

How much bitwise information is in weather and climate data?

Milan Klöwer¹, Miha Razinger², Juanjo Dominguez²
Aaron Spring³, Hauke Schulz³, Peter Düben², Tim Palmer¹

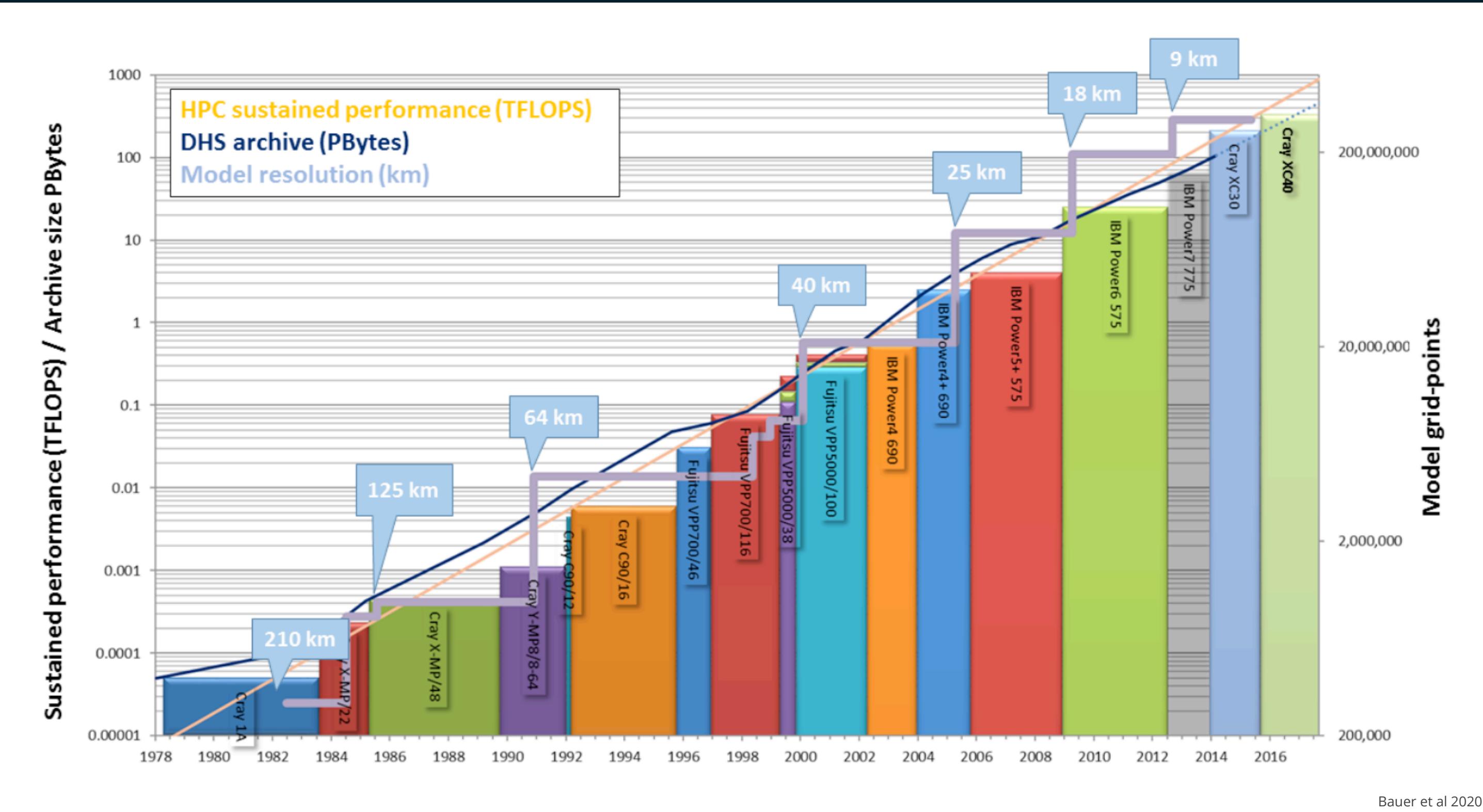
¹University of Oxford, UK

²European Centre for Medium-Range Weather Forecasts, UK

³Max Planck Institute for Meteorology, Hamburg, Germany

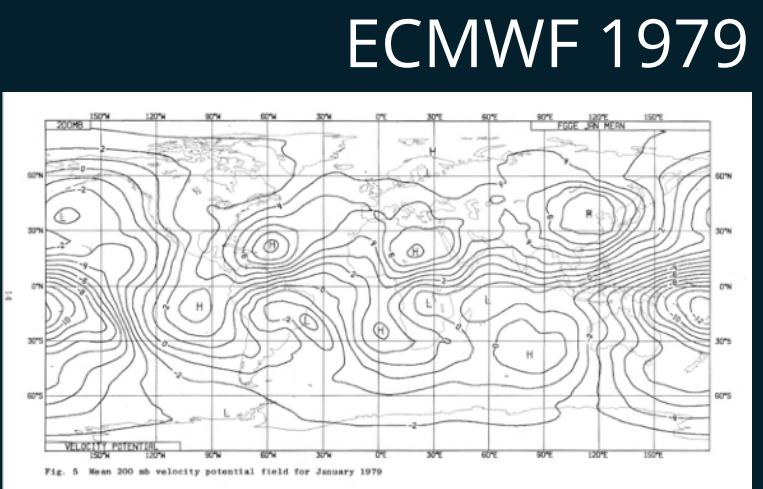
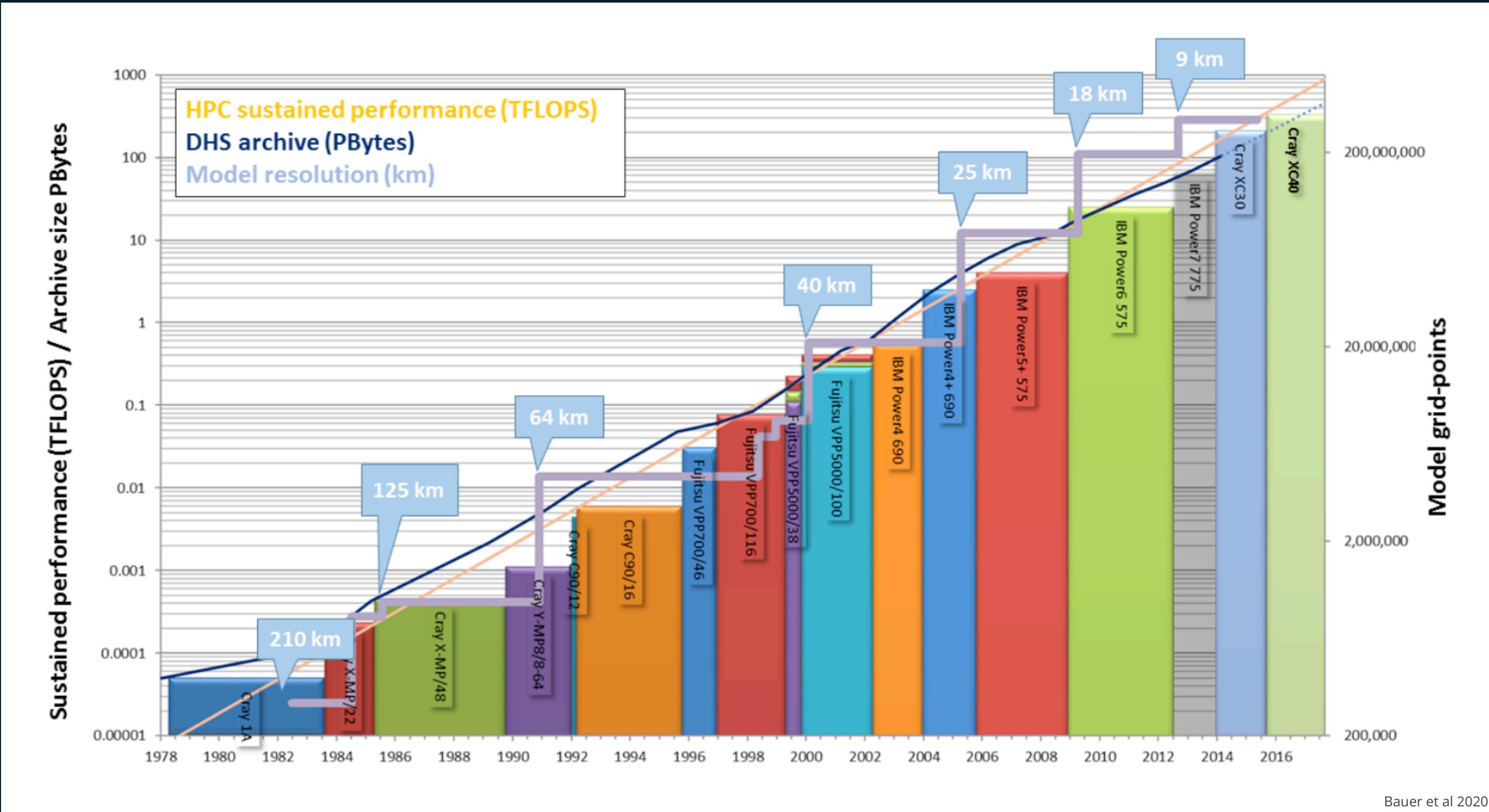


40+ years of numerical weather prediction



- x10⁶ increase in performance and data
- From 200km resolution to 9km
- Ensemble forecasts
- Towards “storm-resolving” models
- Often reliable 7-10 day forecasts

40+ years of numerical weather prediction



IEEE Standard
64-bit floats



32-bit floats

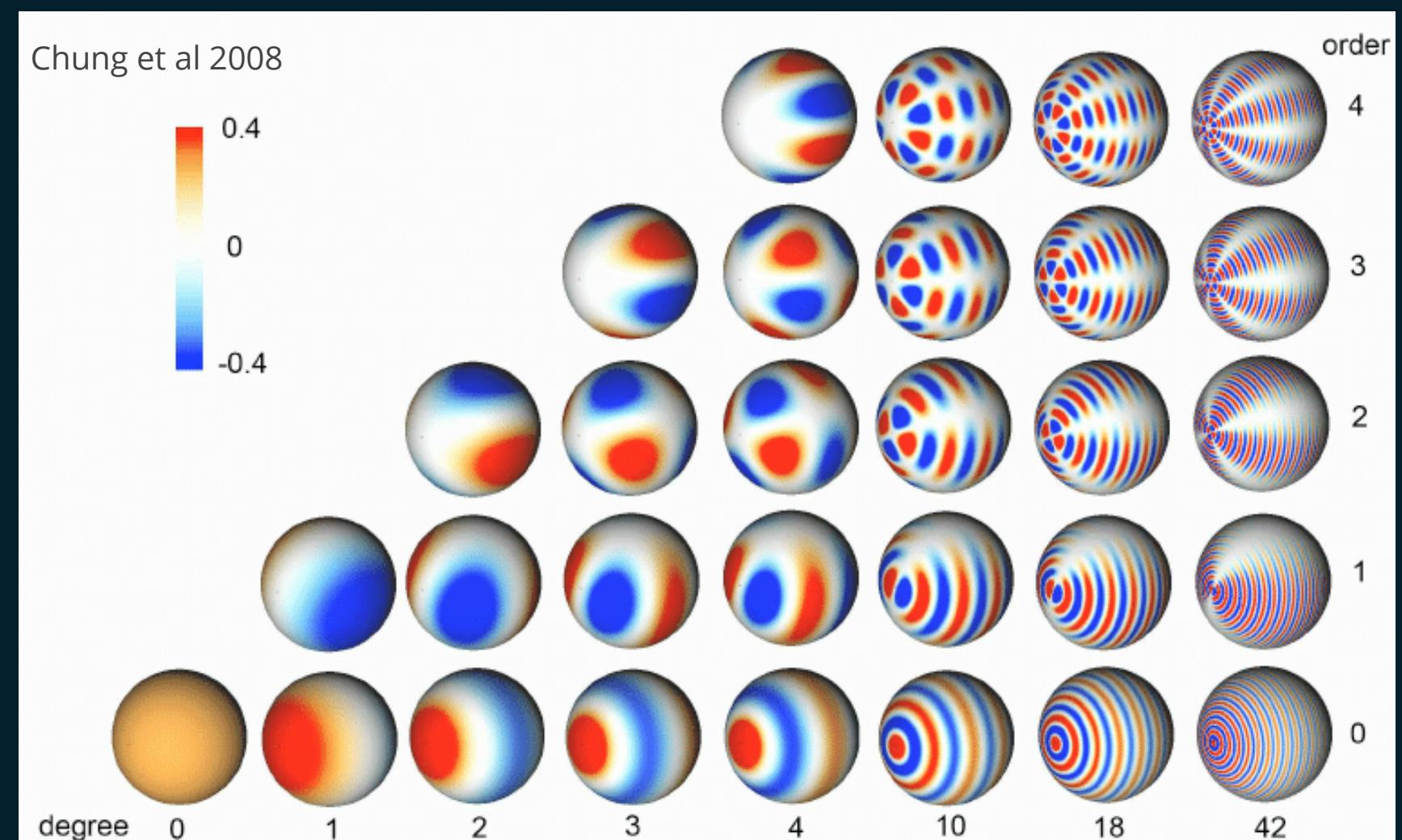
How much have we questioned the bitwise representation of our data?

