

hiopy

Optimizing Model Output for Analysis

Tobias Kölling

Parallel IO NHR Workshop on 07.05.2024

Parallel IO NHR Workshop on 07.05.2024

example: NextGEMS Simulations

- Model resolution: 5km
- Typical output size: **1 PB**
- Typical simulation time: **3 Weeks**

600 MB / s



Parallel IO NHR Workshop on 07.05.2024

high-resolution model output can be slow

Parallel IO NHR Workshop on 07.05.2024

the time it takes until the analysis plot is ready

- **understanding** the data
- **coding** the analysis
- **getting** the data

Useful output is
written once and
read **at least** once.

Parallel IO NHR Workshop on 07.05.2024

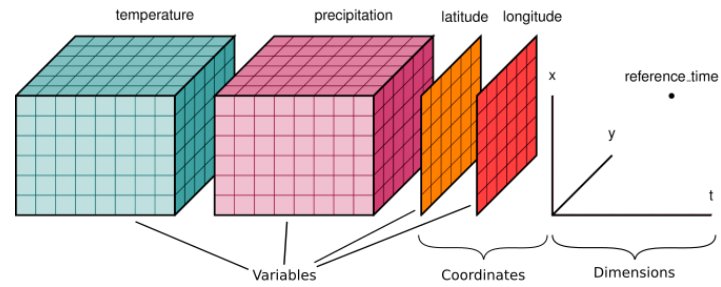
optimize output for analysis

(not write throughput)

Parallel IO NHR Workshop on 07.05.2024

Parallel IO NHR Workshop on 07.05.2024

(for this talk)



- n-dimensional variables
- shared dimensions
- coordinates
- attributes for metadata

- a single file
- a storage format
- shaped by storage & handling

```
$ ls *.nc
ngc2009_atm_mon_20200329T000000Z.nc
ngc2009_oce_2d_1h_inst_20200329T000000Z.nc
ngc2009_atm_pl_6h_inst_20200329T000000Z.nc
ngc2009_lnd_tl_6h_inst_20200329T000000Z.nc
ngc2009_lnd_2d_30min_inst_20200329T000000Z.nc
ngc2009_atm_2d_30min_inst_20200329T000000Z.nc
ngc2009_oce_0-200m_3h_inst_1_20210329T000000Z.nc
ngc2009_oce_0-200m_3h_inst_2_20210329T000000Z.nc
ngc2009_oce_moc_1d_mean_20210329T000000Z.nc
ngc2009_oce_2d_1d_mean_20210329T000000Z.nc
ngc2009_oce_ml_1d_mean_20210329T000000Z.nc
ngc2009_oce_2d_1h_mean_20210329T000000Z.nc
...
$ ls *.nc | wc -l
12695
```

















Parallel IO NHR Workshop on 07.05.2024

- provides an easy-to-understand overview
- forces consistency across output
- cutting things is easier than glueing things



























```
1 ds = cat.ICON.ngc4008.to_dask()
```

► Dimensions: (time: 10958, depth_half: 73, cell: 49152, level_full: 90, crs: 1, depth_full: 72, soil_depth_water_level: 5, level_half: 91, soil_depth_energy_level: 5)

▼ Coordinates:

crs	(crs)	float32	nan	 
depth_full	(depth_full)	float32	1.0 3.1 ... 5.546e+03 ...	 
depth_half	(depth_half)	float32	0.0 2.0 4.2 ... 5.681e...	 
level_full	(level_full)	int32	1 2 3 4 5 6 7 ... 85 8...	 
level_half	(level_half)	int32	1 2 3 4 5 6 7 ... 86 8...	 
soil_depth_en...	(soil_depth_energy_level)	float32	0.0325 0.192 0.7755...	 
soil_depth_wa...	(soil_depth_water_level)	float32	0.0325 0.192 0.7755...	 
time	(time)	datetime64[ns]	2020-01-02 ... 2050...	 

▼ Data variables:

A_tracer_v_to	(time, depth_half, cell)	float32	...	 
FrshFlux_IceSalt	(time, cell)	float32	...	 
FrshFlux_TotalIce	(time, cell)	float32	...	 
Qbot	(time, cell)	float32	...	 
Qtop	(time, cell)	float32	...	 
Wind_Speed_1...	(time, cell)	float32	...	 
atmos_fluxes_F...	(time, cell)	float32	...	 
atmos_fluxes_F...	(time, cell)	float32	...	 
atmos_fluxes_F...	(time, cell)	float32	...	 
atmos_fluxes_F...	(time, cell)	float32	...	 
atmos_fluxes_...	(time, cell)	float32	...	 
atmos_fluxes_...	(time, cell)	float32	...	 
atmos fluxes ...	(time, cell)	float32	...	 

