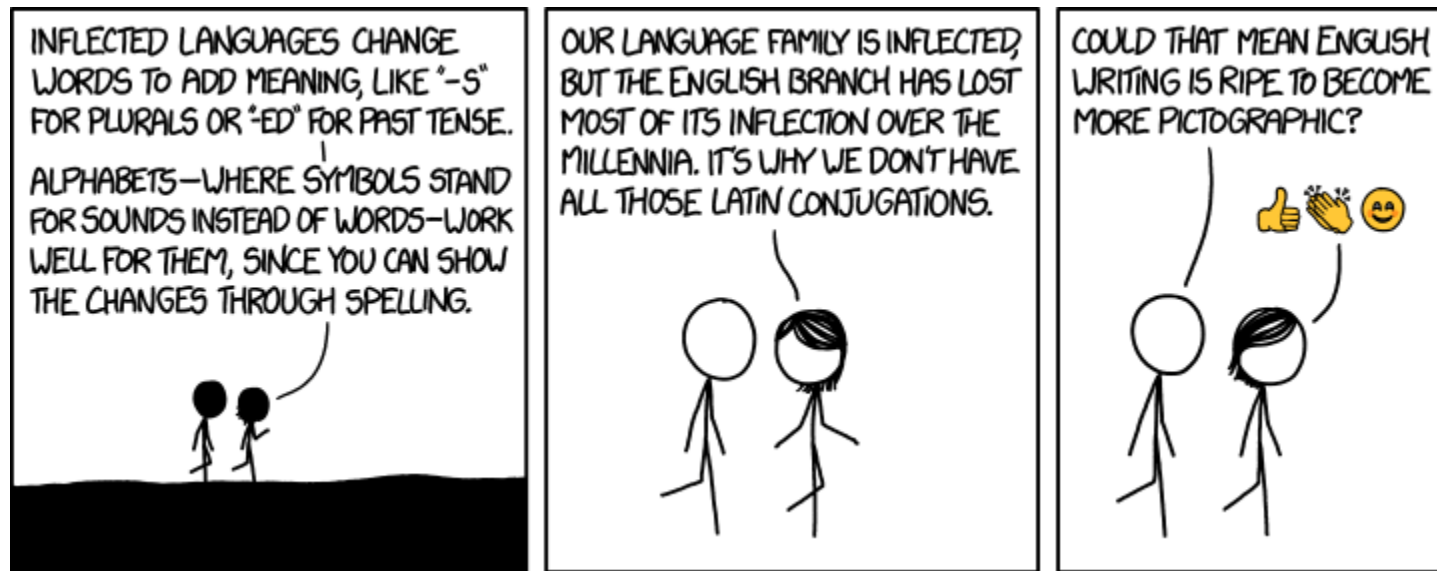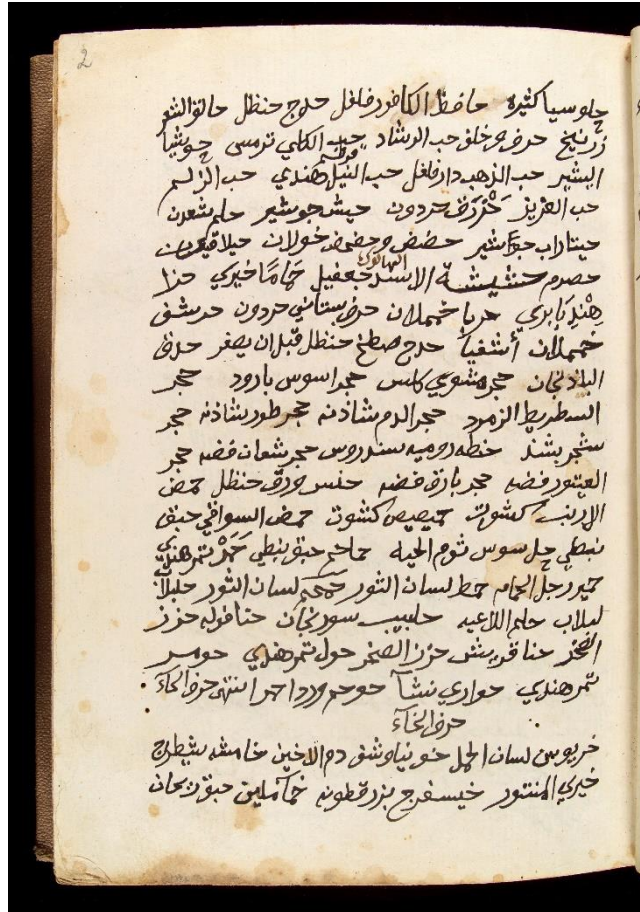# ENGSCI 233 Lecture 2