(computation and archive)

Schools of thoughts: Data compression

School 1: Transformations (the physical perspective)

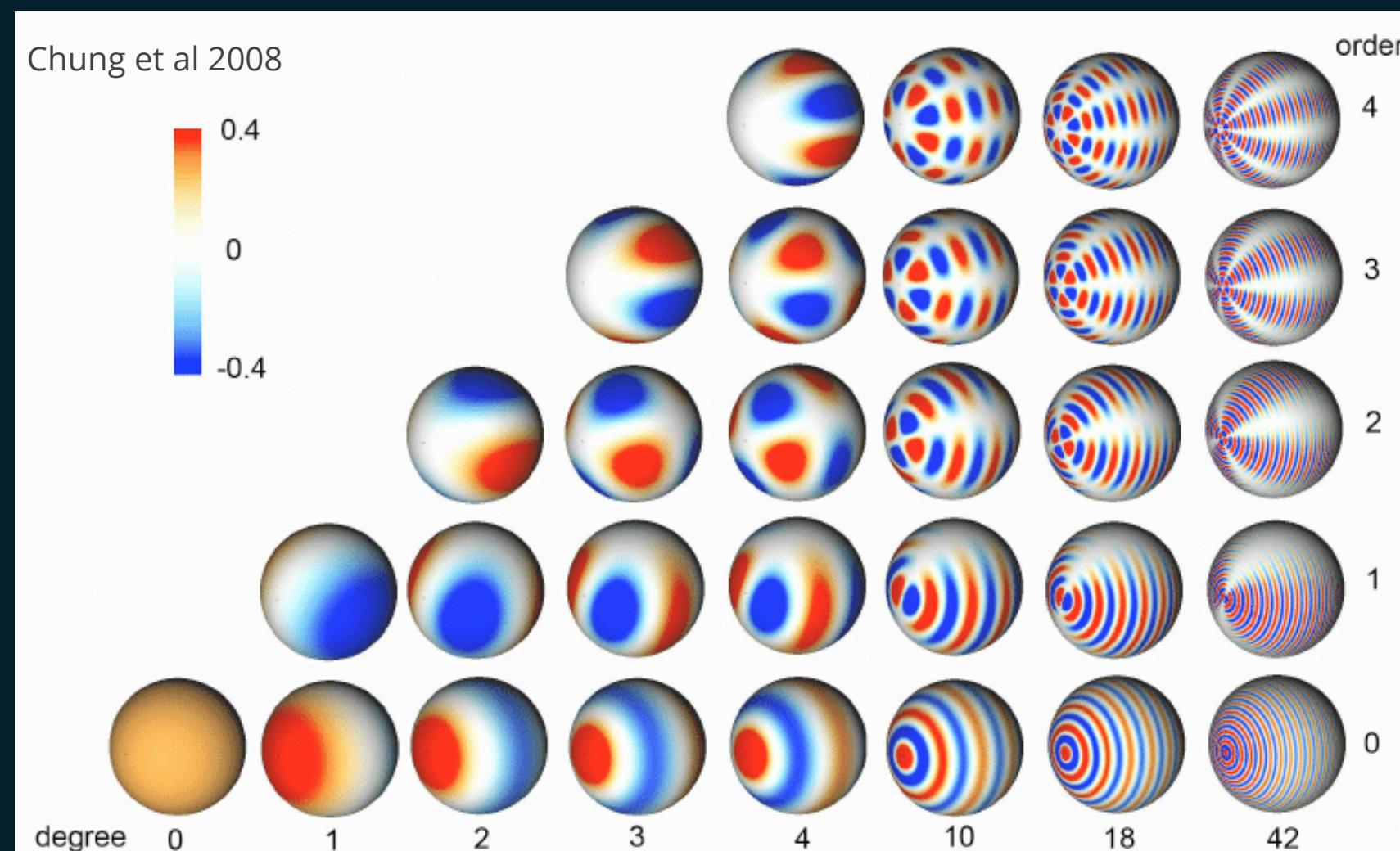
- Spectral
 - EOF
 - Tensor trains
 - ...
- Spatial structure
 - Expensive
 - Error bounds difficult
 - No random access



Schools of thoughts: Data compression

School 1: Transformations (the physical perspective)

- Spectral
- EOF
- Tensor trains
- ...
- Spatial structure
- Expensive
- Error bounds difficult
- No random access



School 2: Precision and information theory (the binary perspective)

- Bitwise encoding
- Floats; linear or logarithmic quantisation
- Entropy coding
- Lossless compression



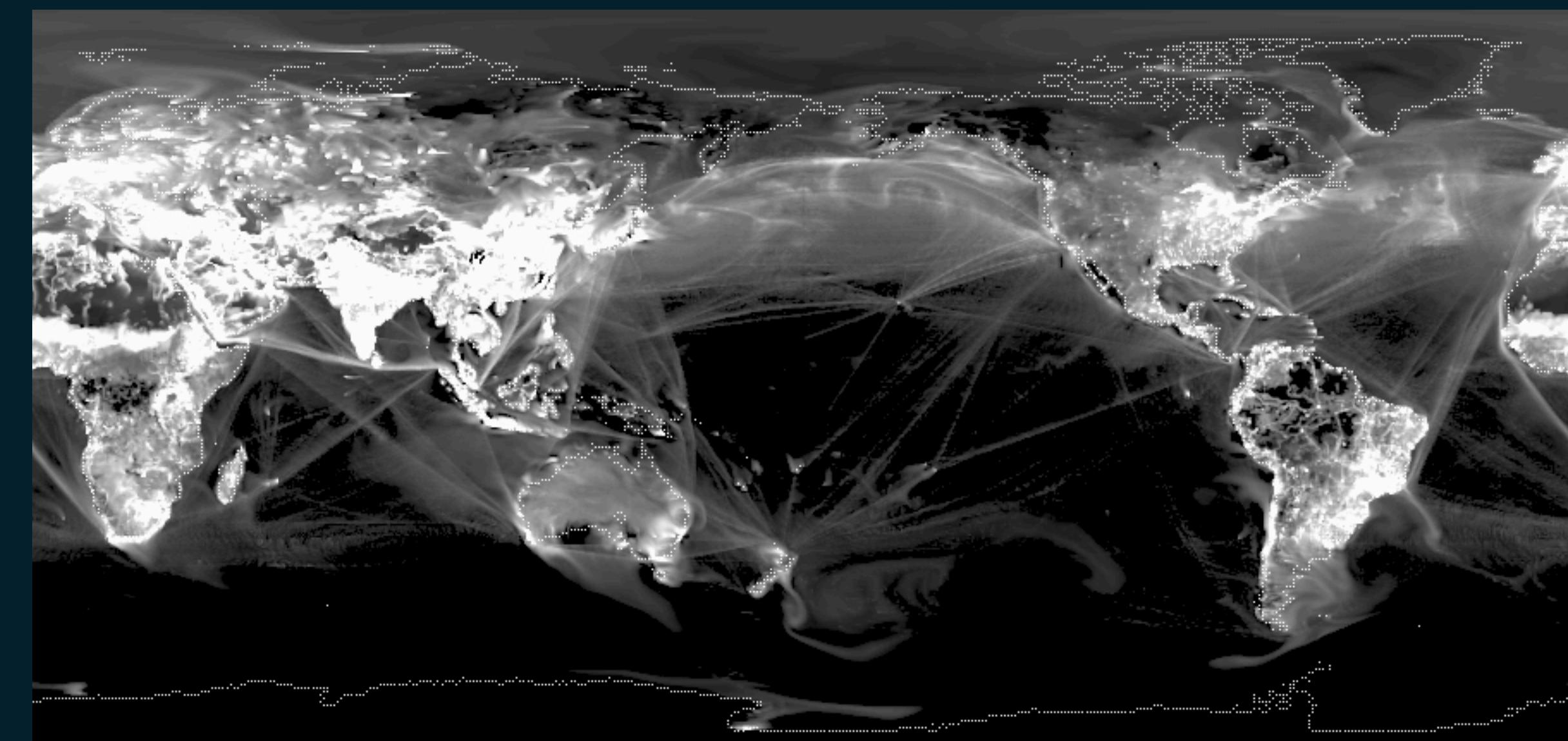
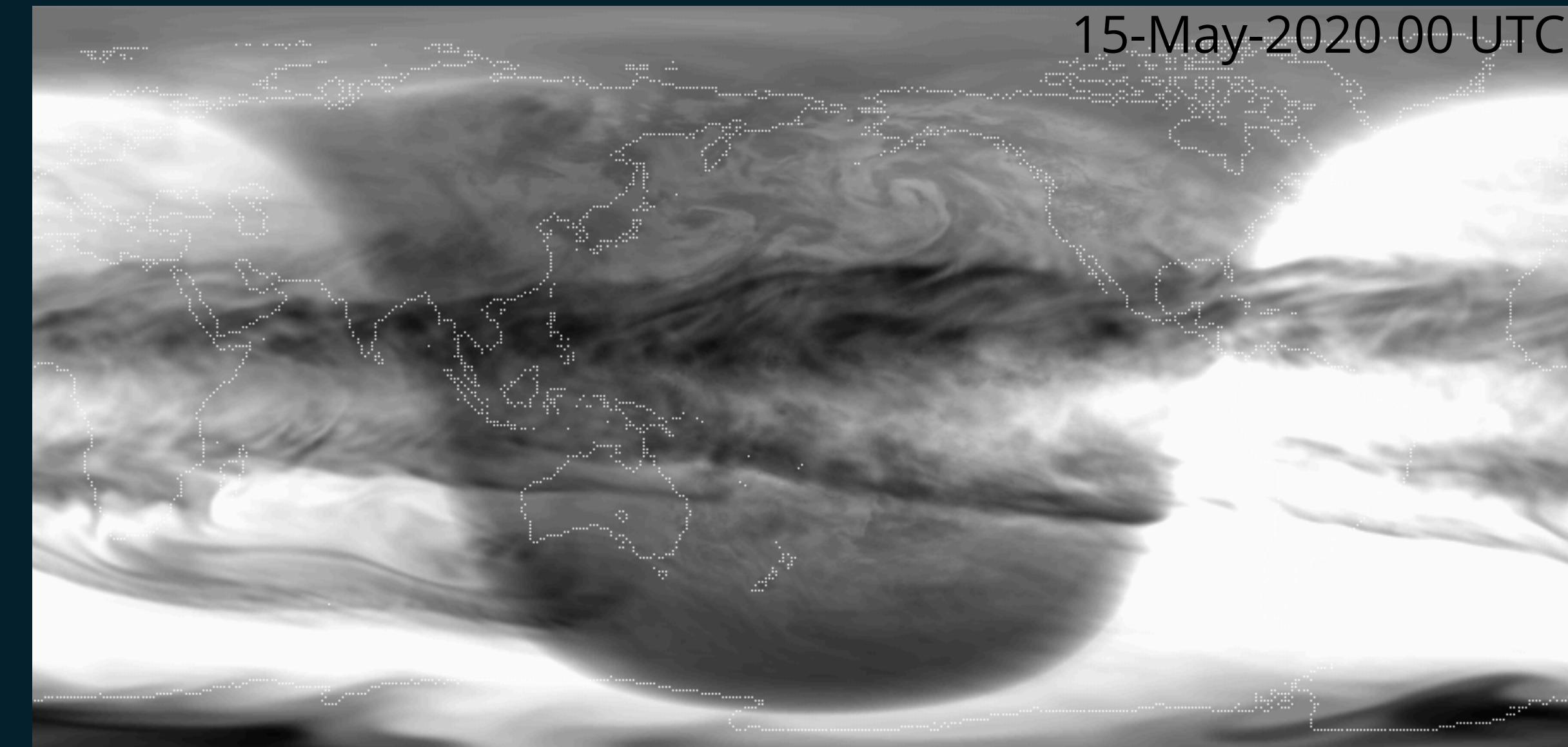
- Spatial structure
- Cheap transforms
- Rigid error bounds
- Random access



Copernicus Atmospheric Monitoring Service (CAMS)

60+ atmospheric variables,
including NO_2 , O_3 , SO_2 , aerosols, ...

