

Data representation & storage

GIGA Doctorate School

Christophe Phillips, Ir Ph.D.



Program

- Bits & bytes
- Data format
- Signal discretization
- File format & compression
- Storage & Safety



Program

- Bits & bytes
- Data format
- Signal discretization
- File format & compression
- Storage & Safety



Bits & bytes

- Bit (for "binary digit") =
 - a basic unit of information used in computing and digital communications.
 - can have only one of two values → physically represented with a twostate device.
 - most commonly represented as either a 0 or 1
- Byte =
 - a unit of digital information
 - most commonly consists of eight bits,
 - representing a binary number



Bytes

Originally,

- number of bits used to encode a single character of text in a computer
- hardware dependent
- convenient as power of $2 \rightarrow$ values from 0 to 255

Now

- *de facto* standard for smallest amount of "memory unit"
- 32- or 64-bit 'words', built of four or eight bytes
- aka. "octet", symbol 'o',



Memory size

Expressed in binary vs. decimal base

| Name | Binary | Decimal | Discrepancy |
|----------------|-----------------------------------|-----------------------|-------------|
| Kilo-byte (kb) | 2^10 = 1.024 o | 1.000 | 2,4% |
| Mega-byte (Mb) | 2^20 = 1.048.576 o | 1.000.000 | 4,8% |
| Giga-byte (Gb) | 2^30 = 1.073.741.824 o | 1.000.000.000 | 7,4% |
| Tera-byte (Tb) | 2^40 = 1.099.511.627.776 o | 1.000.000.000.000 | 9,9% |
| Peta-byte (Pb) | 2^50 = 1.125.899.906.842.674 o | 1.000.000.000.000.000 | 12,6% |



Transfer speed

- Typical bandwidth
- ► RAM, ~10Gb per second → 1Gb of data in ~ 0,1 second
- ▶ Hard drive, ~0,5Gb per second
 → 10Gb of data in ~ 20 second
- Network, ~100Mbps = ~ 0,1GB per second → 1Tb of data in ~ 10.000 seconds = ~2.8 hours !!!

Data transfer can be a bottle neck!



Program

- Bits & bytes
- Data format
- Signal discretization
- File format & compression
- Storage & Safety

USASCII code chart

- = letter, digit, or punctuation
- With 1 byte, 1 simple character, aka. 'char' from ASCII (American Standard Code for Information Interchange) to UTF-8 (Unicode Transformation Format – 8-bit)
- UTF-8 extended up to 4 bytes
 - \rightarrow extension to more characters and alphabets
 - \rightarrow mot common for WWW and emails encoding

| Deb | 5 - | | | •••• | | °°, | °°, | ° 0 | °, | ¹ 0 ₀ | '°, | ¹ 1 | 1 1 1 |
|-----|--------|----------|----------|---------|--------------|-------------|------|-------|----|-----------------------------|-----|----------------|-------------|
| | 6 1 | Þ 3 1 | Þ 2 1 | ۵, ۱ | Row | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| | 0 | 0 | 0 | 0 | 0 | NUL | DLE | 5P | 0 | 0 | P | ``` | P |
| | Ò | 0 | 0 | 1 | 1 | SOH | DC1 | [! | 1 | Α | 0 | ٥ | q |
| | 0 | 0 | Ţ | 0 | 2 | STX | DCZ | • | 2 | 6 | Ĥ | 9 | r |
| | Ö | 0 | | - | 3 | ETX | DC 3 | # | 3 | C | S | c | 5 |
| | 0 | | 0 | 0 | 4 | EOT | DC4 | 1 | 4 | D | Т | d | t |
| | 0 | [ï | 0 | + | 5 | ENQ | NAK | % | 5 | E | υ | ŧ | 2 |
| | 0 | Γι_ | 1 | 0 | 6 | ACK | SYN | 8 | 6 | F | V | 1 | ۷ |
| | 0 | Ĩ | 1 | L | 7 | 8EL | ETĐ | • | 7 | G | * | g | w |
| | 1 | 0 | 0 | 0 | 8 | 85 | CAN | (| 8 | н | × | h | × |
| | | 0 | 0 | 1 | 9 | нт | EM |) | 9 |] | Ŷ | i | y j |
| to | 1 | 0 | Т | 0 | 10 | LF | \$U8 | | : | J | Z | ز | ž |
| ιυ | Т | 0 | I | 1 | | VT | ESC | + | ì | к | C | k | |
| | Ŧ | I I | ٥[| 0 | 12 | FF | FS | | < | L | N | 1 | |
| | I | T | 0 | | <u>[</u> 13] | CR | GS | - | E | м | 3 | m | } |
| | , | T | I | 0 | 14 | 50 | RS | | > | N | . ^ | n | \sim |
| | 1 | | | 1 T | 15 | \$ 1 | US | 1 | ? | 0 | - | 0 | DEL |



Integer, signed or unsigned

- with 1 byte,
 - 'int8', values ∈ [-128 127]
 - 'unit8', values ∈ [0 255]
- with 2 bytes,
 - 'int16' or 'short', values \in [-32,768 32,767] i.e. [-(2¹⁵) 2¹⁵ 1]
 - 'uint16', values \in [0 65,535] i.e. [0 2^{16} 1]
- with 4 bytes,
 - 'int32' or 'long', values $\in [-(2^{31}) \ 2^{31} 1]$
 - 'uint32', values \in [0 4,294,967,295 i.e. [0 $2^{32} 1$]
- with 8 bytes,

...