Parallel IO NHR Workshop on 07.05.2024

Parallel IO NHR Workshop on 07.05.2024

Parallel IO NHR Workshop on 07.05.2024

Parallel IO NHR Workshop on 07.05.2024

Parallel IO NHR Workshop on 07.05.2024

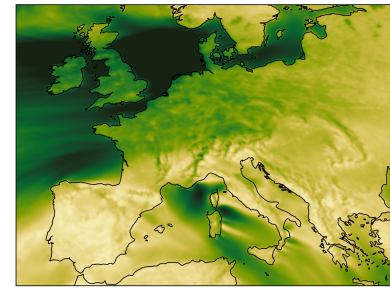
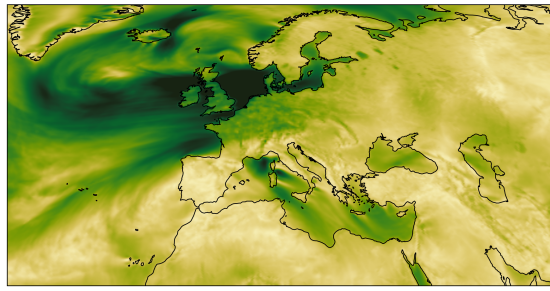
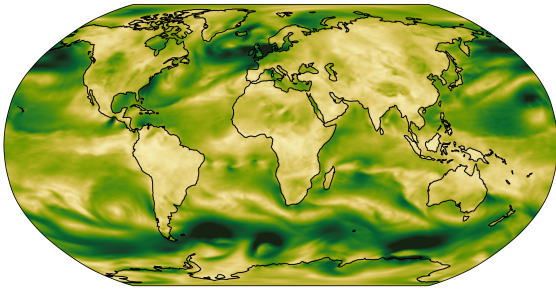
Parallel IO NHR Workshop on 07.05.2024

Parallel IO NHR Workshop on 07.05.2024

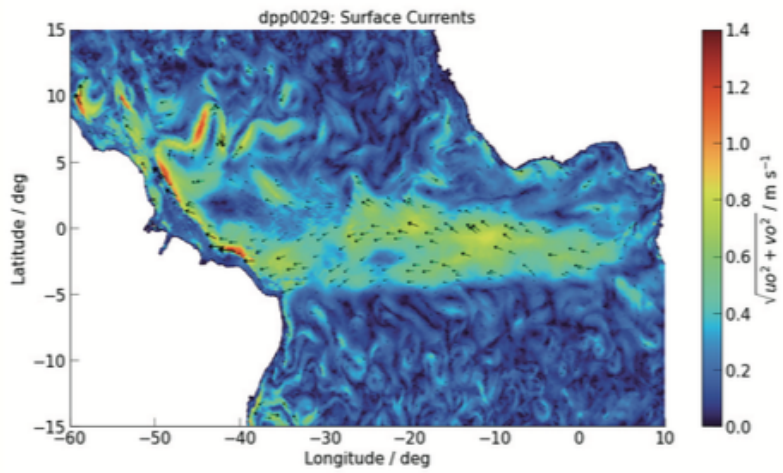
Grid	Cells	Screen	Pixels
1° by 1°	0.06M	VGA	0.3M
10 km	5.1M	Full HD	2.1M
5 km	20M	MacBook 13'	4.1M
1 km	510M	4K	8.8M
200 m	12750M	8K	35.4M

It's **impossible** to look at the entire globe in full resolution.

