# Data challenges in Earth Sciences
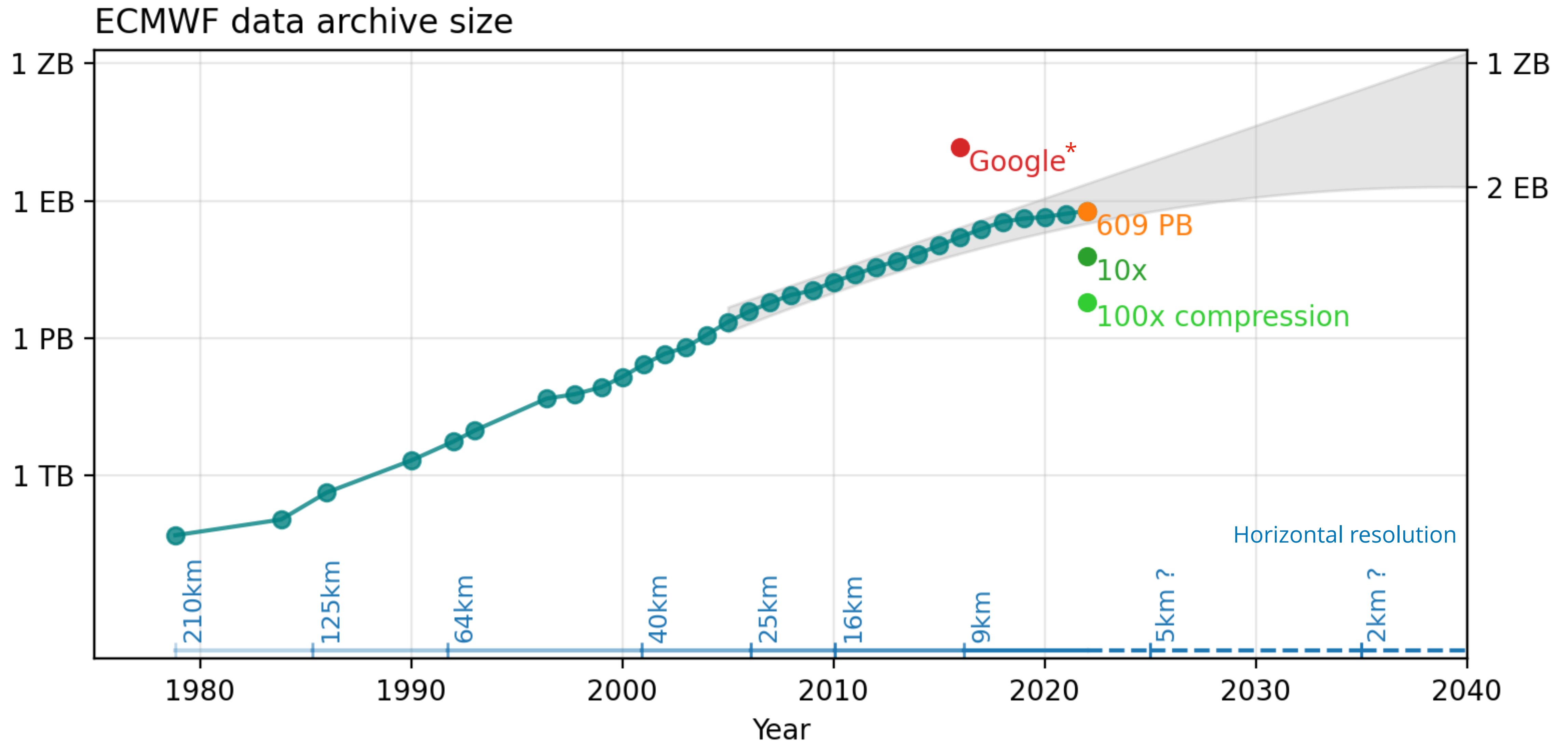
*What is the real information in weather and climate data?*

