# **String Data**

- In computing, much data is plain text, so you will be dealing with _strings_.  A "string" is a sequence of **_zero or more_** characters -- though usually 1+ _characters_

- You will use strings often, in different ways:
  - ➢ Printing as output
  - ➢ Fetching as input
  - ➢ Comparing
  - ➢ Reversing
  - ➢ Converting to/from other types

- Work and practice to become comfortable with this type and its many uses

# <u>Characters</u>

- The *ASCII character set* is older and smaller (**<u>8-bit</u>**) than Unicode, but is still quite popular (in C programs)

- The ASCII characters are a subset of the Unicode character set, including:

| | |
|---|---|
| **uppercase letters** | **A, B, C, …** |
| **lowercase letters** | **a, b, c, …** |
| **punctuation** | **period, semi-colon, …** |
| **digits** | **0, 1, 2, …** |
| **special symbols** | **&, \|, \, …** |
| **control characters** | **carriage return, tab, …** |

- See, for example, `http://www.asciitable.com`

# Characters

- Each character, however, will correspond to an *integer value* in some *character set*, and there are methods to perform conversions.

- The following example uses *Python* based methods
  - ➢ Integer to character: **chr**
  - ➢ Example: **chr(97)** → **a**
  - ➢ Character to integer: **ord**
  - ➢ Example: **ord('a')** → **97**

- This can be useful when you want to do arithmetic with characters, for example.

# Character Sets

- A *character set* is an ordered list of characters, with each character corresponding to a unique number

- Much software today uses the *Unicode character set*

- The Unicode character set uses sixteen bits per character, allowing for **65,536 (2^16)** unique characters

    o Unicode character values are often expressed as quartets of hex digits (4 hex digits equating to 16 bits), such as ☻ (*Char #920, or* **0398** *in hex*)

    o It is an international character set, containing symbols and characters from many world languages

    o Obviously, this is much more expansive than the ASCII character set!

- Reference: `https://unicode-table.com/en/`

# Character Encodings

- A *character <u>encoding</u>*, in contrast, deals with how characters (within a set) are to be represented in other, non-character forms: numerical, electrical, etc..

- For example, in computing, all data is encoded as bytes, which are made of **<u>bits</u>** : <u>ones</u> and <u>zeroes</u>

- A character encoding will entail

  o  a sequence of bits

  o  for each character

  o  within some character set

# Character Encodings

- An encoding for just **ASCII** characters would pretty simple, because char values are limited to *zero* through *127*, which requires only 7 bits (i.e., no more than 1 byte) per character

- For larger sets of characters, more bits would be needed.

- There are some other character encodings, more expansive than **ASCII** but still representable with 1 byte per character, such as **ISO-8859-1** (a.k.a., **Latin-1**)

- **Latin-1** is "ASCII-based" but includes a wider range of characters, such as accented vowels

# Character Encodings

- An encoding for just **ASCII** characters would pretty simple, because char values are limited to *zero* through *127*, which requires only 7 bits (i.e., no more than 1 byte) per character

- For larger sets of characters, more bits would be needed.

- There are some other character encodings, more expansive than **ASCII** but still representable with 1 byte per character, such as **ISO-8859-1** (a.k.a., **Latin-1**)

- **Latin-1** is "ASCII-based" but includes a wider range of characters, such as accented vowels

# Character Encodings

- An encoding for the Unicode character set would, in theory, entail two bytes (16 bits) per character.

- However, that could end up consuming space in memory, when frequently-used characters end up requiring the same space as rarely-used characters

- Also, if a character has a low numerical value, then many of its leading bits would be all zeroes (to make up the whole 16 bits)

- Fortunately, however, a character encoding -- the concrete representation of characters from a set -- can be designed *intelligently*

# Character Encodings

- The most popular (currently) encoding for the Unicode character set is UTF-8

- In UTF-8, a character is represented using from one to four bytes

- For any particular character, some of its leading bits will signal whether it is going to take up one, two, three, or four bytes

  - Characters *zero* through *127* are 1 byte: **0**_xxxxxxx_

  - Chars 128 through 2047 are 2 bytes: **110**_xxxxx_ **10**_xxxxxx_

- This allows for more efficient usage of space for storing characters as textual data

# Character Encodings

- The catch is that one must be mindful, to some extent, about which encoding is being used to...