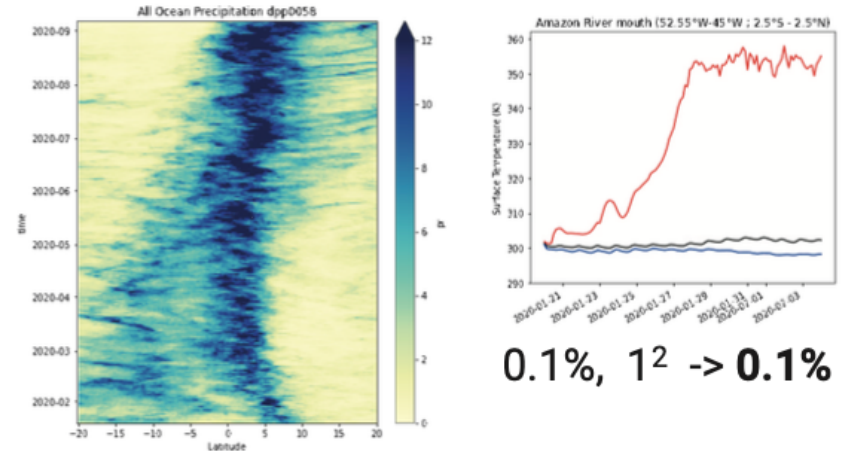
Parallel IO NHR Workshop on 07.05.2024



Parallel IO NHR Workshop on 07.05.2024



5%, $1/4^2 \rightarrow 0.3\%$



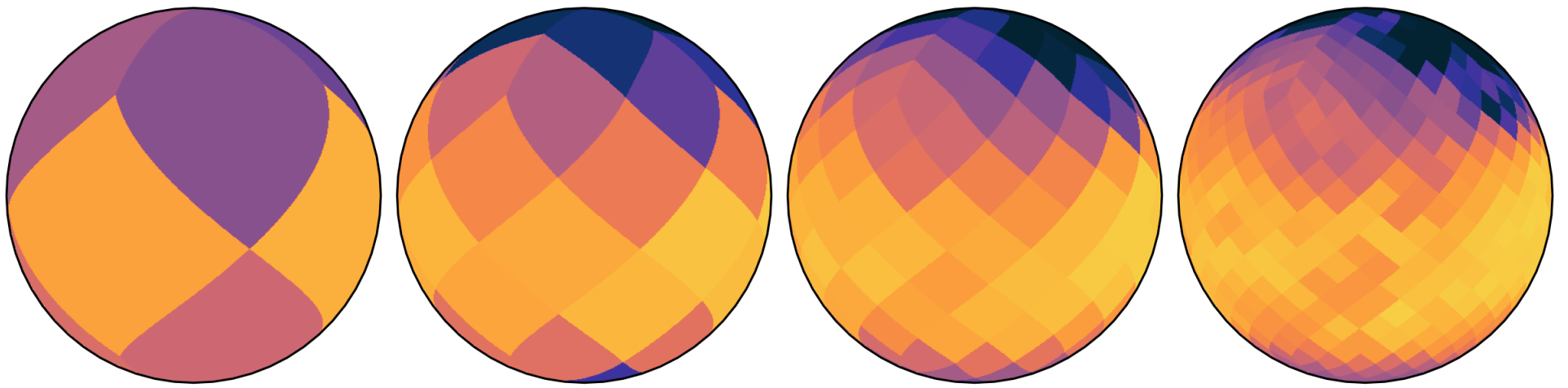
0.1%, $1^2 \rightarrow 0.1\%$

34%, $1/4^2 \rightarrow 2\%$

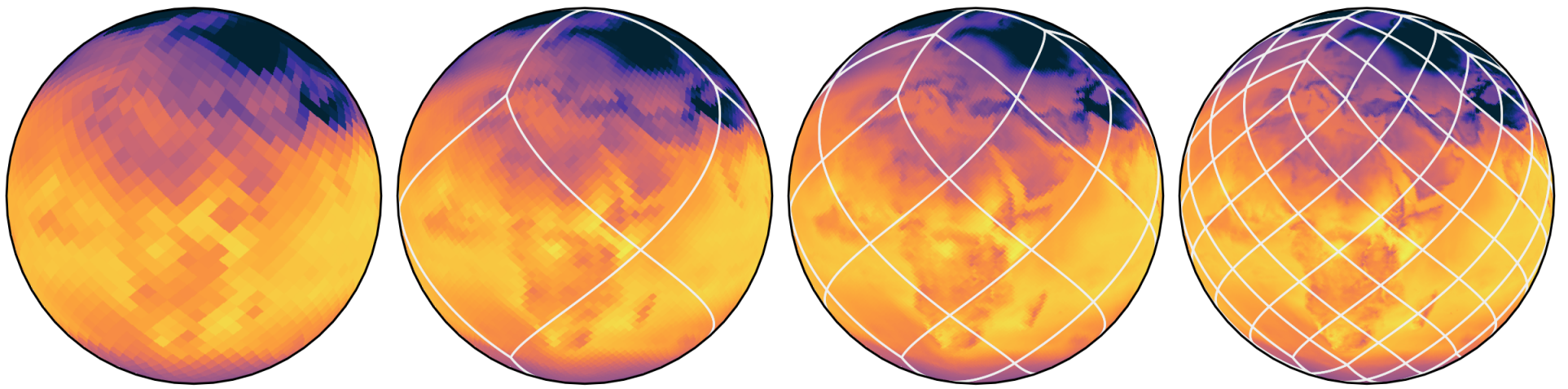
fraction of globe covered
 required resolution
 useful data

Analysis scripts are **forced** to load way too much data.