**Milan Klöwer**
Massachusetts Institute of Technology
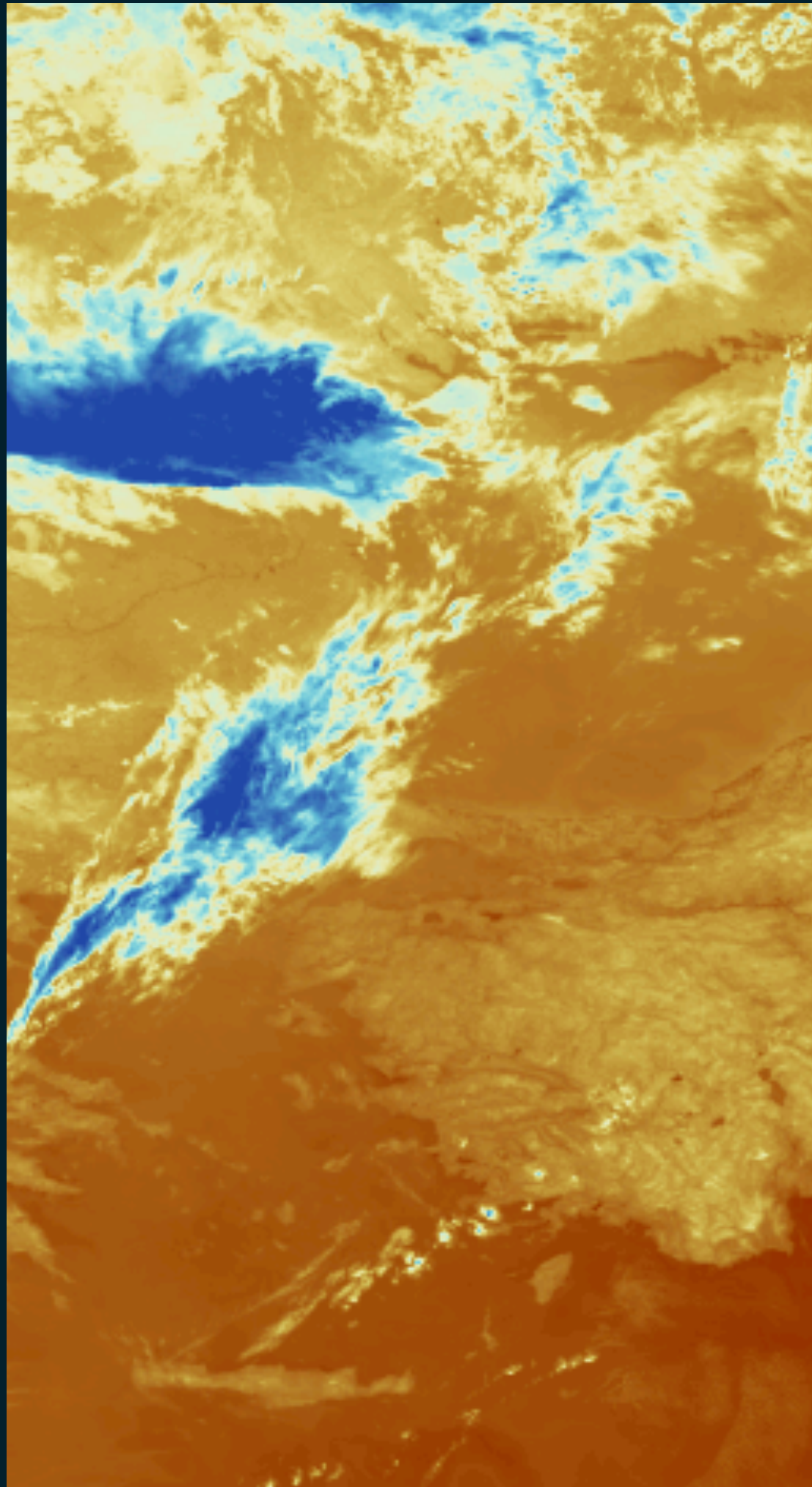