- 60+ atmospheric variables
- including NO_2 , O_3 , SO_2 , aerosols, ...
- Varying spatial structures (physics, chemistry)

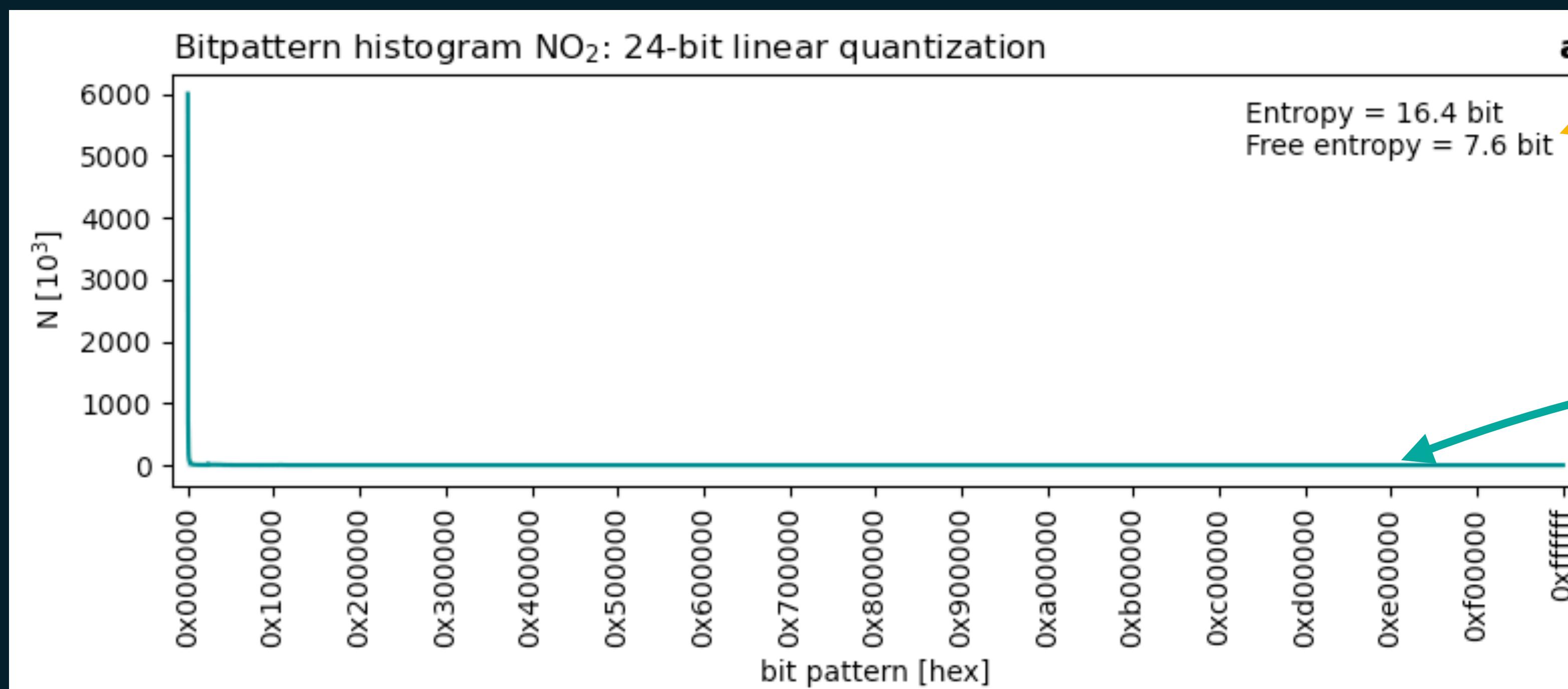


Current compression method: 24-bit linear quantization

split into $2^{n\text{bits}}$ equidistant intervals



Use $n\text{bits} = 24$ instead of 64 bit per number



>7 bits are effectively unused

Most bit patterns very rarely used

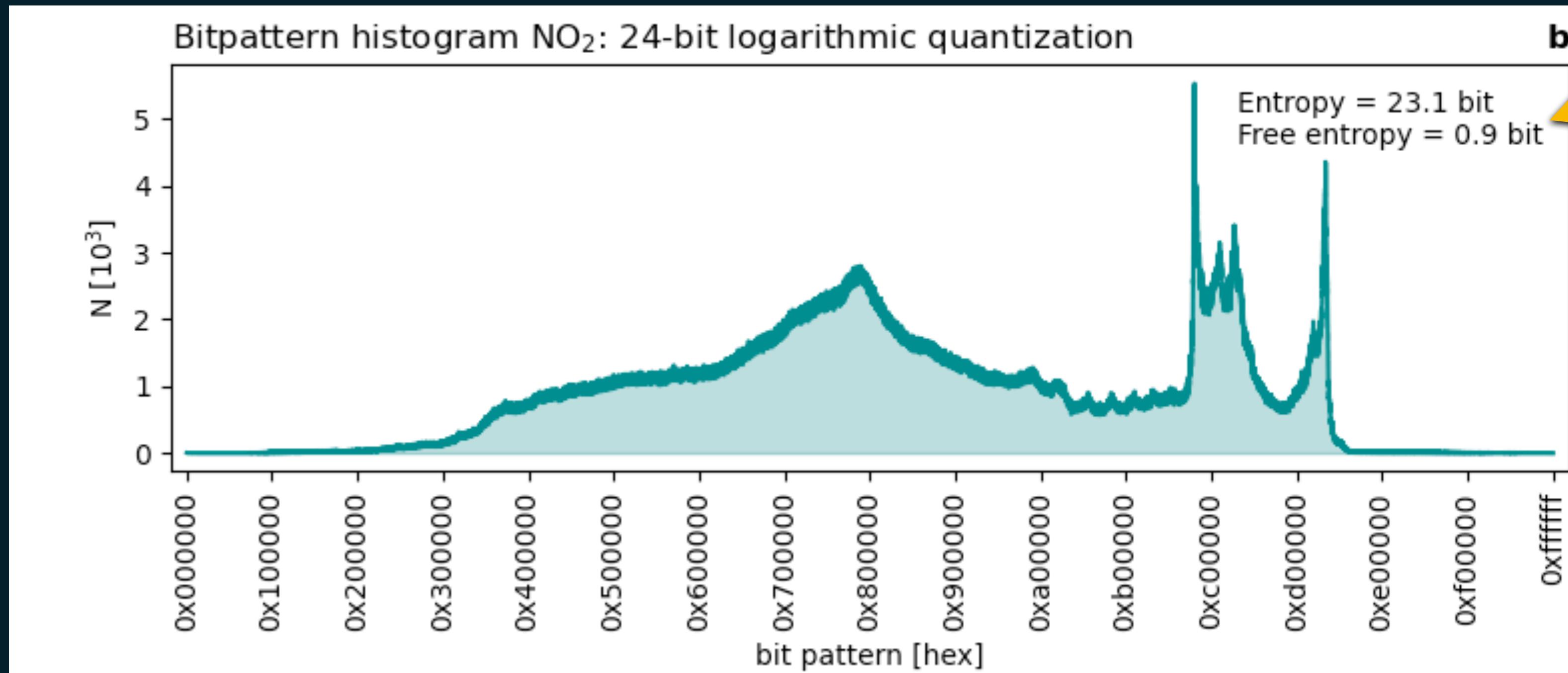
How much information is in the first bit?
~0 bit. Information <= Entropy.

Alternative: Logarithmic quantization

split into $2^{n\text{bits}}$ log-spaced intervals

min

max

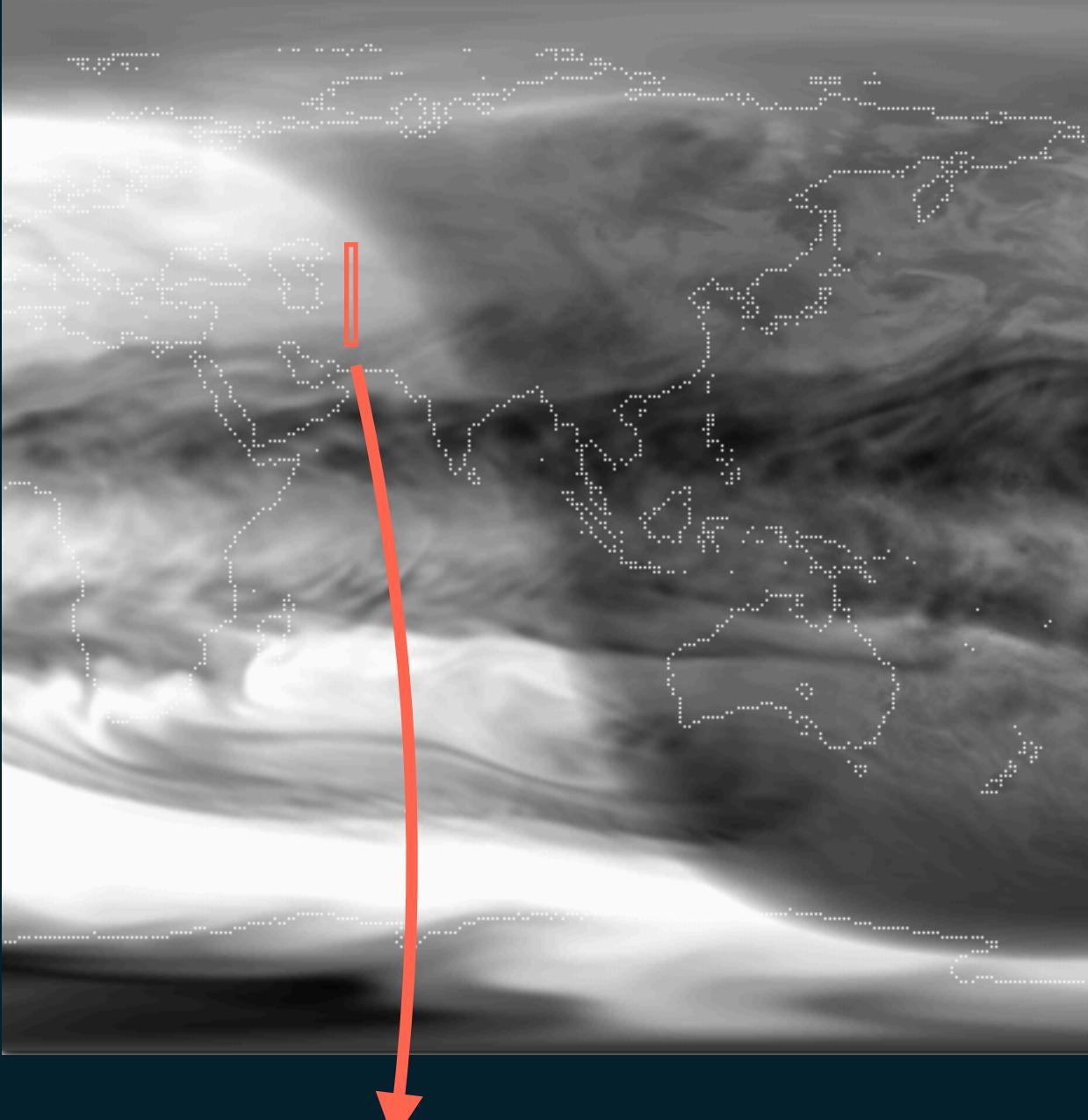


Most bits are effectively used

Floats are also approximately log-distributed

How much entropy is in the first bit?