  - _Write_ the text _to_ storage as bytes

  - _Read_ the bytes _from_ storage as text

  - (Here, consider "write" and "read" as roughly analogous to "save" and "open" -- in that they involve operations to and from disk)

- There are usually default encodings (within a program) for both writing and reading textual data

- These may be subject to user preference, to some extent

# Escape Sequences

- What if we want to include the quote character itself?
- The following line would confuse the interpreter because it would interpret the two pairs of quotes as two strings and the text between the strings as a syntax error:

```
print ("I said "Hello" to you.")
```

**A String**      **Syntax Error**      **A String**

- One option would be to replace the beginning and ending double-quote symbols with single-quotes:

```
print ('I said "Hello" to you.')
```

- The *reverse* would also be valid

```
print ("I said 'Hello' to you.")
```

# Escape Sequences

- Another option is to use *escape sequences*, which are character combinations that have a special meaning within a string

- Some Escape Sequences:

| Escape Sequence | Meaning |
|---|---|
| \t | tab |
| \n | newline |
| \r | carriage return |
| \" | double quote |
| \' | single quote |
| \\ | backslash |

Example:
```
print ("Hello,\n\tworld")
Hello,
    world
```

# Understanding and Working With Data

- No matter which route you take in the IT field, you will be dealing with data, in some form

- You can go with this definition for now: Data are pieces of information about the real world...

  o That are gathered and maintained -- as well as...

  o Made expressible and readable in some form(at) or another

  o For one or more purposes:

    ▪ Knowledge      ▪ Analysis          ▪ Decision-making

    ▪ Reporting      ▪ Interpretation    ▪ Problem-solving!

- There are many kinds of data....

# Numeric Data

- First, we have two type of **real numbers**:
  - *Integers* are whole numbers (No fractional component):
    **7**, **-358**, **0**, **-10**, **12398**
  - *Decimals* (or "floating-point") numbers do have a fractional component: **7.6**, **-35.8**, **-1.09**

- A **complex number** has an imaginary component.
  - In other words, some non-zero multiple of the constant *i*
  - We define *i* as the square root of **-1**
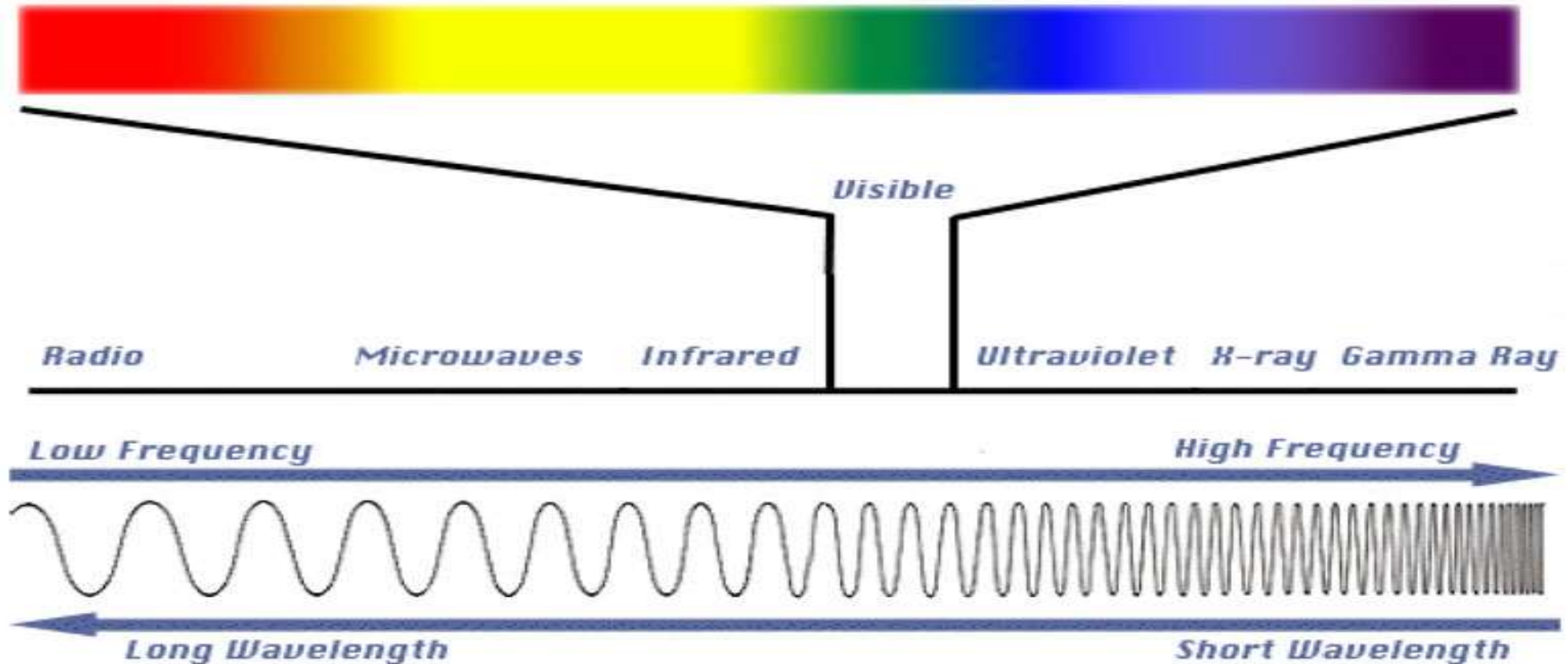  - We call *i* "imaginary" because no two real numbers can be squared to produce a negative result