## Structured Data

# Today's objectives:

- Understand role of endianness in data storage
- Understand how text data are represented in computers
- Know how to handle situations where other forms of data need to be stored
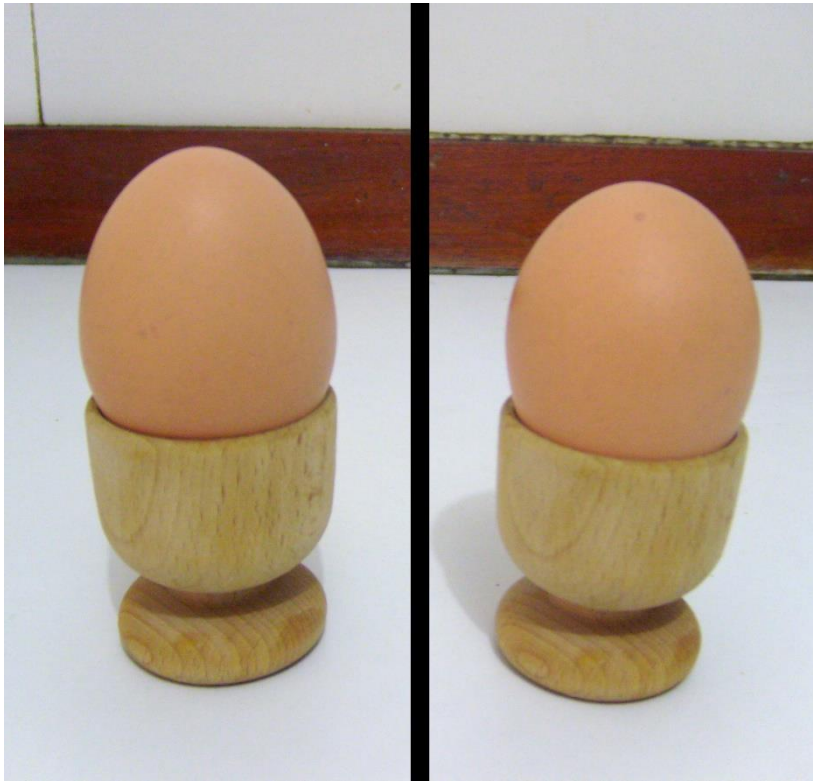
# What happens when we have more than 8 bits?