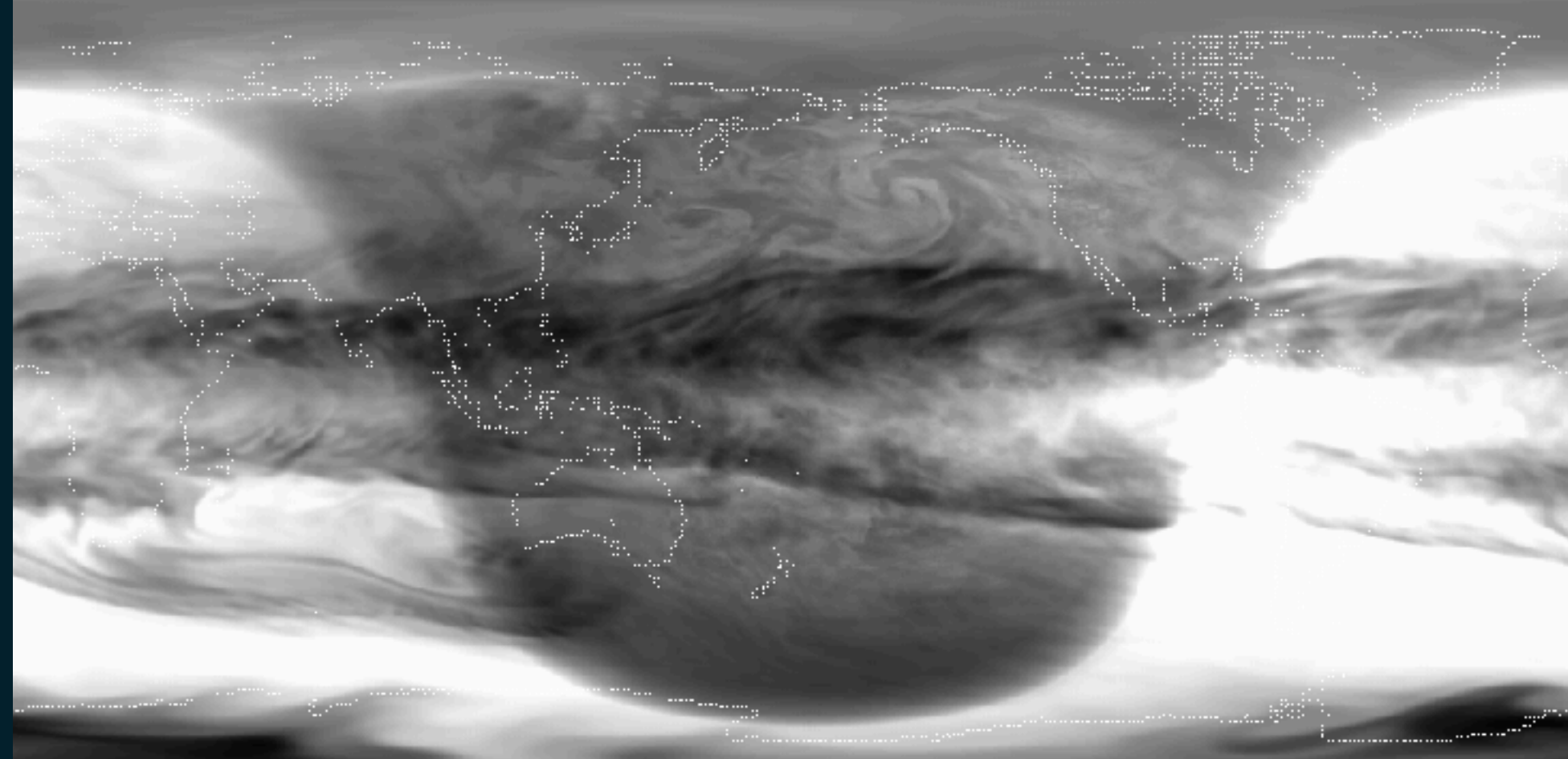
# Earth system data: Will we enter the *Google regime*?



