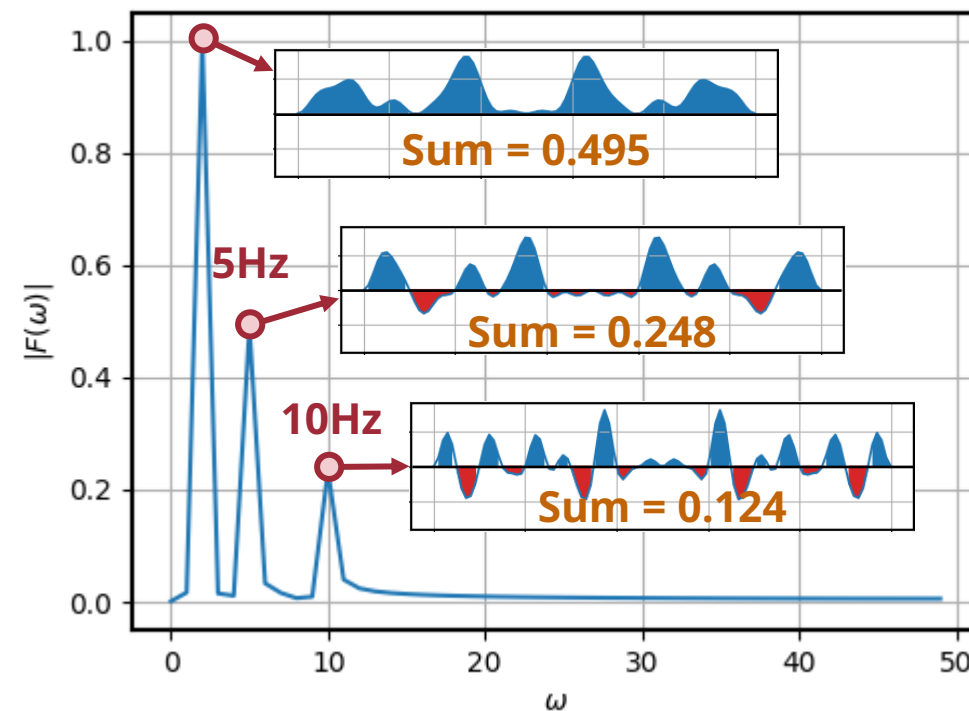
(time-domain)