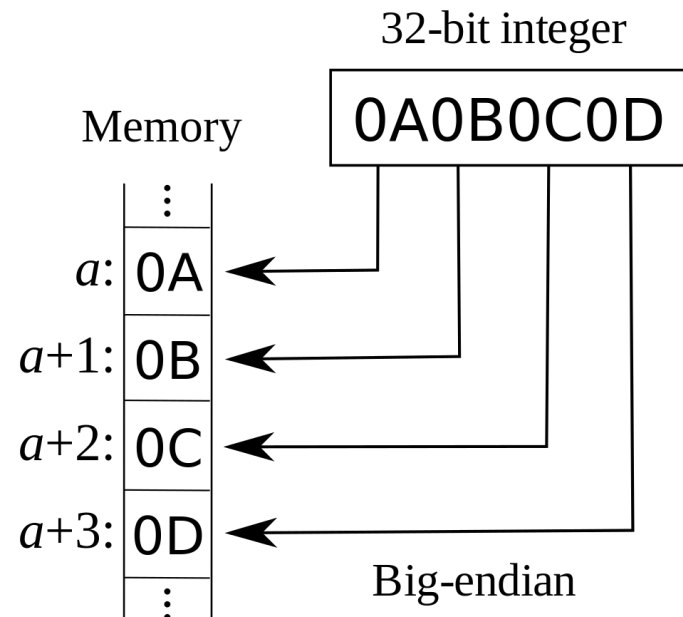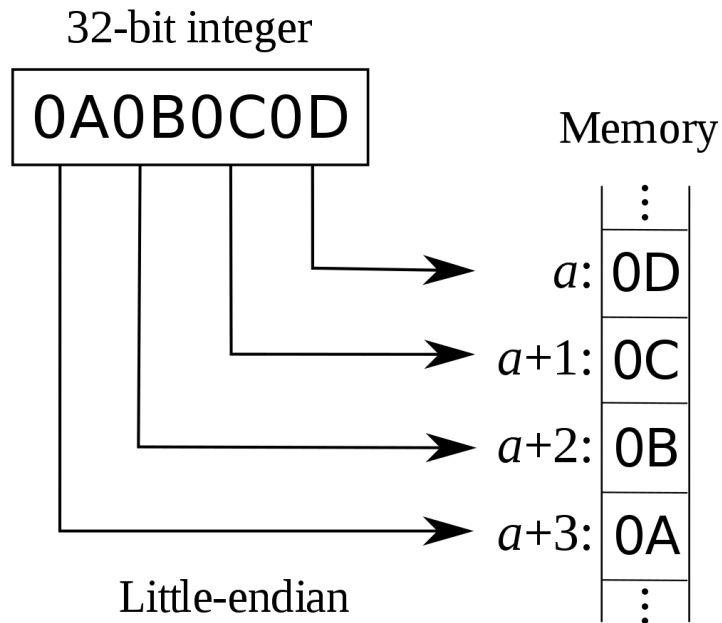# Which way does a computer read?



- We have been reading left-to-right

- Not everyone does!

- Computers are no different…

# Which way does a computer read?



- Big-endian
  - First byte is highest
  - Used in mainframes and network equipment

- Little-endian
  - First byte is lowest
  - Used on most PCs

# What does this look like?

32-bit integer

| 0A0B0C0D |

Memory

| | |
|---|---|
| | ⋮ |
| $a$: | 0D |
| $a+1$: | 0C |
| $a+2$: | 0B |
| $a+3$: | 0A |
| | ⋮ |

Little-endian

32-bit integer

| 0A0B0C0D |

Memory

| | |
|---|---|
| | ⋮ |
| $a$: | 0A |
| $a+1$: | 0B |
| $a+2$: | 0C |
| $a+3$: | 0D |
| | ⋮ |

Big-endian

# What does this look like?

- 0xFACEFEED

# Why does endianness matter?

- Your code normally can't see it!
  - Programming languages do a good job of hiding it.

- What happens when you store data, though?

# What about non-numbers?

- Encode each character as a number.
- But how?
- English fits in seven bits:

## ASCII Code Chart

|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | A | B | C | D | E | F |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | NUL | SOH | STX | ETX | EOT | ENQ | ACK | BEL | BS | HT | LF | VT | FF | CR | SO | SI |
| 1 | DLE | DC1 | DC2 | DC3 | DC4 | NAK | SYN | ETB | CAN | EM | SUB | ESC | FS | GS | RS | US |
| 2 |   | ! | " | # | $ | % | & | ' | ( | ) | * | + | , | - | . | / |
| 3 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | : | ; | < | = | > | ? |
| 4 | @ | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O |
| 5 | P | Q | R | S | T | U | V | W | X | Y | Z | [ | \ | ] | ^ | _ |
| 6 | ` | a | b | c | d | e | f | g | h | i | j | k | l | m | n | o |
| 7 | p | q | r | s | t | u | v | w | x | y | z | { | \| | } | ~ | DEL |