# Boolean Data

- A `boolean` datum (*singular of "data"*) can have either of two possible values: **True** or **False**

- This is applicable to many ***either/or*** scenarios:
  - *Yes* or *No*
  - **1** or **0**
  - *Open* or *Closed*
  - *Up* or *Down*
  - *On* or *OffA*

# EMR, and The Science Behind It...

- The behavior of fiber optic cabling is based upon the transmission of ***electromagnetic radiation (EMR)***:
    - What is EMR?
    - What are some technologies that make use of it?
        - Radios
        - Microwave Ovens
        - X-Ray machines
- You will hear the term "light" used much more generally to refer to EMR – versus our more common definition of *visible* light.
- These varieties of EMR are all located on the ***electromagnetic spectrum***.
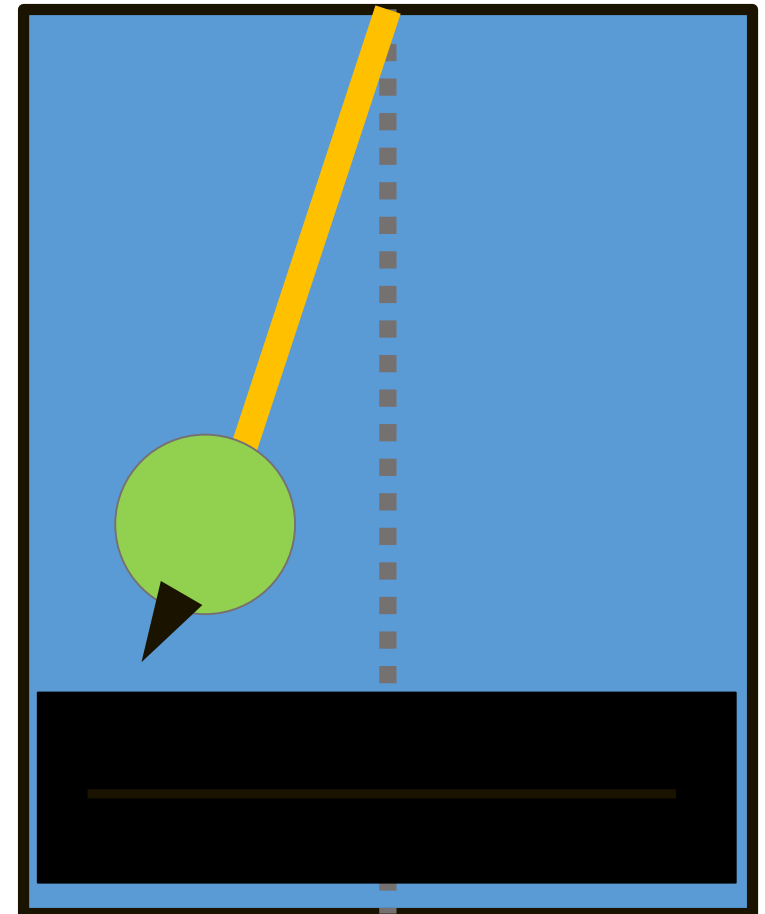
# The EM Spectrum



By Katarina Stevanovic (Own work) [CC BY-SA 3.0 (http://creativecommons.org/licenses/by-sa/3.0)], via Wikimedia Commons