Floating-point

- Single-precision = 32 bits = 4 bytes
- wide dynamic range of values with "floating radix point":
 - sign bit : 1 bit
 - exponent width: 8 bits
 - significand precision: 24 bits (23 explicitly stored)

 $ightarrow (-1)^{b_{31}} imes 2^{(b_{30}b_{29}\dots b_{23})_2 - 127} imes (1.b_{22}b_{21}\dots b_0)_2,$

- values up to (2 − 2⁻²³) × 2¹²⁷ ≈ 3.402823 × 10³⁸
- still limited (relative) precision, e.g. estimating (v+1) -v can be 0 !
- Half-/double-precision with 16/64 bits = 2/4 bytes





Program

- Bits & bytes
- Data format
- Signal discretization
- File format & compression
- Storage & Safety



Signal discretization

Some continuous values =

- 1. measured somewhere, and
- 2. stored numerically
- \rightarrow discretized value with **finite resolution**!

Two faces of "resolution" \rightarrow Different file weight!

- time/space \rightarrow sampling rate
- amplitude \rightarrow encoding precision



Encoding precision

How is the value represented on disk?

- Integer vs. float?
- Number of bytes?
- \rightarrow Different resolution





Sampling rate

How sparse/coarse are data sampled?

- \rightarrow sampling rate
- \rightarrow Nyquist theorem:

"Sampling Rate > 2 x highest frequency of signal"





Example for 3D image

Consider a 3D image with 256 x 256 x 128 = 2^{23} voxels

- ▶ 1 int16 per voxel → 16 Mb
- ▶ 1 float32 per voxel \rightarrow 32 Mb

Coloured image

→ 3 RGB values par voxel, e.g. 3 int8 per voxel → 24 Mb

Resample at half the resolution, i.e. 128 x 128 x 64 voxels \rightarrow divide sizes by 8



Program

- Bits & bytes
- Data format
- Signal discretization
- File format & compression
- Storage & Safety



File format

Open vs. closed file format:

- fully described vs. proprietary
- openly readable vs. requiring specific software
- community supported vs. software/company dependent

- \rightarrow Stick to open format whenever possible
- \rightarrow More flexibility to use with homemade software



The case of MS Word & Excel

Both are proprietary and cost €€€ + files are "binarized"

Word & .doc files, replace by

 \rightarrow 'MarkDown' (. md) files

 \rightarrow open editor/reader, e.g. Typora (<u>https://typora.io/</u>)

Excel & .xls files , replace by

 \rightarrow 'comma-separated value' or 'tab-separated value' (.csv/.tsv) files

 \rightarrow open editor/reader, e.g. CSVed (<u>https://csved.sjfrancke.nl/</u>)

Whenever possible and appropriate

| T DataComments.md - Typora - 🗆 X | T DataComments.md - Typora — |
|--|--|
| <u>F</u> ile <u>E</u> dit <u>P</u> aragraph F <u>o</u> rmat <u>V</u> iew <u>T</u> hemes <u>H</u> elp | <u>F</u> ile <u>E</u> dit <u>P</u> aragraph F <u>o</u> rmat <u>V</u> iew <u>T</u> hemes <u>H</u> elp |
| Some comments about the data. | ## Some comments about the data. |
| Overall ~79Gb: (~58k files & 208 folders) | Overall ~79Gb: (~58k files & 208 folders) |
| MSHS, 37Gb, 37 subjects | - MSHS, 37Gb, 37 subjects |
| • MSPA, 40Gb, 40 subjects | - MSPA, 40Gb, 40 subjects |
| MSP FLAIR/mask, 2.5Gb, 40 subjects | - MSP FLAIR/mask, 2.5Gb, 40 subjects |
| MSPA: possibly to exclude. s08825. Rather visible movement artefacts. Poor positioning in scanner -> cerebellum out of FOV? s00349. Some movement artefact + hyper-instensities (artefact) in orbito-frontal area for MT. s00356. hyper-instensities (artefact) in orbito-frontal area for MT + small meningiome between the frontal hemispheres. | <pre>#### MSPA: possibly to exclude. **s08825**. Rather visible movement artefacts. Poor positioning in scanner -> cerebellum out of FOV? **s00349**. Some movement artefact + hyper-instensities (artefact) in orbito-frontal area for MT.</pre> |
| | **s00356**. hyper-instensities (artefact) in orbito-frontal area for MT + small meningiome between the frontal hemispheres. |
| | |
| | |



Excel in Genetics

"Gene name errors are widespread in the scientific literature"

<u>Abstract</u>:

The spreadsheet software Microsoft Excel, when used with default settings, is known to **convert gene names to dates and floating-point numbers**. A programmatic scan of leading genomics journals reveals that **approximately one-fifth of papers with supplementary Excel gene lists contain erroneous gene name conversions**.