# How is ASCII used?

## ASCII Code Chart

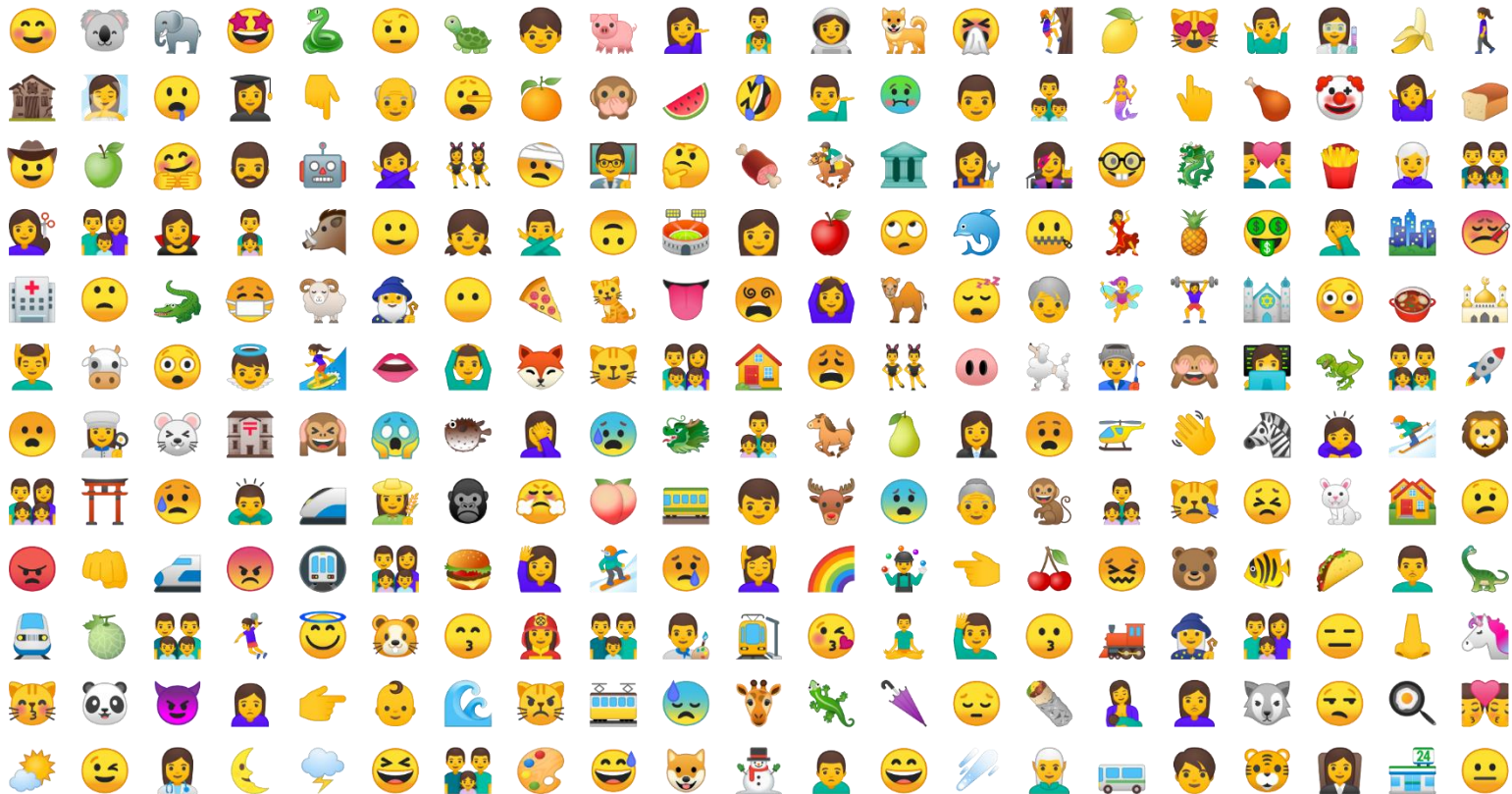|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | A | B | C | D | E | F |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | NUL | SOH | STX | ETX | EOT | ENQ | ACK | BEL | BS | HT | LF | VT | FF | CR | SO | SI |
| 1 | DLE | DC1 | DC2 | DC3 | DC4 | NAK | SYN | ETB | CAN | EM | SUB | ESC | FS | GS | RS | US |
| 2 |   | ! | " | # | $ | % | & | ' | ( | ) | * | + | , | - | . | / |
| 3 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | : | ; | < | = | > | ? |
| 4 | @ | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O |
| 5 | P | Q | R | S | T | U | V | W | X | Y | Z | [ | \ | ] | ^ | _ |
| 6 | ` | a | b | c | d | e | f | g | h | i | j | k | l | m | n | o |
| 7 | p | q | r | s | t | u | v | w | x | y | z | { | | | } | ~ | DEL |