~1bit



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!?

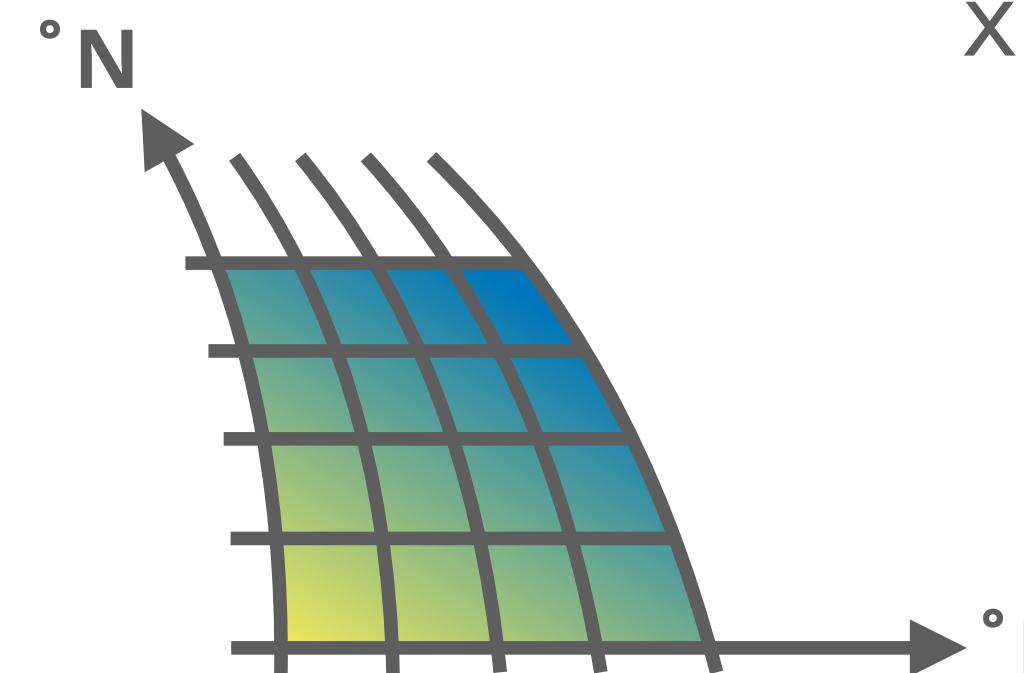
False!?

```
00111101010011100110000110010110  
00111101010011100110011011000010  
00111101010011100111000011011001  
0011110101001110100110111111100  
001111010100111010110010010000  
0011110101001110101100000011100  
0011110101001110101100000011100  
001111010100111010110010011001001  
0011110101001110101100100110101011
```

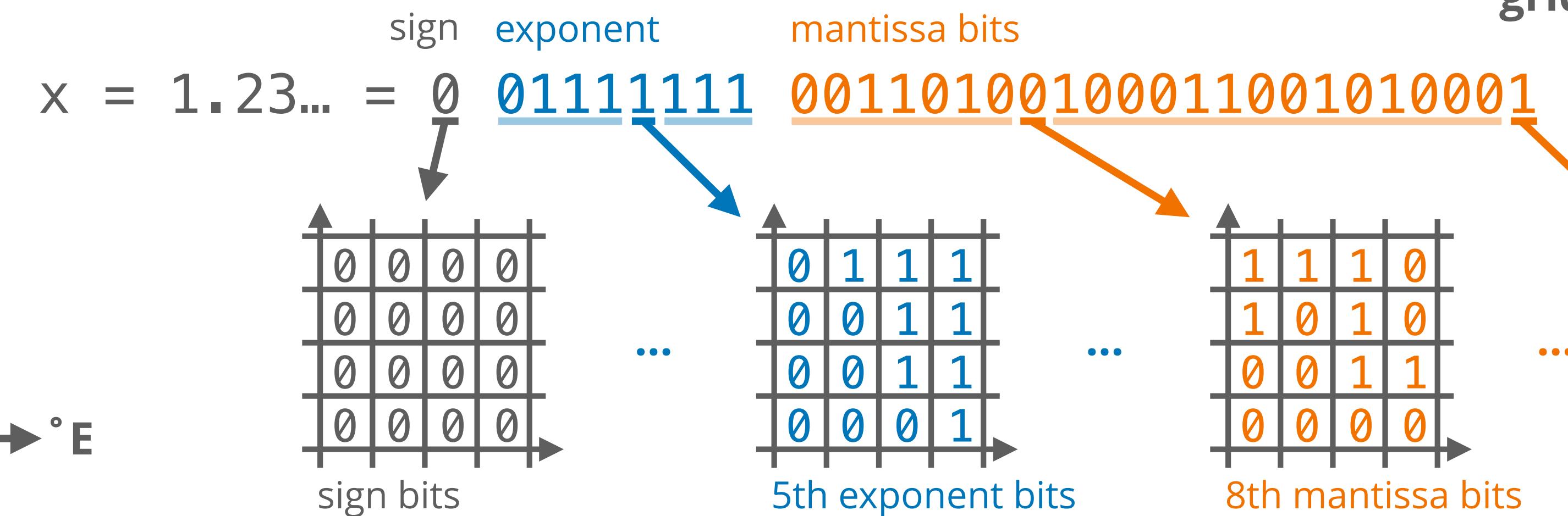
Bitwise real information content

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

4 The mutual information M between bits in adjacent grid points

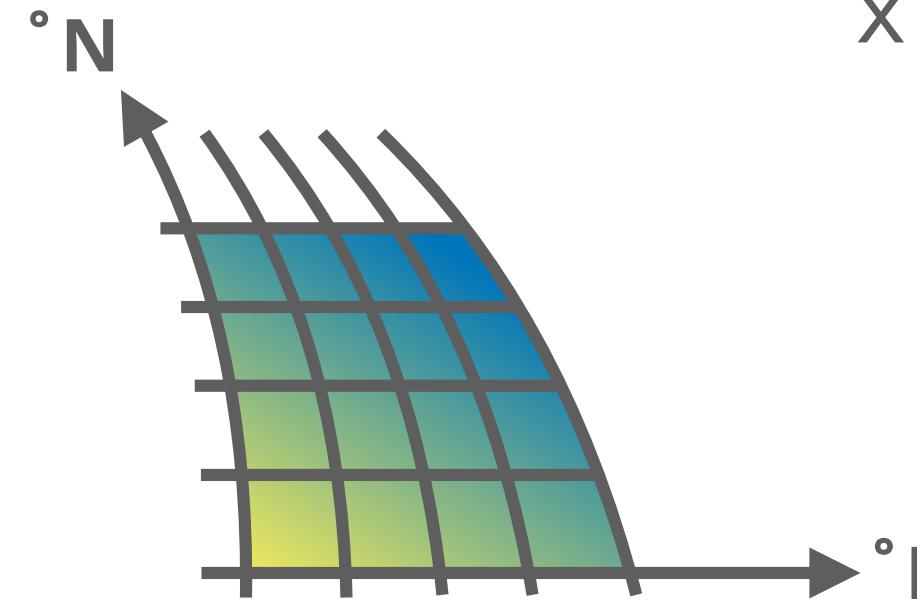
M = 0 bits
If all bits are identical

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

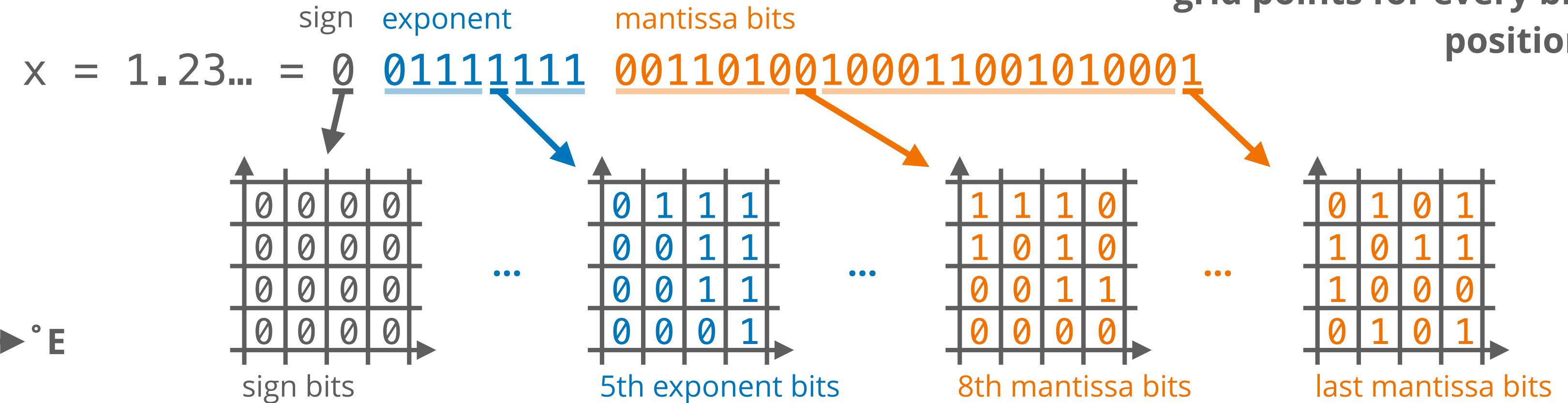
M > 0 bit
If 0 is **likely adjacent** to 0
and 1 is likely adjacent to 1