ECMWF data archive size

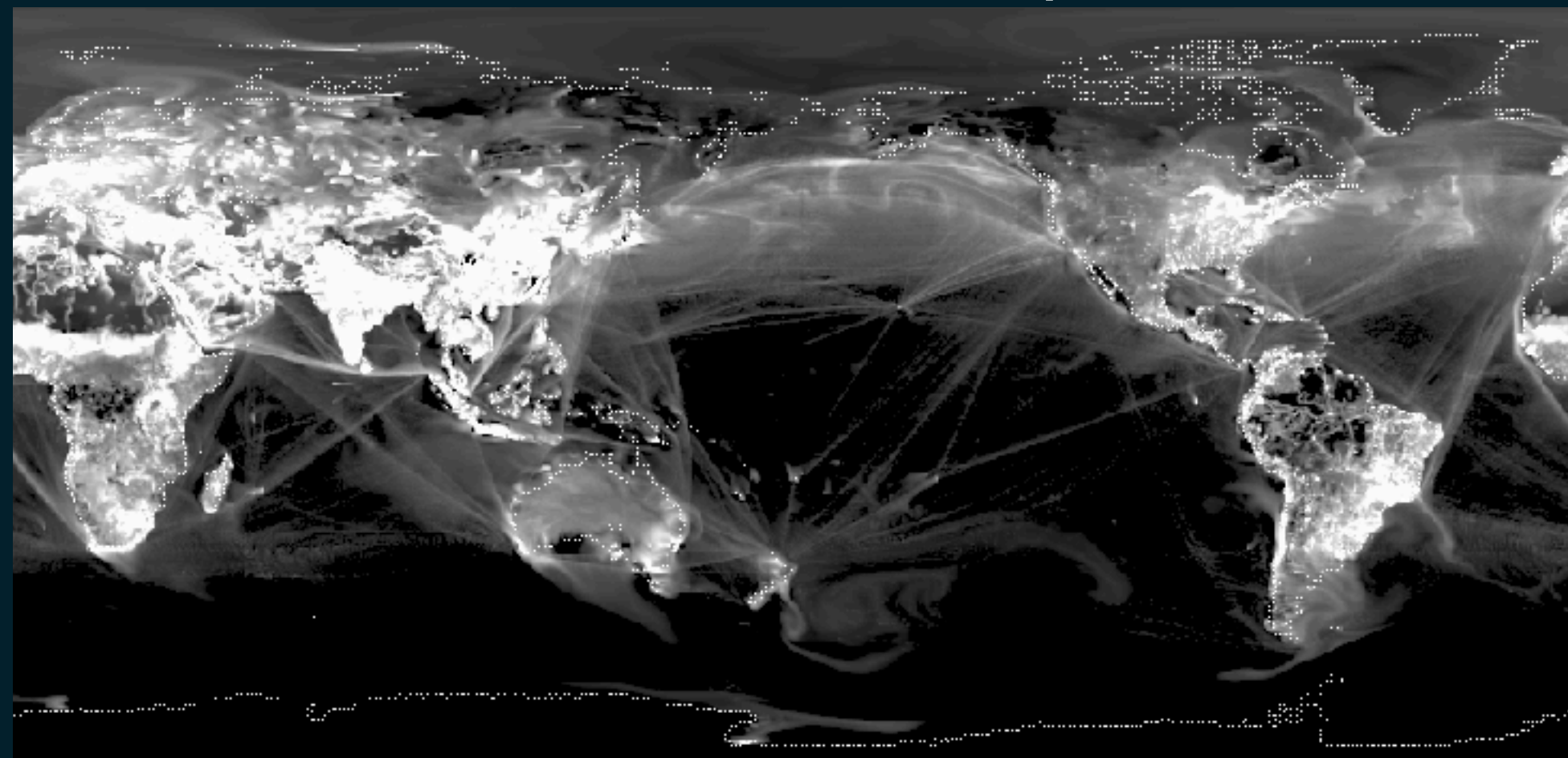*estimate from XKCD, 2016

# What data are we dealing with?



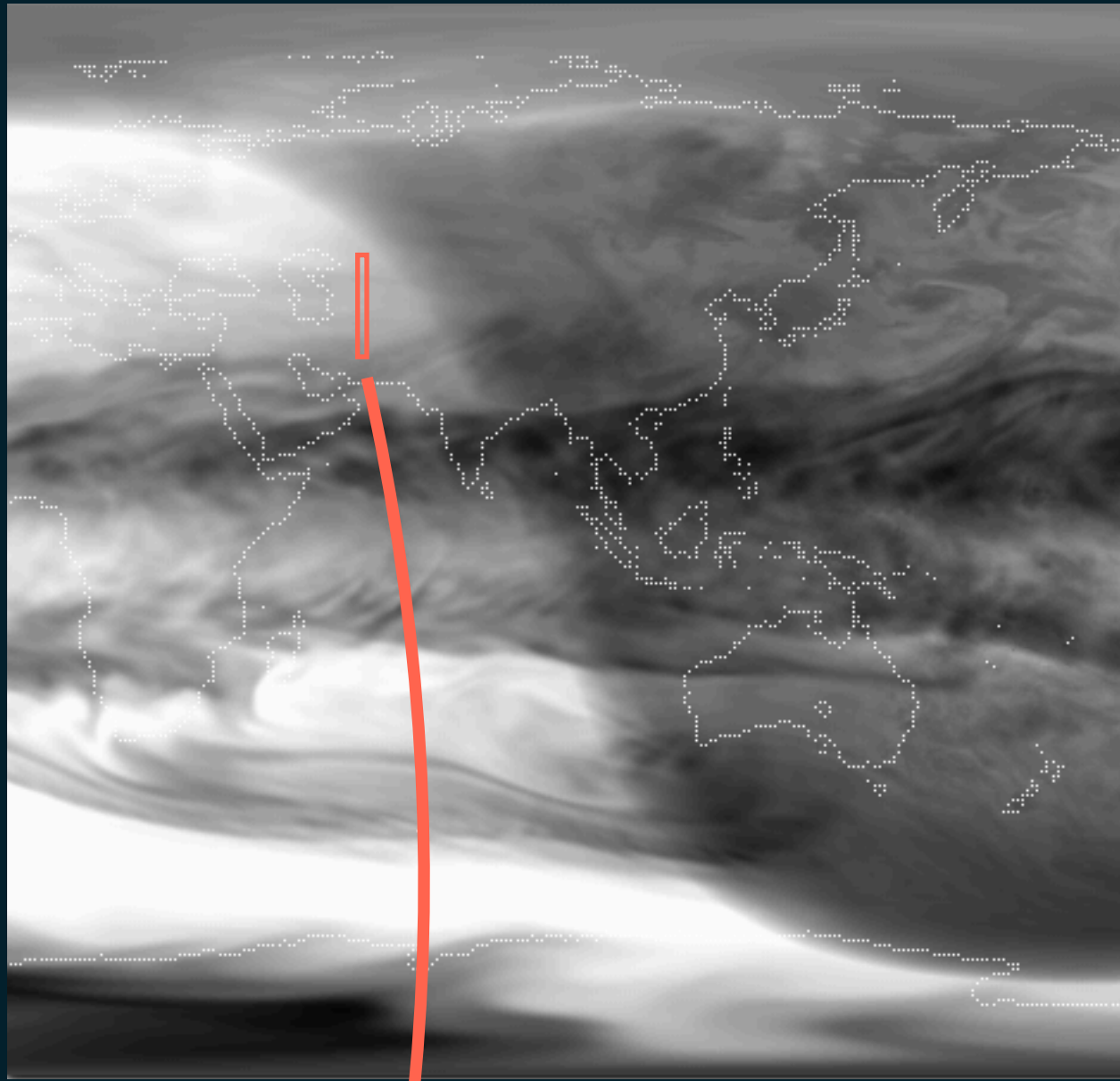Brightness temperature



NO$_2$ in the stratosphere



NO$_2$ at the surface

- **Many** different variables
- Varying uncertainties
- Linear or log-distributed