Fourier Transform
→

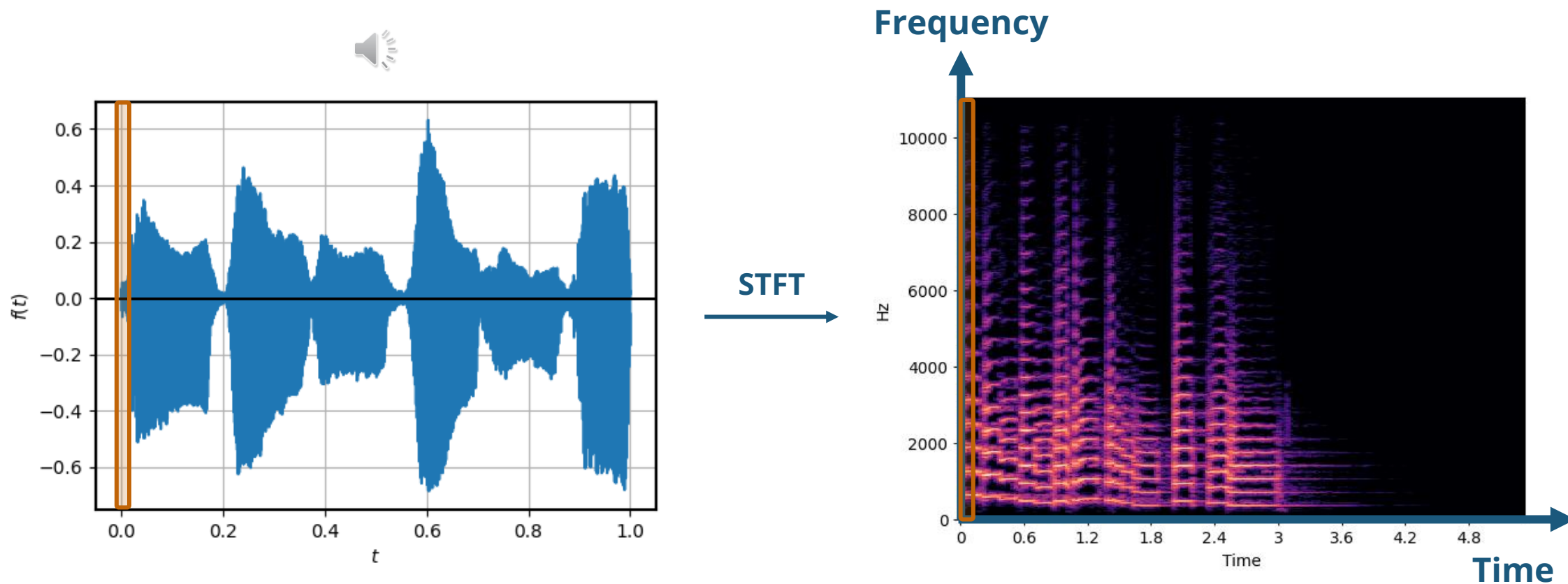
Spectrum

(frequency-domain)

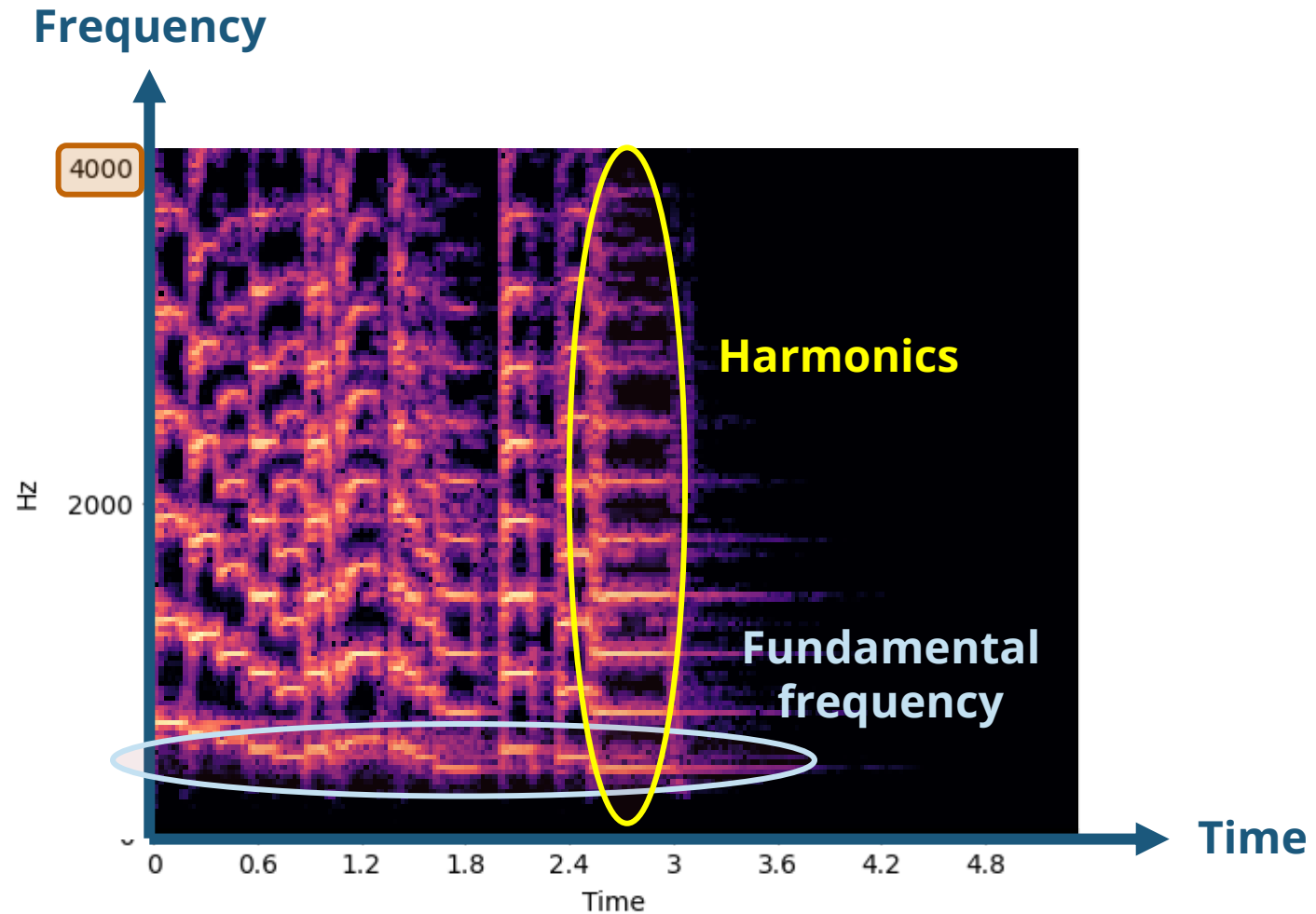


Short-Time Fourier Transform (STFT)