**Video:** https://www.youtube.com/watch?v=cfXzwh3KadE
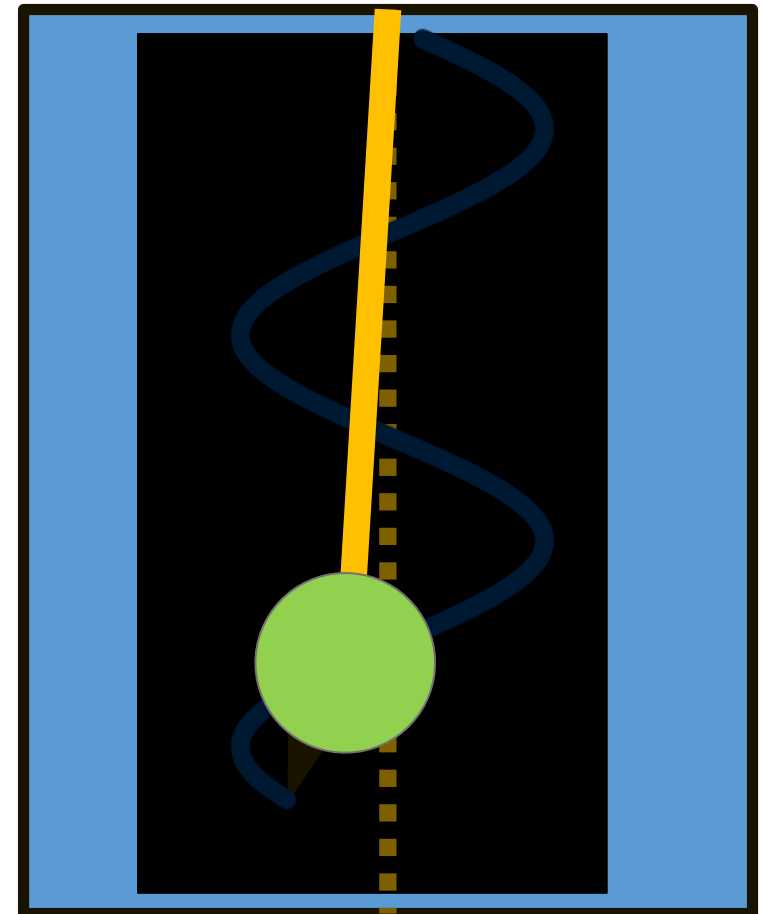
# The EM Spectrum

- EMR is generally conceptualized as __waves__ – cyclic variations about a "center", in which energy is transferred.  (In this case, the energy comes from various forms of EM activity.)

- To envision this, imagine a __pendulum__ swinging to and fro...

  o With a _pencil lead_ at the lowest point

  o In constant contact with _paper_

- If the paper is _stationary_, all the marks will remain within a single, one-dimensional space.
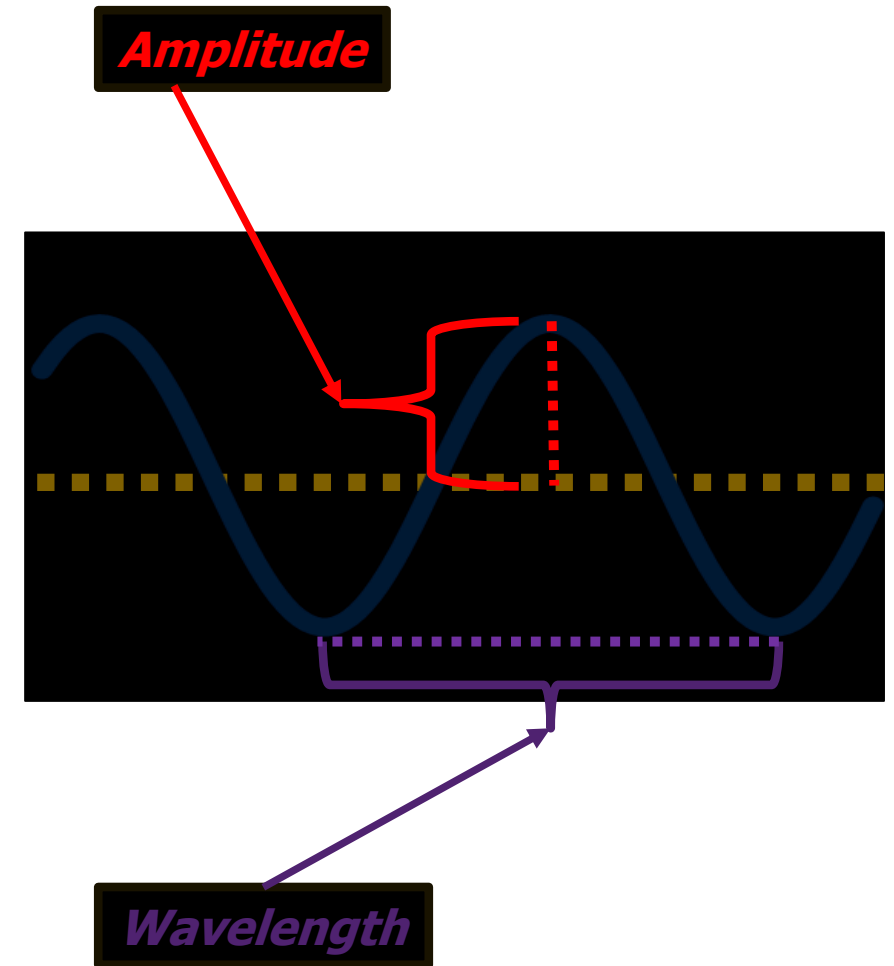
# The EM Spectrum