- Possibly many zeros

- Smoothness varies spatially
- Strong gradients
- Point sources

- Unstructured grids
- Spectral coefficients
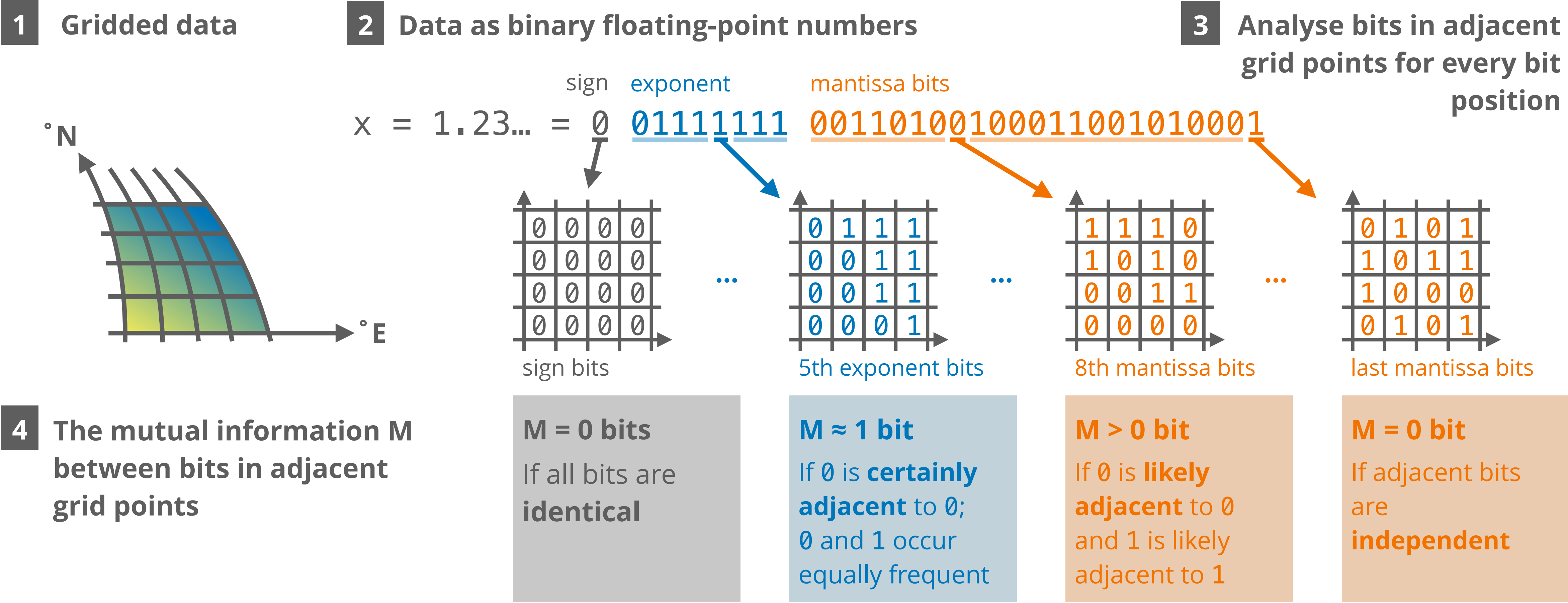- Masked data

# What is real and false information in data?

*Problem: What is the uncertainty in data and how can it be estimated if unknown?*