Hello, world!

# The world no longer uses ASCII.

# The world no longer uses ASCII.

# How can we represent every language?

- One standard with 137,439 characters
- Room for over 1 million!
- We can't fit every character in one byte any more…

# How do we store this?

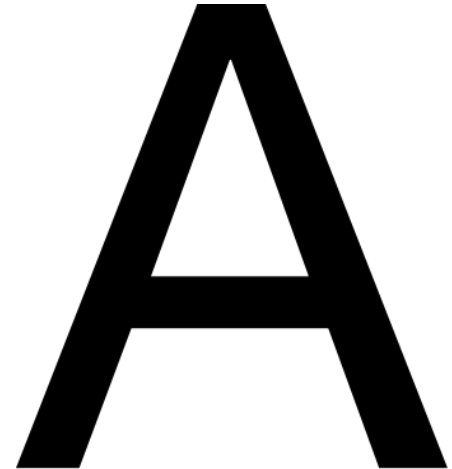- The simple way: UTF-32
- 4 bytes per character



U+1F4A9

0x0001F4A9

U+2615

0x00002615

U+41

0x00000041

# How do we store this?

- The Windows way: UTF-16
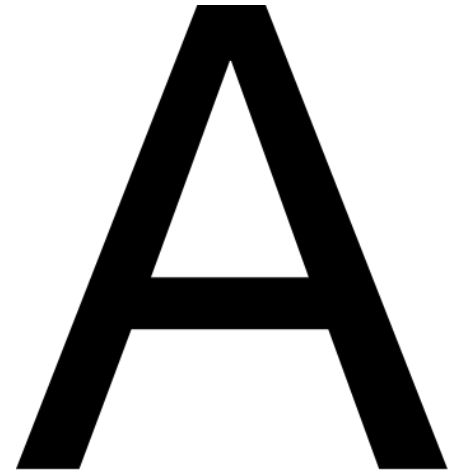- 2 bytes per character, if possible