- However, if the paper – here being of indefinite length – is moving…

  - in one direction

  - at a constant speed

  - perpendicular to the pendulum's plane of motion

- …then you will see a graph of the pendulum's position, relative to the center, *over time*

- In other words, a ***wave***

# The EM Spectrum

- A wave – of EM radiation or another type – will feature a number of properties:

    o **Amplitude (A):** The height of a wave, from the center. **Unit:** *meters*.

    o **Wavelength (λ):** The distance between two analogous points on the wave graph. **Unit:** *meters*.

    o **Frequency (f):** The number of waves passing a given point during a given time. **Unit:** *Hertz* or *$s^{-1}$*

    o **Wave speed (v):** wavelength multiplied by frequency. **Unit:** *$m*s^{-1}$* or *m/s*

# The EM Spectrum

- We will encounter rather large and small numbers, that would normally have a lot of zeros.

- This would quickly become confusing, so we use a method called *scientific notation* to simplify this problem.

  - 300000000 = **$3 \times 10^8$**

  - **.**000000000000000000000602 = **$6.02 \times 10^{-23}$**

  - In computing context, these may be expressed as **3e8** and **6.02e-23**, respectively

  - We have *prefixes* for the different exponents of 10...

| MULTIPLICATION FACTOR | PREFIX | SYMBOL |
|---|---|---|
| $1\ 000\ 000\ 000\ 000\ 000\ 000 = 10^{18}$ | exa | E |
| $1\ 000\ 000\ 000\ 000\ 000 = 10^{15}$ | peta | P |
| $1\ 000\ 000\ 000\ 000 = 10^{12}$ | tera | T |
| $1\ 000\ 000\ 000 = 10^{9}$ | giga | G |
| $1\ 000\ 000 = 10^{6}$ | mega | M |
| $1\ 000 = 10^{3}$ | kilo | k |
| $100 = 10^{2}$ | hecto | h |
| $10 = 10^{1}$ | deka | da |
| $0.1 = 10^{-1}$ | deci | d |
| $0.01 = 10^{-2}$ | centi | c |
| $0.001 = 10^{-3}$ | milli | m |
| $0.000\ 001 = 10^{-6}$ | micro | m |
| $0.000\ 000\ 001 = 10^{-9}$ | nano | n |
| $0.000\ 000\ 000\ 001 = 10^{-12}$ | pico | p |
| $0.000\ 000\ 000\ 000\ 001 = 10^{-15}$ | femto | f |
| $0.000\ 000\ 000\ 000\ 000\ 001 = 10^{-18}$ | atto | a |

# The EM Spectrum

- In a vacuum, light (i.e., electromagnetic radiation) travels at a speed of <u>$3 \times 10^8$ m/s</u> (or **m*s$^{-1}$**).

- **Recall:** v = f * λ.  Therefore, if you know a signal's frequency, then you can calculate its <u>*wavelength*</u> by dividing the <u>*speed of light* (c)</u> by the <u>*frequency*</u>.  **<span style="color:red">λ = c / f</span>**

- Consider WBZ, a Boston radio station broadcasting at a frequency of 1030 kHz – or 1.03 Mhz:

$$\lambda = \frac{3.00 \times 10^8 \text{ m*s}^{-1}}{1.03 \times 10^6 \text{ s}^{-1}} \longrightarrow \lambda = \frac{3.00 \times 10^8 \text{ m*s}^{-1}}{1.03 \times 10^6 \text{ s}^{-1}}^{\;2}$$

# The EM Spectrum

$$\lambda = \frac{3.00 \times 10^8 \, m*s^{-1}}{1.03 \times 10^6 \, s^{-1}} \quad \rightarrow \quad \lambda = \frac{3.00 \times 10^2 \, m}{1.03} \quad \rightarrow \quad \lambda \approx 2.91 \times 10^2 \, m$$

- In doing this, it is important to keep track of your units, including which ones combine or cancel out in the arithmetic.

- Understanding these concepts will be helpful not only for this chapter but also for subsequent chapters, such as wireless networking, which also uses EMR

# The EM Spectrum