Possible solution:
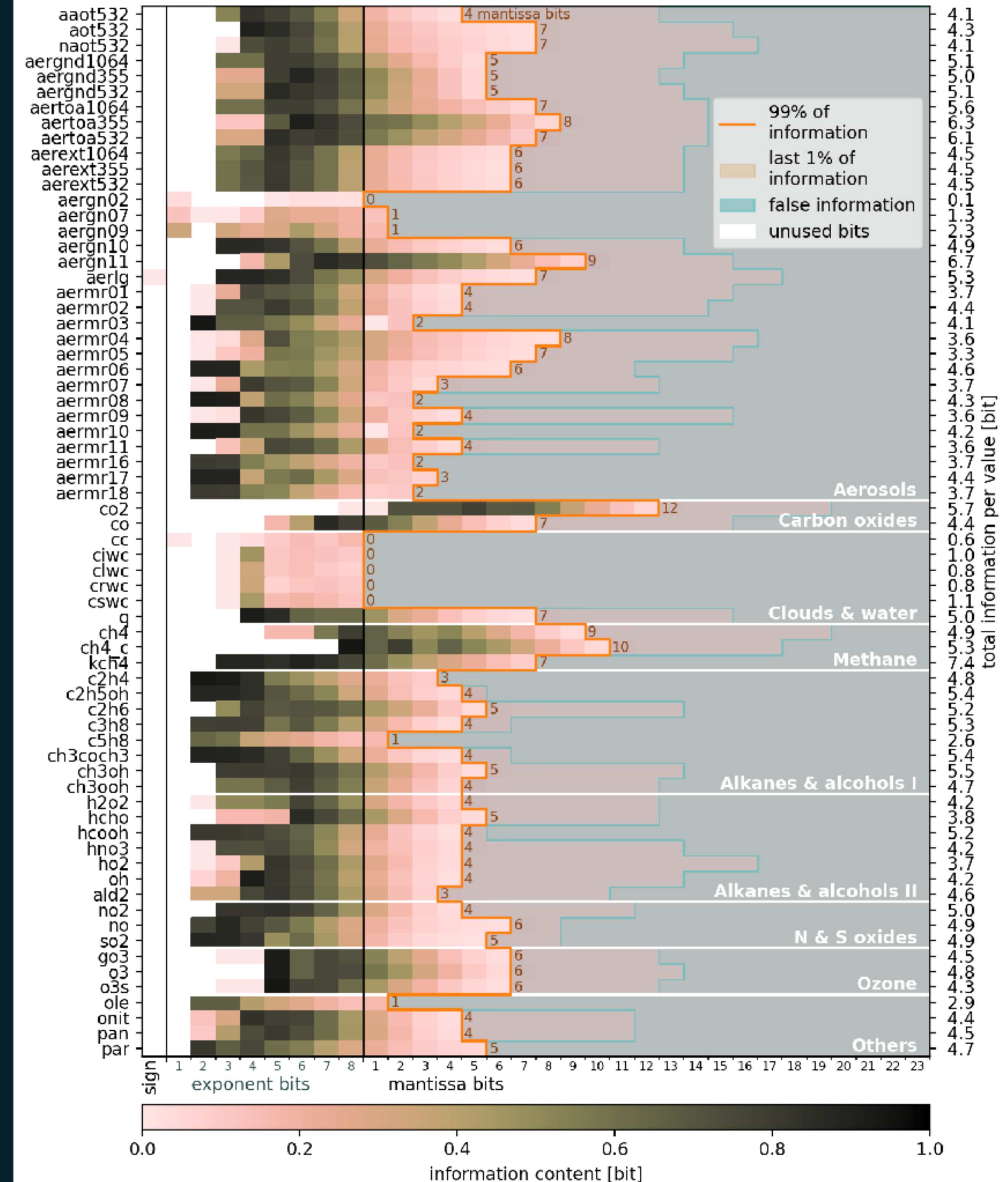**Find an objective way to separate real and false information!**

0.050386034
0.050390966
0.05040059
0.050441727
0.05046302
0.05046855
0.050488267
0.05050127
0.050520953
0.05052939
0.050532646

Encoded in bits

Real!?                                          False!?

00111101010011100110000110010110
00111101010011100110011011000010
00111101010011100111000011011001
00111101010011101001101111111100
00111101010011101011001001010000
00111101010011101011100000011100
00111101010011101100110011001001
00111101010011101101101001101011