U+1F4A9

0xD83D DCA9

U+2615

0x2615

U+41

0x0041

# How do we store this?

- The standard way: UTF-8
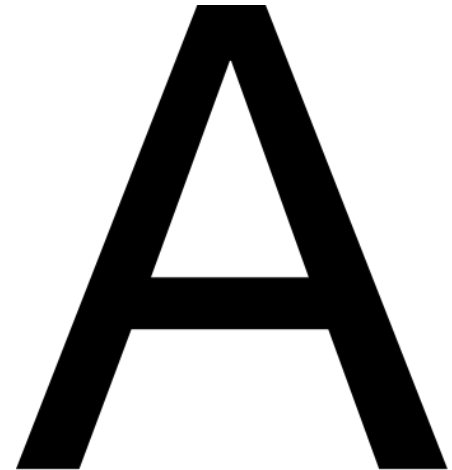- As few bytes as unambiguously possible
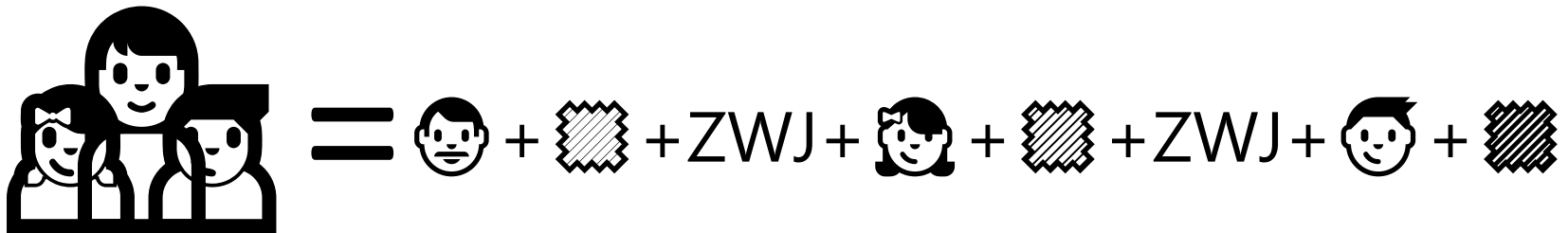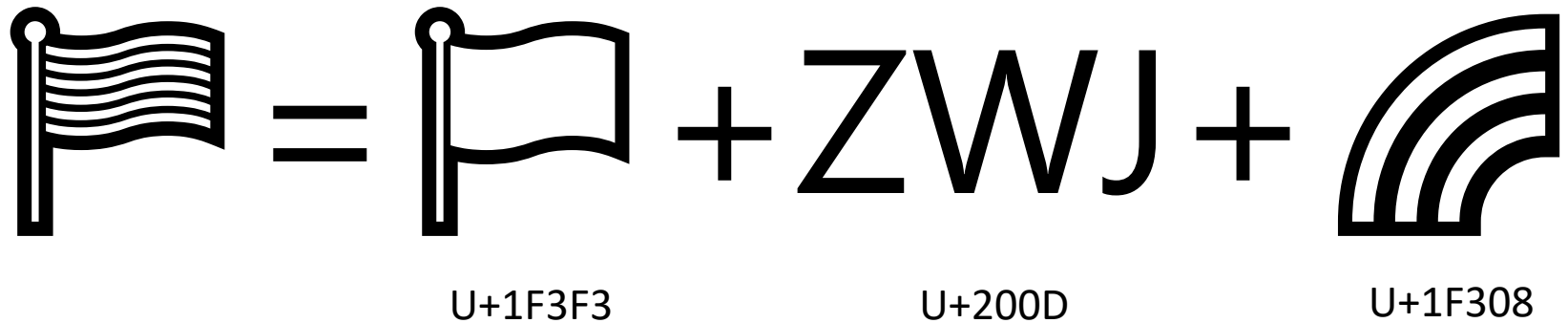
U+1F4A9

0xF0 9F 92 A9

U+2615

0xE2 98 95

U+41

0x41

# Some entities use more characters.

🏳️ = 🏳 + ZWJ + 🌈

U+1F3F3    U+200D    U+1F308

👪 = 👨 + ⬚ + ZWJ + 👧 + ⬚ + ZWJ + 👦 + ⬚

# What does this mean for you?

- ASCII strings are always 1 byte per character
- UTF-8 and UTF-16 have a variable character size
  - How do you figure out the length of a string in memory?
- Python cheats and uses the smallest constant-size representation that will work…
- You can't assume that displayed characters and numbers are interchangeable any more!
- You can't predict required display space from string length.

# What else is there to store?

- Images?

- Formatted text?

- Executable code?

# Use a standard!

- Images?

- Formatted text?

- Executable code?