Ziemann et al., Genome Biology 201617:177



Structured data

Data as

- key/value pairs
- hierarchical structure
- → use 'JavaScript Object Notation', i.e. .json, files

Example, task-Nback_bold.json

```
{
```

```
"RepetitionTime": 3.0,
"EchoTime": 0.0003,
"FlipAngle": 78,
"SliceTiming": [0.0, 0.2, 0.4, 0.6, 0.8, 1.0,
1.2, 1.4, 1.6, 1.8, 2.0, 2.2, 2.4, 2.6, 2.8],
"MultibandAccellerationFactor": 4,
"ParallelReductionFactorInPlane": 2
```



Data compression

Lossless:

- ▶ no data/signal lost → replace "patterns" with fewer bytes (RLE).
- 2-4x compression rate, depending on data
- e.g. ZIP, PNG, JPEG2000

Lossy:

- Removes some signal \rightarrow irreversible loss!
- ► quality factor from 0 to 100 → >10x compression rate
- e.g. JPEG





Program

- Bits & bytes
- Data format
- Signal discretization
- File format & compression
- Storage & Safety

Hard-disk drive

HDD = electromechanical data storage device:

- magnetic storage to read/write data
- on one (or more rigid) rapidly rotating disks
- cheap and storage density increases (Moore's law)
- ► latency = ~a few ms,
- ▶ transfer rate up to ~1 Gb/s
- risk of failure increases with time but...

End of Life Wear-Out Increasing Failure Rate Decreasing Failure Rate Normal Life (Useful Life) Low "Constant" Failure Rate Time

eased Failure

The Bathtub Curve

Hypothetical Failure Rate versus Time



Solid-state drive

SSD = integrated circuit data storage device:

- non-volatile NAND flash memory to read/write data
- no mechanical or moving part
- latency < ms,</p>
- transfer rate up to a few Gb/s
- compared to HDD
 - more expensive and more reliable
 - less power consumption





ULiège mass-storage

- Personal space \rightarrow your own stuff
- ▶ Platform space \rightarrow raw data access
- Team space \rightarrow shared data & results

Keep in mind access time

 \rightarrow no direct processing of data!



Backup vs. Archive

Backup

- copy of current data/system
- includes files which are currently being accessed/changed
- → Restoring data/system to a previous point in time, if they are lost or become corrupted

Archive

- store data/information to be kept for a long period of time
- includes files which should not be modified, accidentaly or purposely
- \rightarrow Restoring the 'original' data/information, e.g. to re-analyse them



Local vs. Remote storage

Local, e.g. USB drive

- Cheap and easy
- Can be lost or corrupted with the rest of the computer
- \rightarrow Better than nothing but not so safe!

Remote, e.g. institutional mass-storage

- More expensive (for the institution/users) and more constraining (network access)
- Little risk of losing anything (tapes, redundant disks, multi-sites,...)
- \rightarrow Safest option, if available

For code, use versioning \rightarrow more on Thursday!



References

- https://en.wikipedia.org/wiki/Bit
- https://en.wikipedia.org/wiki/Byte
- https://en.wikipedia.org/wiki/Character (computing)
- https://en.wikipedia.org/wiki/ASCII
- https://en.wikipedia.org/wiki/UTF-8
- https://en.wikipedia.org/wiki/Integer_(computer_science)
- https://en.wikipedia.org/wiki/Single-precision_floating-point_format
- https://en.wikipedia.org/wiki/Nyquist%E2%80%93Shannon_sampling_theorem



References

- https://en.wikipedia.org/wiki/Markdown
- https://typora.io/
- https://en.wikipedia.org/wiki/Comma-separated_values
- https://en.wikipedia.org/wiki/Tab-separated_values
- https://csved.sjfrancke.nl/
- https://en.wikipedia.org/wiki/JSON
- https://doi.org/10.1186/s13059-016-1044-7
- https://en.wikipedia.org/wiki/Run-length_encoding
- https://en.wikipedia.org/wiki/JPEG
- https://en.wikipedia.org/wiki/Hard_disk_drive
- https://en.wikipedia.org/wiki/Solid-state_drive





Thank you for your attention!