Compressible                              Incompressible

# Bitwise real information content

defined here as the **mutual bitwise information in adjacent grid points**

**1** **Gridded data**

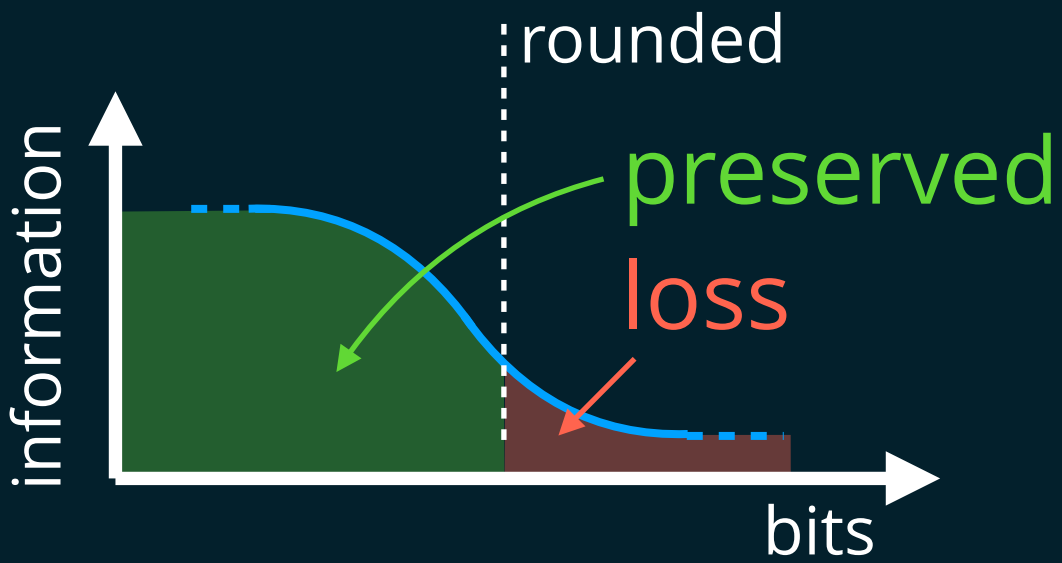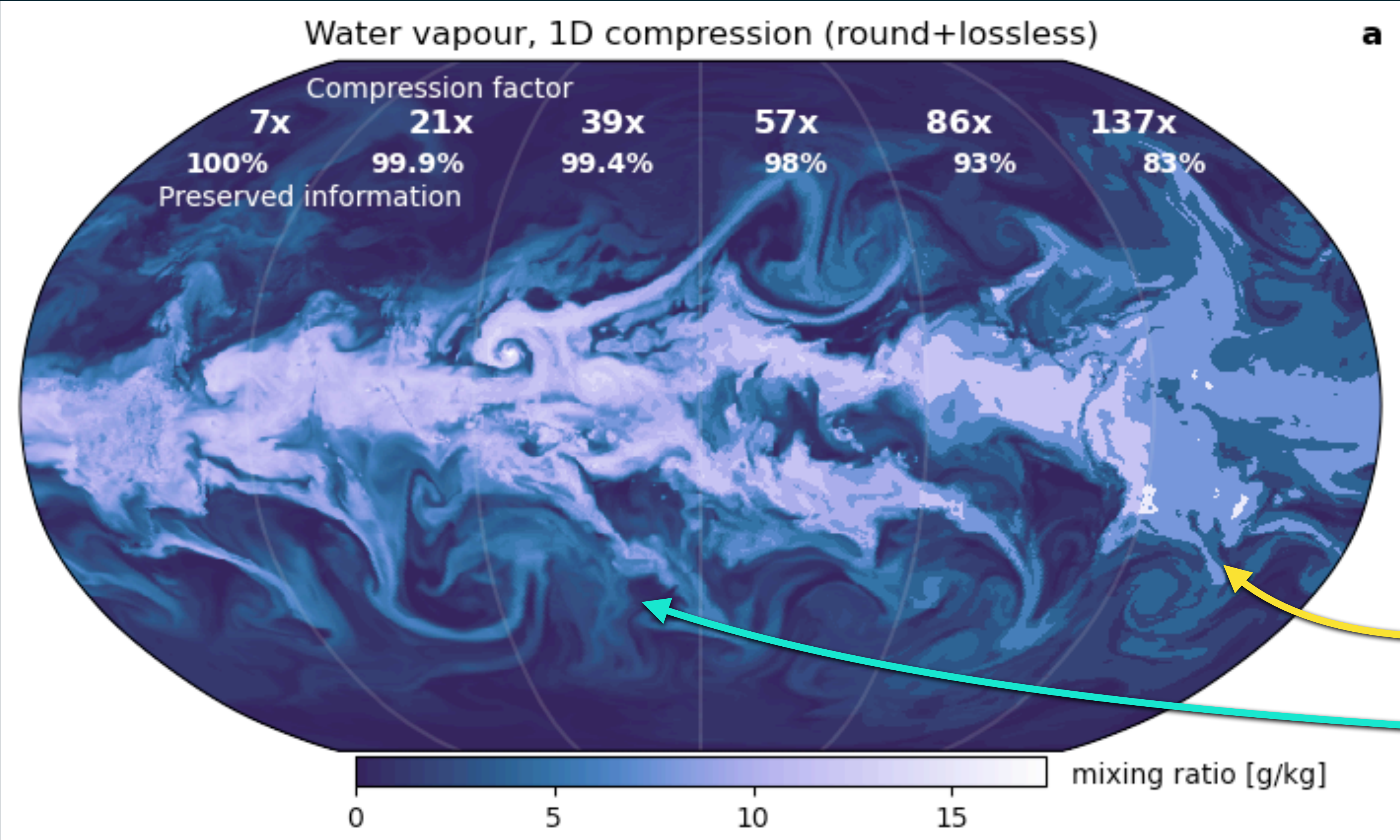**2** **Data as binary floating-point numbers**