INTERNATIONAL
STANDARD

**ISO/IEC
15948:2004**

Information technology -- Computer
graphics and image processing -- Portable
Network Graphics (PNG): Functional
specification

# Use a standard!

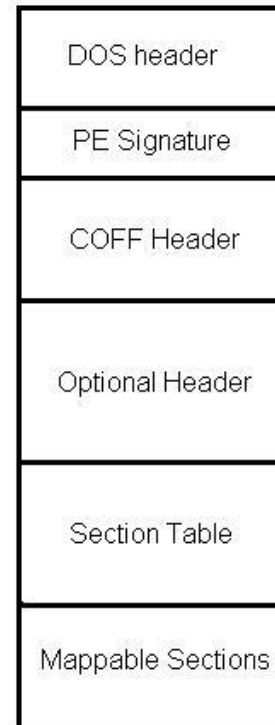- Images?

- Formatted text?

- Executable code?

# Use a standard!

- Images?

- Formatted text?

- Executable code?

| |
|---|
| DOS header |
| PE Signature |
| COFF Header |
| Optional Header |
| Section Table |
| Mappable Sections |

# How do standards work?

- Often contain a magic number or signature
  - PNG image has ASCII "PNG" plus bytes to detect mis-processing
  - PDF starts with ASCII "%PDF-"
- Fixed size fields or variable size stored with length
- Metadata stored with data
- Documentation made available to all
  - Not always freely!

# What about pointers?

- Don't keep pointers around for posterity!
- A pointer has the size in memory dictated by the computer architecture.

8-bit

32-bit

16-bit

64-bit

# Next time: Computer Hardware

# Image References

Slide 1: *Inflection*, by Randall Munroe, from https://xkcd.com/1709/ (CC BY-NC 2.5)

Slides 3 and 19: by Christiaan Colen, from https://www.flickr.com/photos/christiaancolen/20607150556 (CC BY-SA 2.0)

Slides 4 and 11 (left): from the Wellcome Collection, https://wellcomecollection.org/works/vzjsdw4c?query=L0033679 (CC BY 4.0)

Slide 5: by Cirofono, from https://www.flickr.com/photos/ciroduran/2060860569 (CC BY 2.0)

Slide 6 (left): by R. S. Shaw, from https://commons.wikimedia.org/wiki/File:Little-Endian.svg (Public domain)

Slide 6 (right): by R. S. Shaw, from https://commons.wikimedia.org/wiki/File:Big-Endian.svg (Public domain)

Slide 8: by Evan-Amos, from https://commons.wikimedia.org/wiki/File:SanDisk-Cruzer-USB-4GB-ThumbDrive.jpg (Public domain)

Slides 9 and 10: by Anomie, from https://commons.wikimedia.org/wiki/File:ASCII_Code_Chart.svg (Not subject to copyright)

Slide 11 (right): by John Vetterli, from https://www.flickr.com/photos/51824383@N00/935697405 (CC BY-SA 2.0)

Slide 12: from https://techcrunch.com/2017/08/21/meet-android-oreos-all-new-emoji/

Slides 14-16 (left): by Google, from https://commons.wikimedia.org/wiki/File:Emoji_u1f4a9.svg (Apache license)

Slides 14-16 (centre): by Google, from https://commons.wikimedia.org/wiki/File:Emoji_u2615.svg (Apache license)

Slides 14-16 (right): by Eirik1231, from https://commons.wikimedia.org/wiki/File:Latin_alphabet_Aa.svg (Public domain)

Slide 22: by Whiteknight, from https://commons.wikimedia.org/wiki/File:RevEngPEFile.JPG (Public domain)

Slide 24 (top left): by Evan-Amos, from https://commons.wikimedia.org/wiki/File:NES-Console-Set.jpg (Public domain)

Slide 24 (top right): by Evan-Amos, from https://commons.wikimedia.org/wiki/File:PSX-Console-wController.jpg (Public domain)

Slide 24 (bottom left): by Evan-Amos, from https://commons.wikimedia.org/wiki/File:SNES-Mod1-Console-Set.jpg (CC BY-SA 3.0)

Slide 24 (bottom right): by Evan-Amos, from https://commons.wikimedia.org/wiki/File:Nintendo-64-wController-L.jpg (Public domain)