M = 0 bit
If adjacent bits are **independent**

1 Gridded data



2 Data as binary floating-point numbers



4 The mutual information M between bits in adjacent grid points

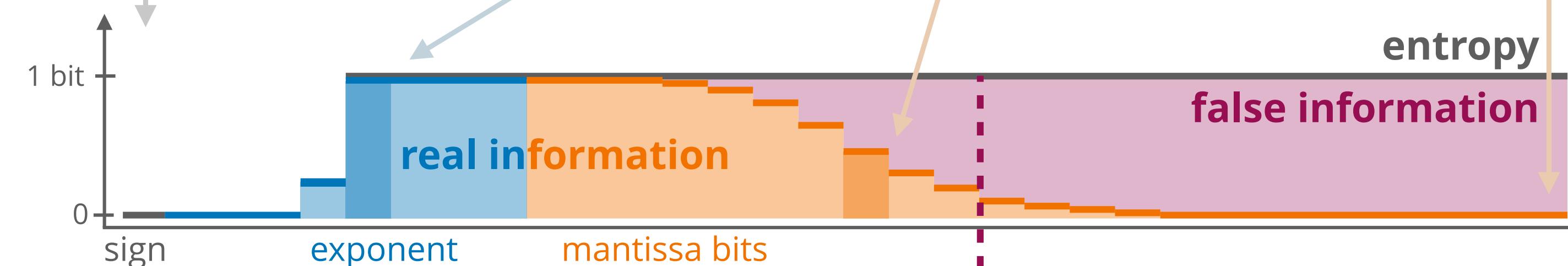
M = 0 bits
If all bits are identical

$M \approx 1$ 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**

5 Bitwise real information is the mutual information between bits



6 Rounding to remove the false information

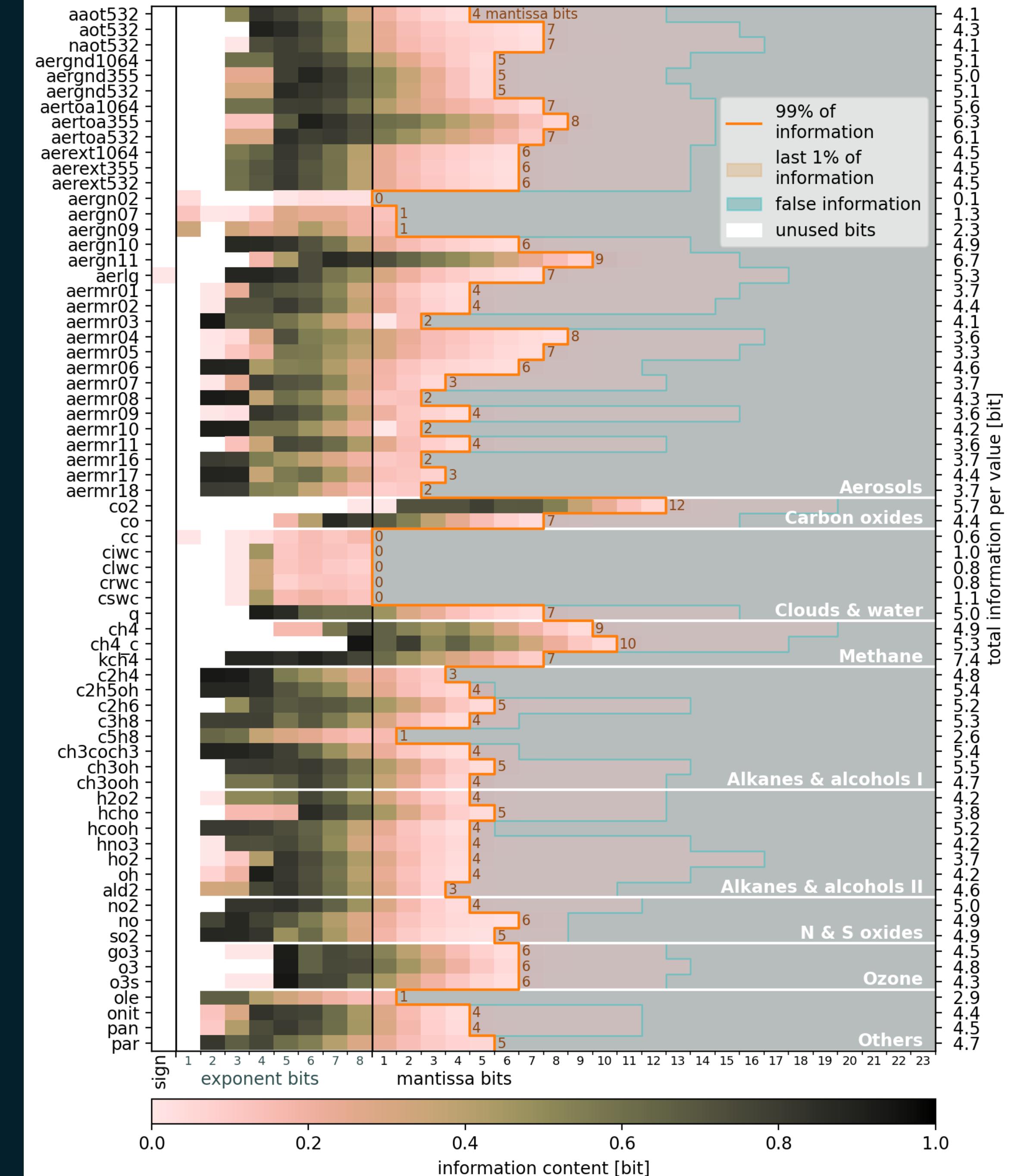
$$x \approx 0 \underline{01111111} \underline{00110100100}0000000000000000$$

Retain bits with **>99% of real information** in total

Remove false information
by rounding trailing bits to 0

Bitwise real information content

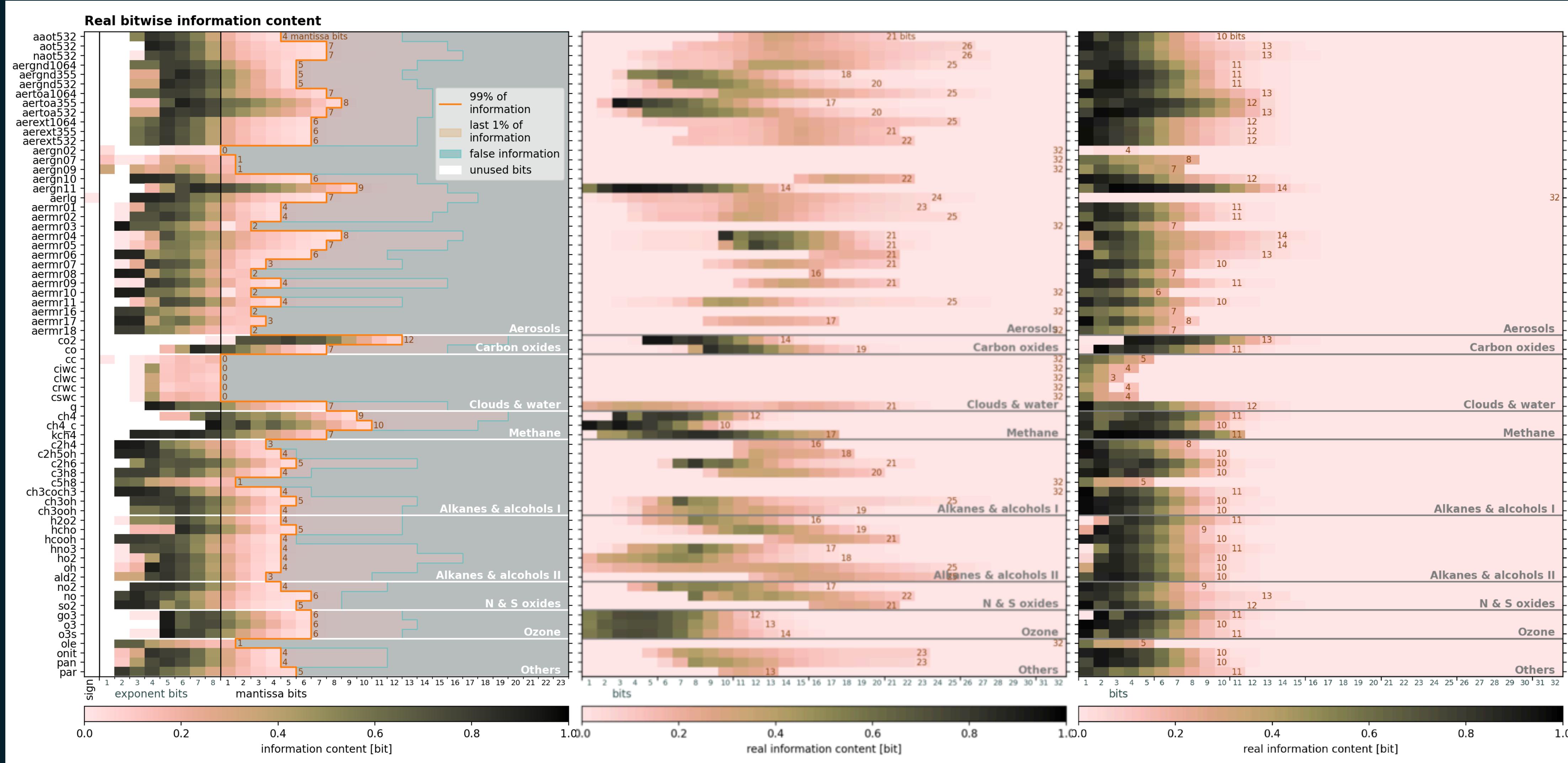
- Every variable requires a different precision
 - Many bits do not contain real information
 - Preserve only the bits with real information



Encoding: Floats

Linear quantisation

Logarithmic quantisation

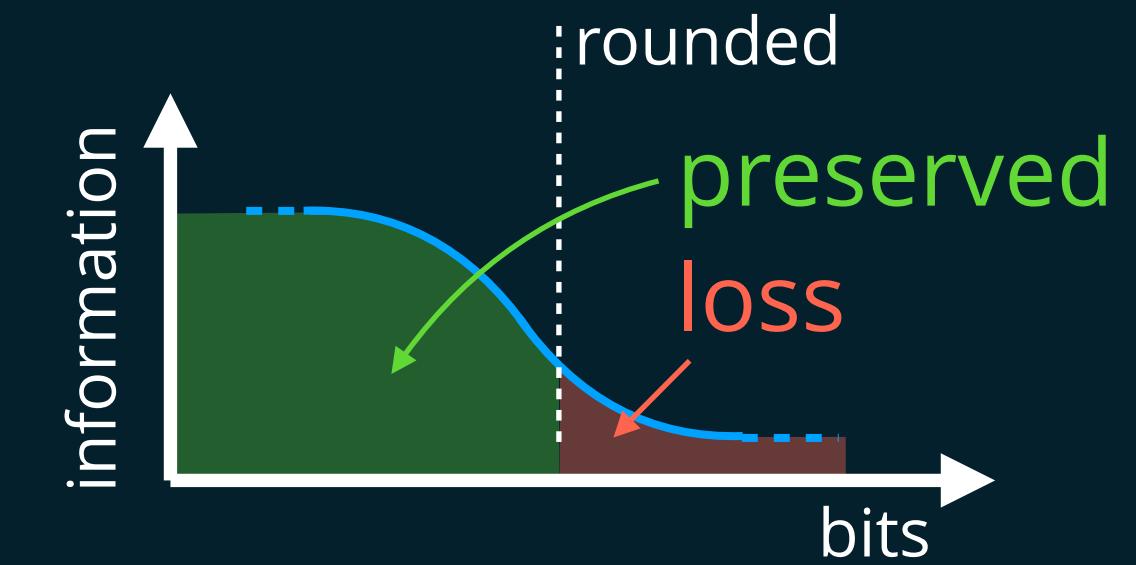
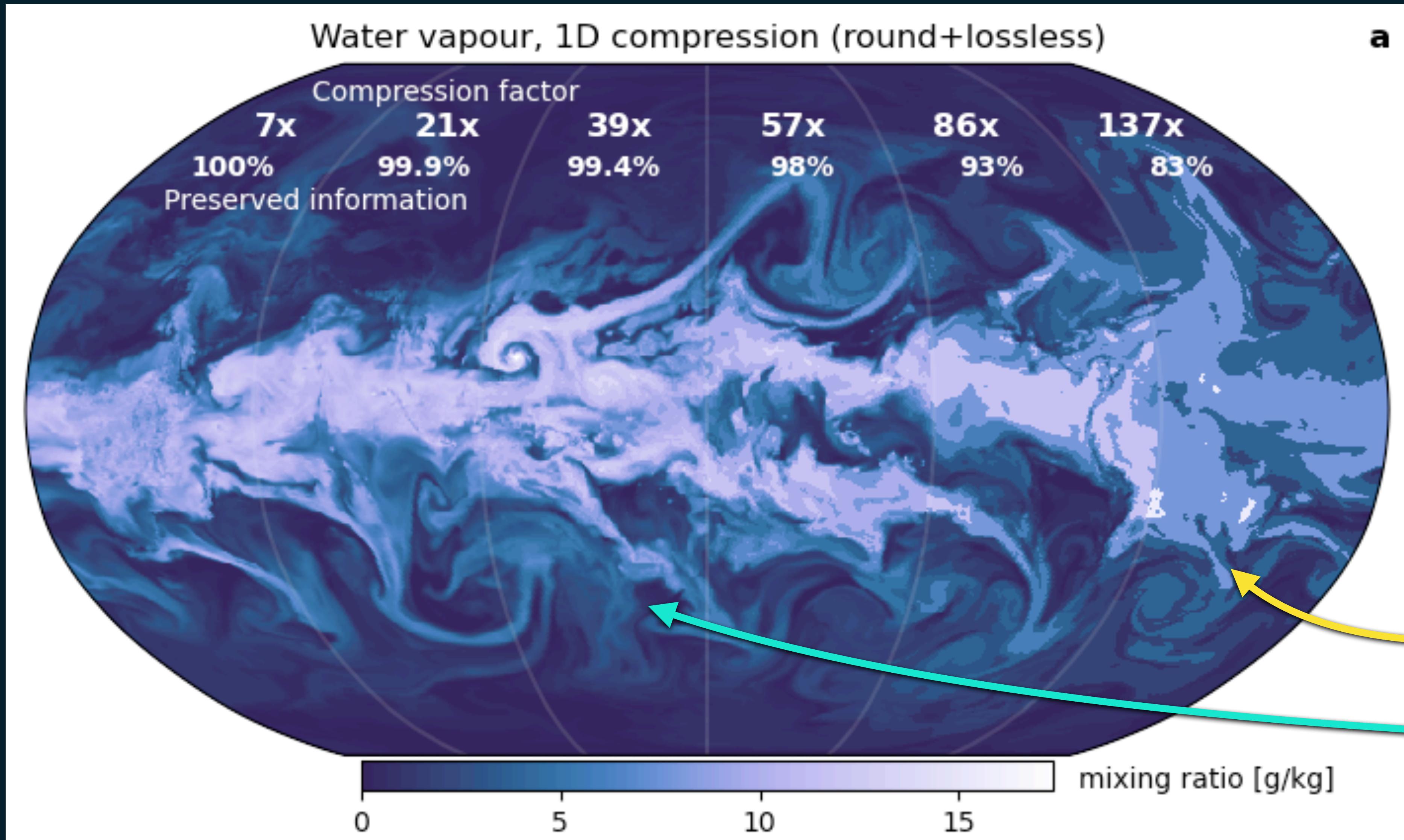