**3** **Analyse bits in adjacent grid points for every bit position**

°N

°E

sign   exponent     mantissa bits

$x = 1.23... = 0\ 01111111\ 00110100100011001010001$

| 0 0 0 0 |
| 0 0 0 0 |
| 0 0 0 0 |
| 0 0 0 0 |

sign bits

...

| 0 1 1 1 |
| 0 0 1 1 |
| 0 0 1 1 |
| 0 0 0 1 |

5th exponent bits

...

| 1 1 1 0 |
| 1 0 1 0 |
| 0 0 1 1 |
| 0 0 0 0 |

8th mantissa bits

...

| 0 1 0 1 |
| 1 0 1 1 |
| 1 0 0 0 |
| 0 1 0 1 |

last mantissa bits

**4** **The mutual information M between bits in adjacent grid points**

**M = 0 bits**

If all bits are **identical**

**M ≈ 1 bit**

If 0 is **certainly adjacent** to 0; 0 and 1 occur equally frequent

**M > 0 bit**

If 0 is **likely adjacent** to 0 and 1 is likely adjacent to 1

**M = 0 bit**

If adjacent bits are **independent**

# Bitwise
# real information content

- Every variable requires a different precision
- Many bits do not contain real information
- Preserve only the bits with real information

# Information-preserving compression



Water vapour, 1D compression (round+lossless)

Compression factor

| 7x | 21x | 39x | 57x | 86x | 137x |
|----|-----|-----|-----|-----|------|
| 100% | 99.9% | 99.4% | 98% | 93% | 83% |

Preserved information

mixing ratio [g/kg]

0    5    10    15

a

rounded
information
preserved
loss
bits

303.25**19283**
303.25**00000**

**Visual artefacts**

**99% preserved information suggested as sweet spot**

# Implementations

| Julia: BitInformation.jl | Python & xarray: xbitinfo | Python, netCDF+GRIB | netCDF+HDF: NCO |
|---|---|---|---|



Milan Klöwer

Schulz, Spring

David Meyer

Zender et al

# BitInformation.jl

## 1. Analyse bitinformation
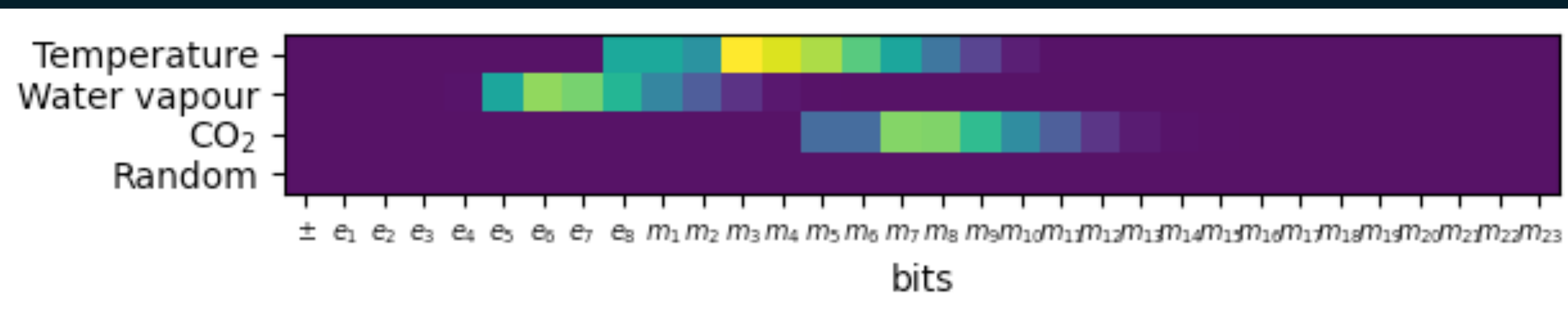
```julia
julia> using BitInformation

julia> A = …

julia> bitinformation(A)
32-element Vector{Float64}:
 0.9999975605213736
 0.2667830274416382
 0.2934084155989065
 0.5817502292436738
 0.9331021403275142
 ⋮
```

## 2. Remove noise

```julia
julia> keepbits = 7

julia> A_round = round(A,keepbits)

julia> bitstring.(A_round)
10000000-element Vector{String}:
 "11000000101010110000000000000000"
 "11000000101000100000000000000000"
 "11000000100111110000000000000000"
 "11000000100110110000000000000000"
 "11000000100110110000000000000000"
 ⋮
```



## 3. Compress

```julia
julia> using JLD2

julia> jldopen("test.jld2","w"; compress=true) do f
           f["A"] = A_round
       end
```