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



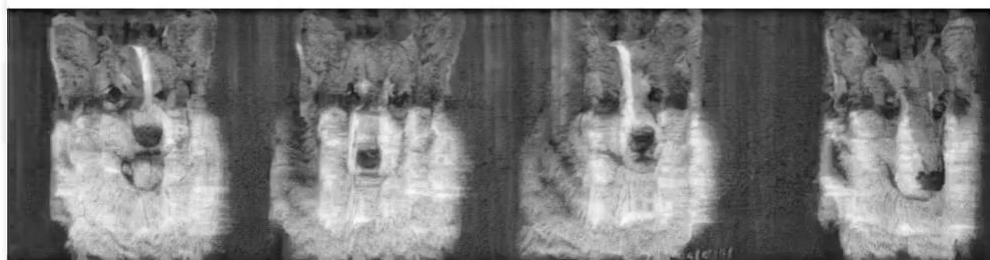
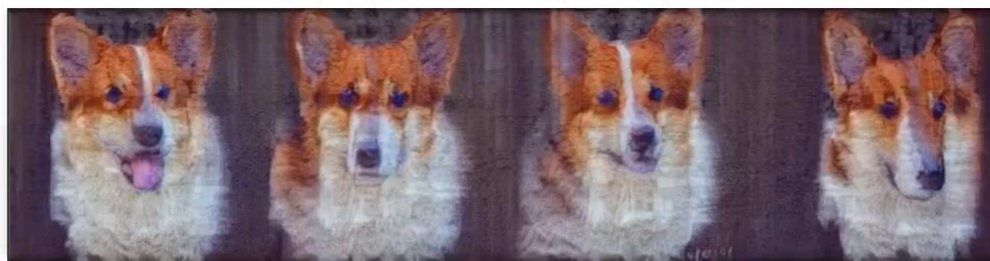
| Spectrogram



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

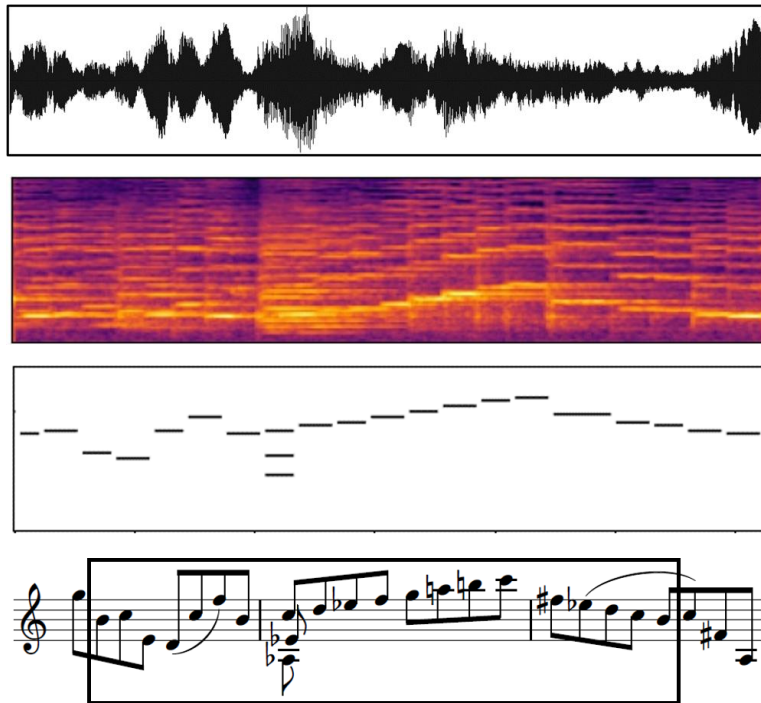


Audio prompt: tiger growling

(Source: Chen et al., 2024)

Next Lecture

Music Analysis



(Source: Dong et al., 2022)