- Many unused exponent bits, lossless compression essential

- Information is often spread across many bit positions

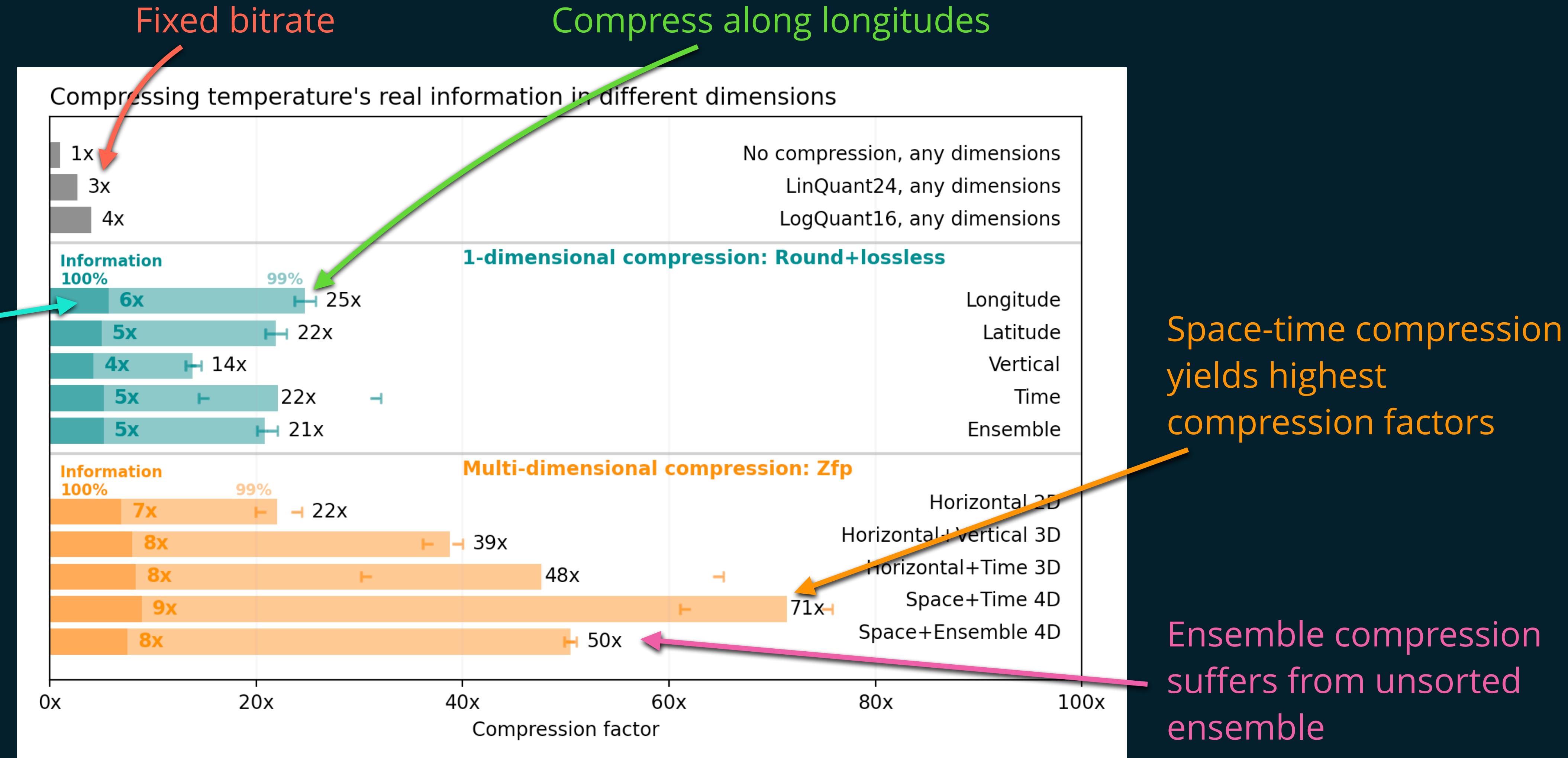
- Information is densely packed in the first bits

Information-preserving compression



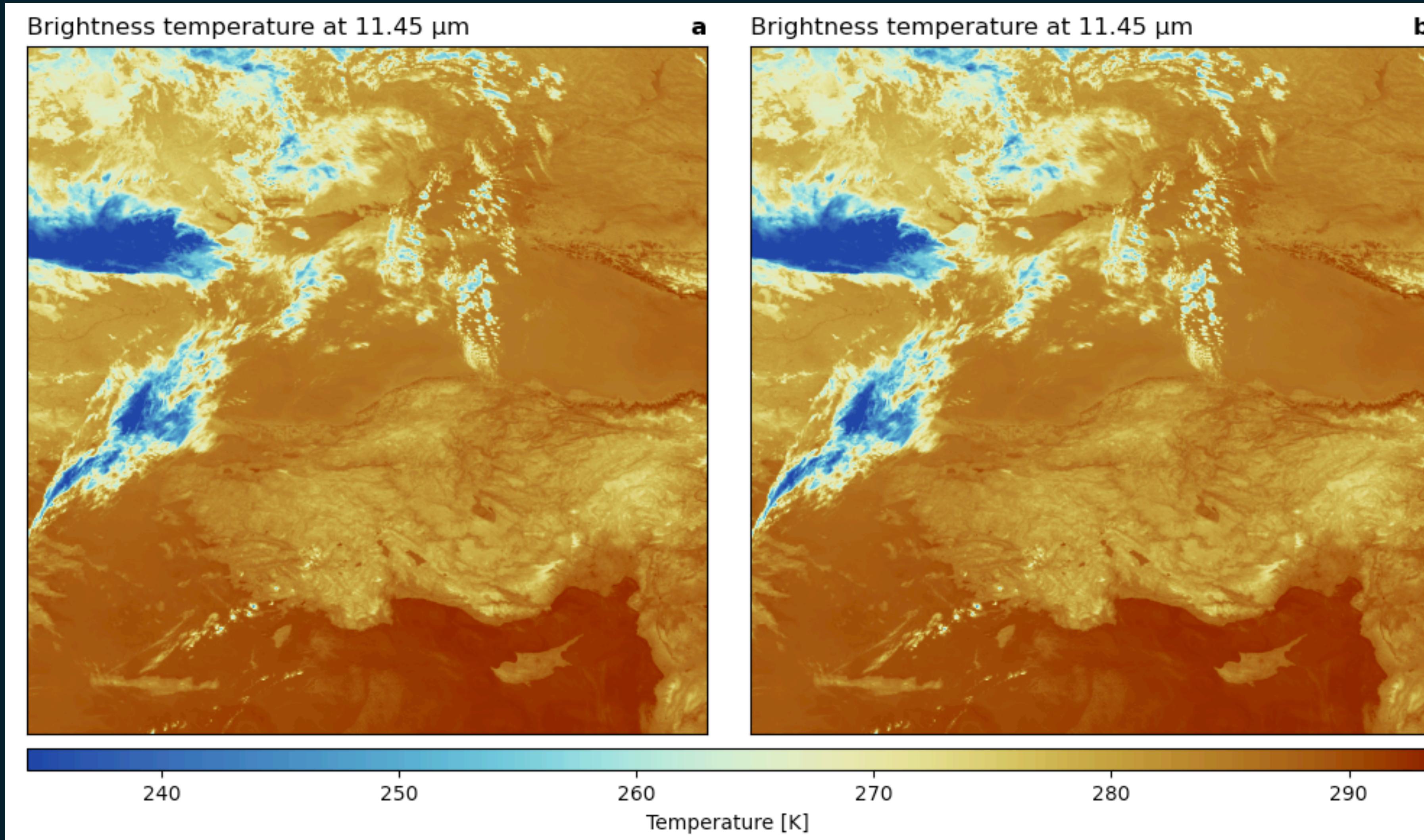
Multi-dimensional compression

Discarding 1%
of information
allows much
higher
compression



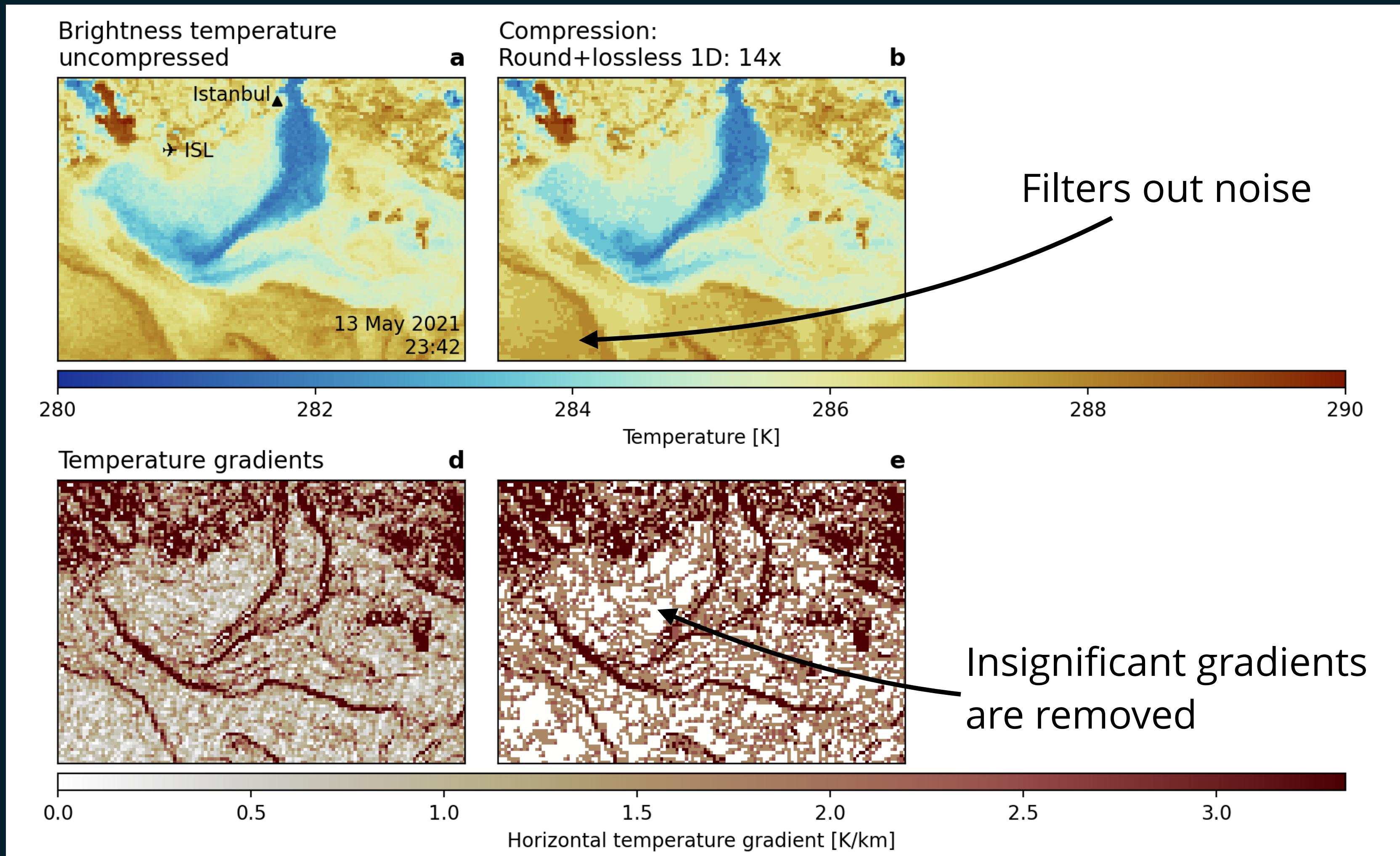
Information-preserving compression

One is 15-20x smaller than the other, which has been compressed?



Information-preserving compression

Filtering out insignificant gradients



Workflow

1. Determine the real bitwise information content

Round+lossless compression

2. Remove insignificant bits via rounding
3. Lossless compression

Workflow

1. Determine the real bitwise information content

Round+lossless compression

2. Remove insignificant bits via rounding

3. Lossless compression

Zstandard, ...

BitInformation.jl v0.5.1

BitInformation.jl
Numcodecs
nco

Workflow

1. Determine the real bitwise information content



Round+lossless compression

2. Remove insignificant bits via rounding



3. Lossless compression



BitInformation.jl v0.5.1

BitInformation.jl
Numcodecs
nco

- Information analysis only once offline
- Rounding is very fast
- Any lossless codec can be used
- Flexibility size/performance
- Works with any bit encoding
- Error bounds from floats
- Flexible chunking/granularity possible

Advantages

Workflow

1. Determine the real bitwise information content

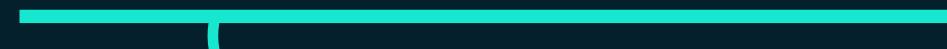


Round+lossless compression

2. Remove insignificant bits via rounding



3. Lossless compression



BitInformation.jl v0.5.1

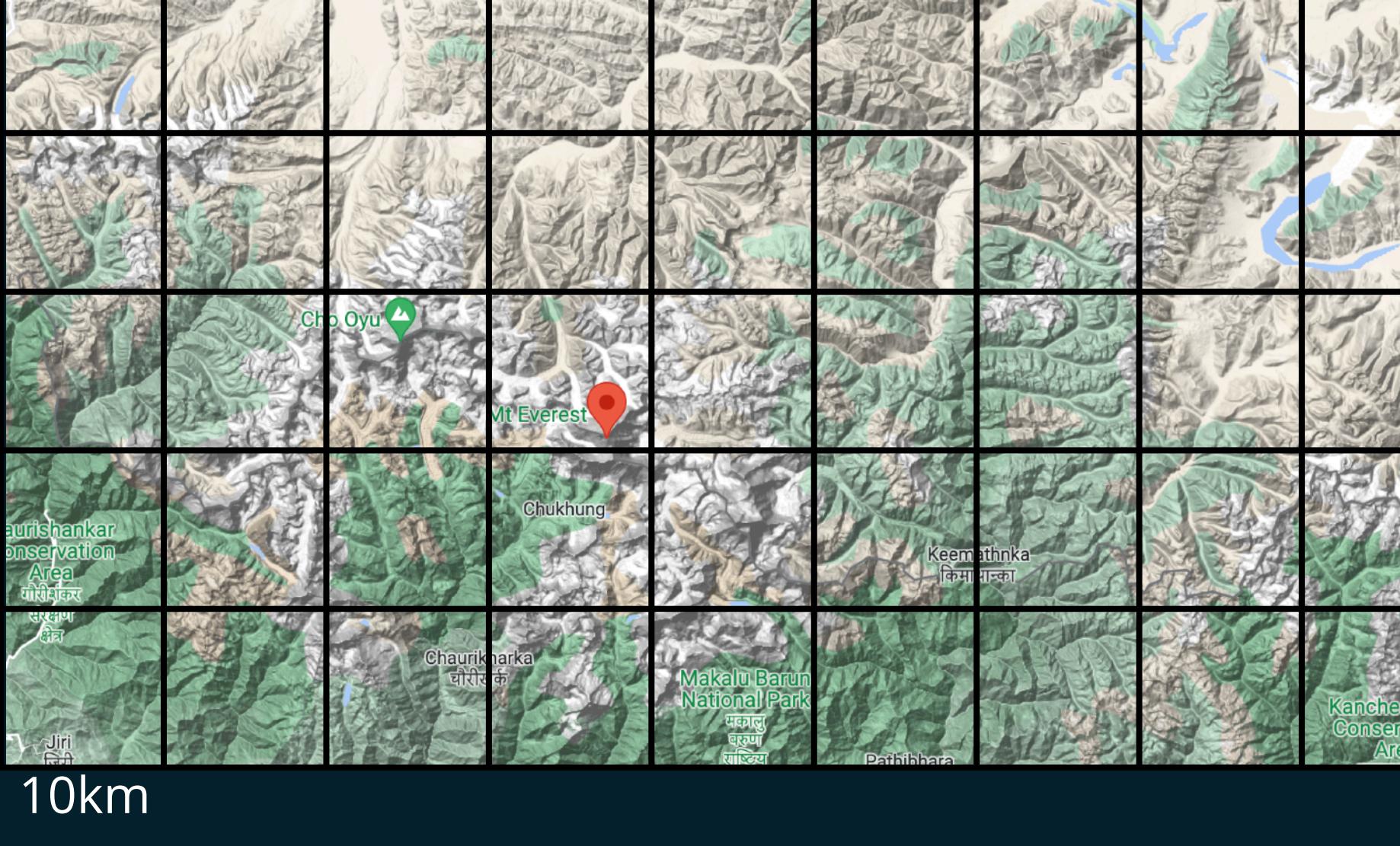
BitInformation.jl
Numcodecs
nco

- Information analysis only once offline
- Rounding is very fast
- Any lossless codec can be used
- Flexibility size/performance
- Works with any bit encoding
- Error bounds from floats
- Flexible chunking/granularity possible

Disadvantages

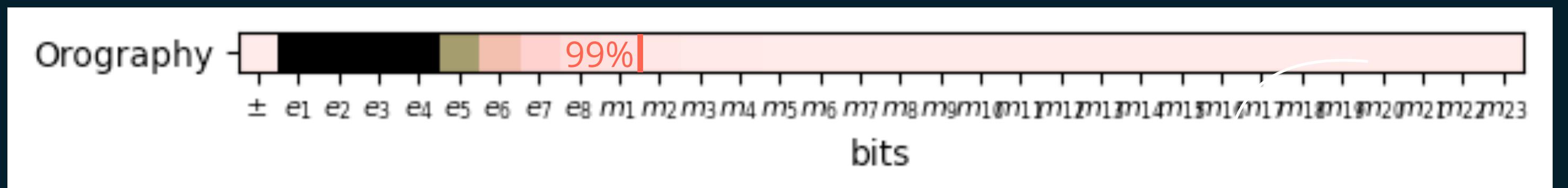
- Is real information well defined for any data?
- Mostly only 1D lossless compression

Orography



Julia workflow example

```
julia> using NetCDF, BitInformation  
julia> ncfile = NetCDF.open("orography.nc")  
julia> orog = ncfile.vars["orog"][:, :]  
julia> bi = bitinformation(orog)
```



```
julia> round(8848.86,1)  
8192.0
```

```
julia> round(60.0,1)  
64.0
```

Using NCO

```
> ncks -L 9 --baa=8 --ppc orog=1 orography.nc orography_round.nc
```

8MB

1.2MB

Reference implementation

BitInformation.jl is a package for bitwise information analysis and manipulation in Julia arrays. Based on counting the occurrences of bits in floats (or generally any bits type) across various dimensions, this package calculates quantities like the bitwise real information content, the mutual information, the joint entropy, and the conditional entropy.

The repository has 5 stars, 1 fork, and 1 pull request. It is licensed under MIT and is 100% compatible with Julia. The last commit was on April 1, 2023.

BitInformation.jl

CI passing docs dev DOI 10.5281/zenodo.4774191

BitInformation.jl is a package for bitwise information analysis and manipulation in Julia arrays. Based on counting the occurrences of bits in floats (or generally any bits type) across various dimensions, this package calculates quantities like the bitwise real information content, the mutual information, the joint entropy, and the conditional entropy.

Python workflow example

github.com/observingClouds/xbitinfo
xarray-wrapper around BitInformation.jl

A Spring & H Schulz, MPI

How to use

```
import xarray as xr
import xbitinfo as xb
ds = xr.tutorial.load_dataset(inpath)
bitinfo = xb.get_bitinformation(ds, dim="lon") # calling bitinformation.jl
keepbits = xb.get_keepbits(bitinfo, inflevel=0.99) # get number of mantissa bits
ds_bitrounded = xb.xr_bitround(ds, keepbits) # bitrounding keeping only kept bits
ds_bitrounded.to_compressed_netcdf(outpath) # save to netcdf with compressed
```

The screenshot shows the GitHub repository page for `xbitinfo`. The repository is owned by `observingClouds` and has a public status. The description states: "Python wrapper of BitInformation.jl to easily compress xarray datasets based on their information content". The repository has 6 stars and 0 forks. The language usage chart shows 97.9% Python and 2.1% Julia. The navigation bar includes tabs for Code (selected), Issues (7), Pull requests, Discussions, Actions, and Projects. A sidebar on the left lists README.md, xbitinfo.readthedocs.io, MIT license, and a pull request. The main content area features a large heading "xbitinfo: Retrieve information content and compress accordingly". Below it are status indicators for launch (binder), CI (passing), pre-commit.ci (passed), docs (passing), and pypi (package or version not found). A descriptive text block explains the purpose of the package: "This is an xarray -wrapper around BitInformation.jl to retrieve and apply bitrounding from within python. The package intends to present an easy pipeline to compress (climate) datasets based on the real information content."

observingClouds / `xbitinfo` Public

Python wrapper of BitInformation.jl to easily compress xarray datasets based on their information content

Languages

Python 97.9% Julia 2.1%

Starred

Watch

Code Issues 7 Pull requests Discussions Actions Projects

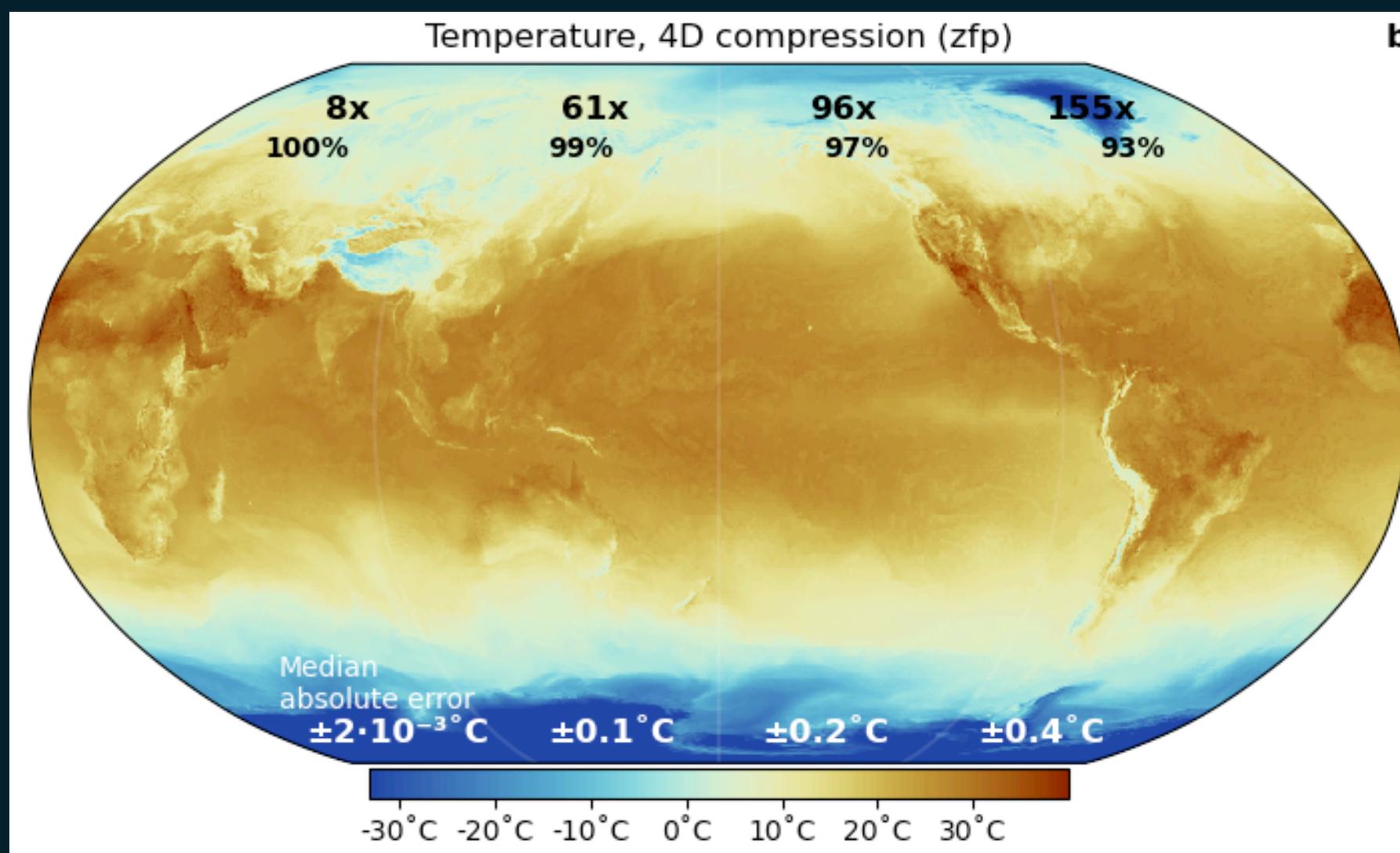
README.md

xbitinfo: Retrieve information content and compress accordingly

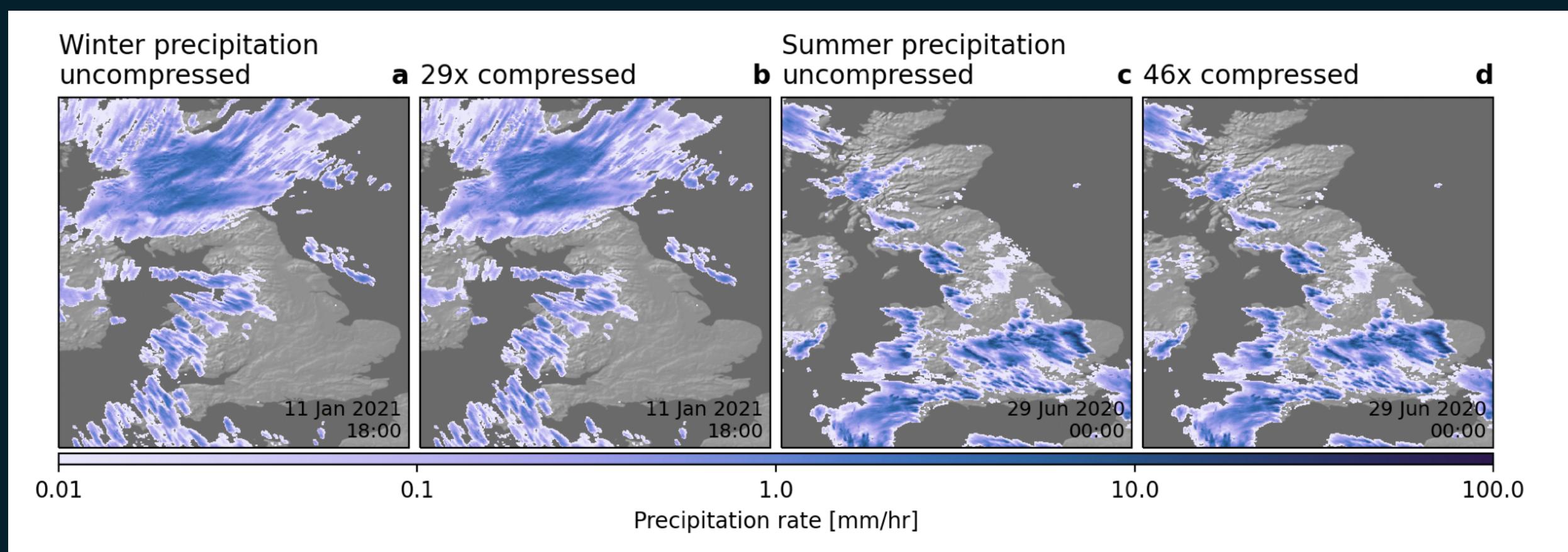
launch binder CI passing pre-commit.ci passed docs passing pypi package or version not found

This is an xarray -wrapper around BitInformation.jl to retrieve and apply bitrounding from within python. The package intends to present an easy pipeline to compress (climate) datasets based on the real information content.

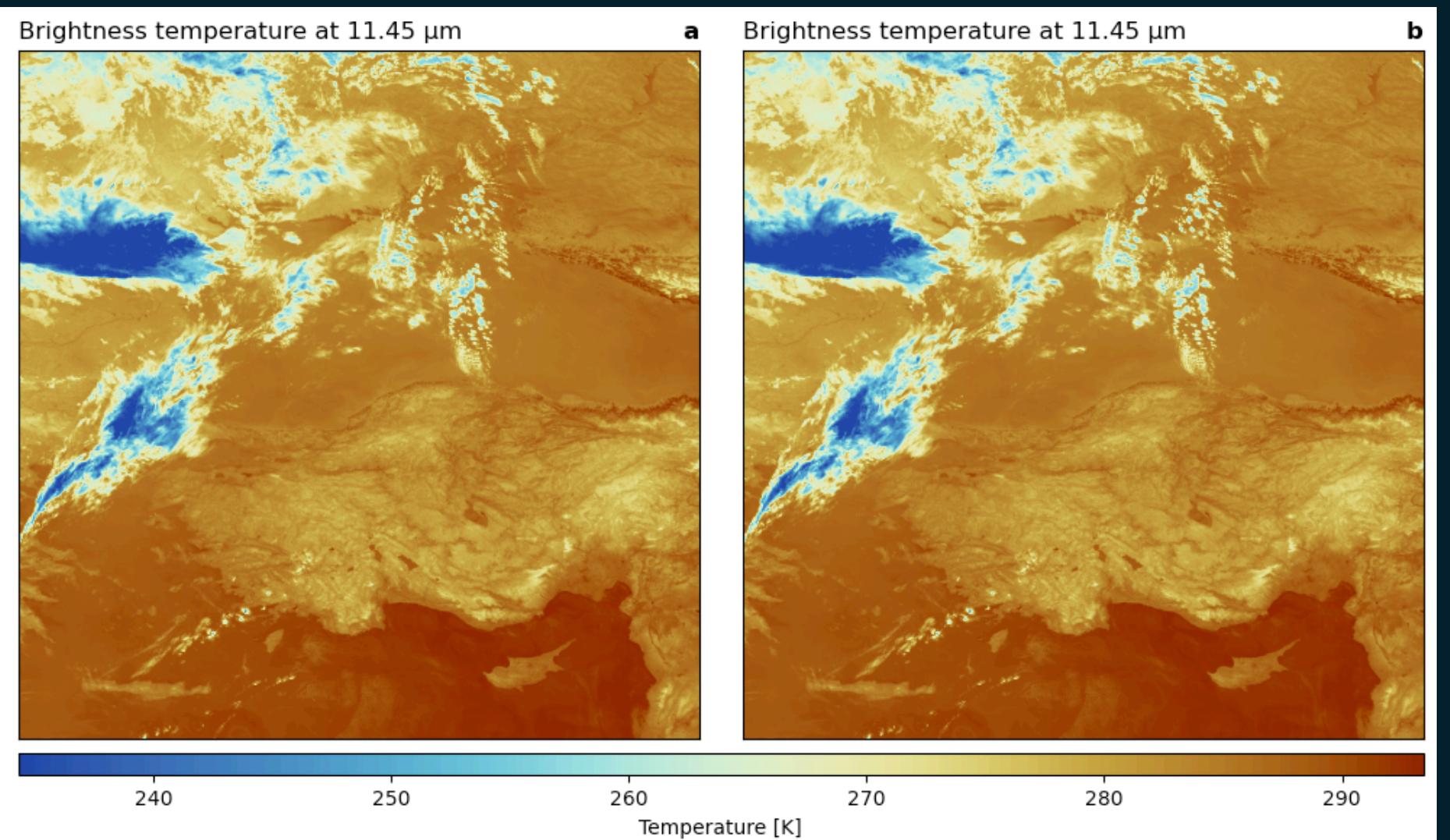
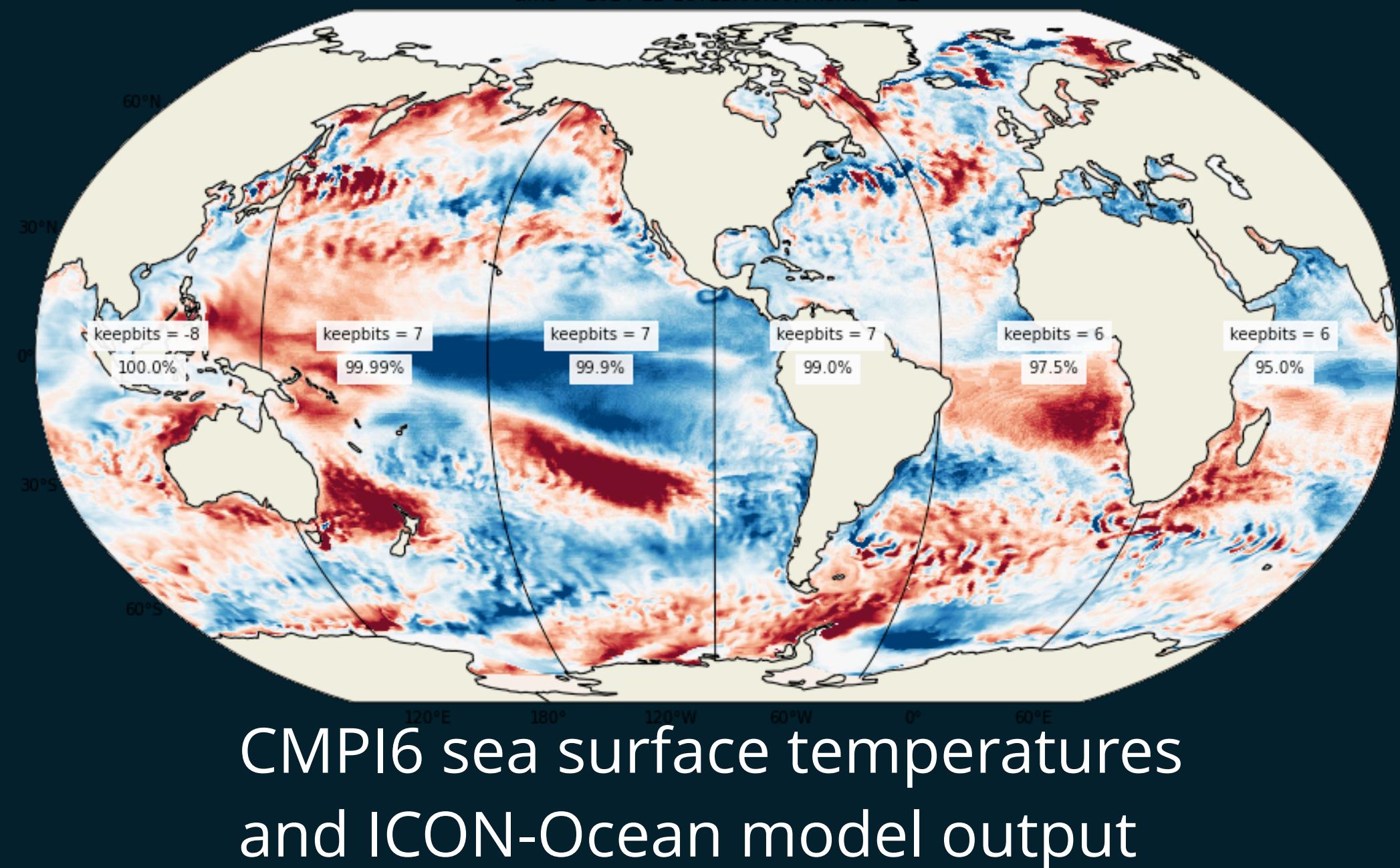
Application examples



ECMWF's ensemble temperature forecast



Precipitation from radar data



Brightness temperature from satellites