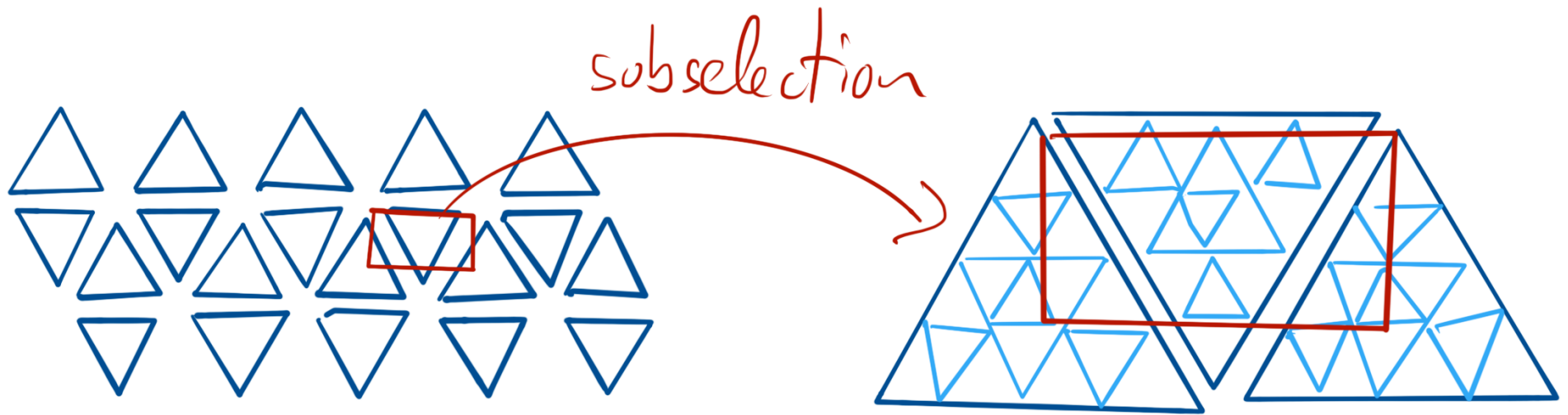
Parallel IO NHR Workshop on 07.05.2024



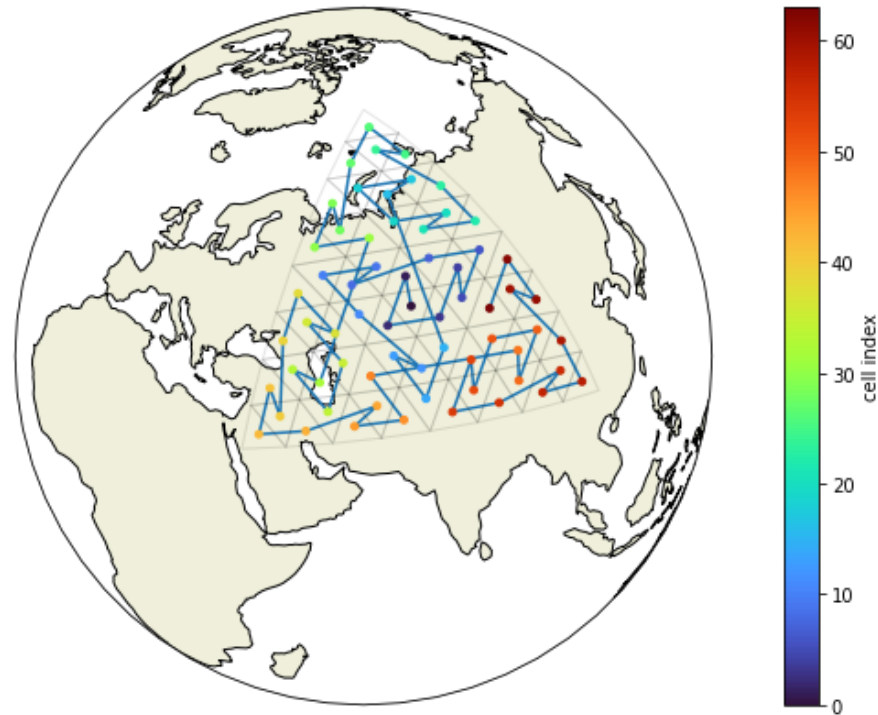
Parallel IO NHR Workshop on 07.05.2024



Parallel IO NHR Workshop on 07.05.2024



Parallel IO NHR Workshop on 07.05.2024

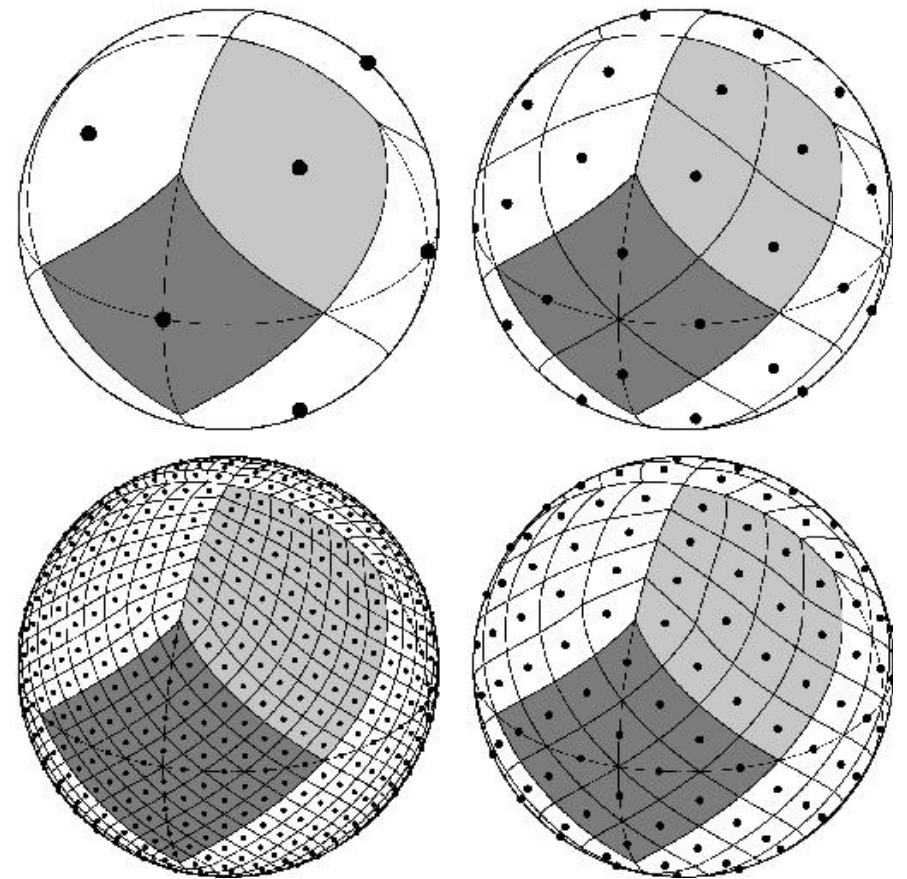


- spatially close \approx close in index space
- recursive pattern (repeats every 4^n)

Parallel IO NHR Workshop on 07.05.2024

Parallel IO NHR Workshop on 07.05.2024

- hierarchical
- equal area
- iso latitude



3 11_2	1 01_2
2 10_2	0 00_2

Prograde 

 Degrade

15 $11 11_2$	13 $11 01_2$	7 $01 11_2$	5 $01 01_2$
14 $11 10_2$	12 $11 00_2$	6 $01 10_2$	4 $01 00_2$
11 $10 11_2$	9 $10 01_2$	3 $00 11_2$	1 $00 01_2$
10 $10 10_2$	8 $10 00_2$	2 $00 10_2$	0 $00 00_2$

aggregation can be exact and is easy to implement

scale analysis with screen size
(instead of with model size)

Parallel IO NHR Workshop on 07.05.2024

Parallel IO NHR Workshop on 07.05.2024

Select ICON model output at all dropsonde locations during EUREC4A field campaign:

```
1 sonde_pix = healpy.ang2pix(  
2     icon.crs.healpix_nside, joanne.flight_lon, joanne.flight_lat,  
3     lonlat=True, nest=True  
4 )  
5  
6 icon_sondes = (  
7     icon[["ua", "va", "ta", "hus"]]  
8     .sel(time=joanne.launch_time, method="nearest")  
9     .isel(cell=sonde_pix)  
10    .compute()  
11 )
```

(55 sec, 1GB, single thread, full code at [easy.gems](https://easygems.org))

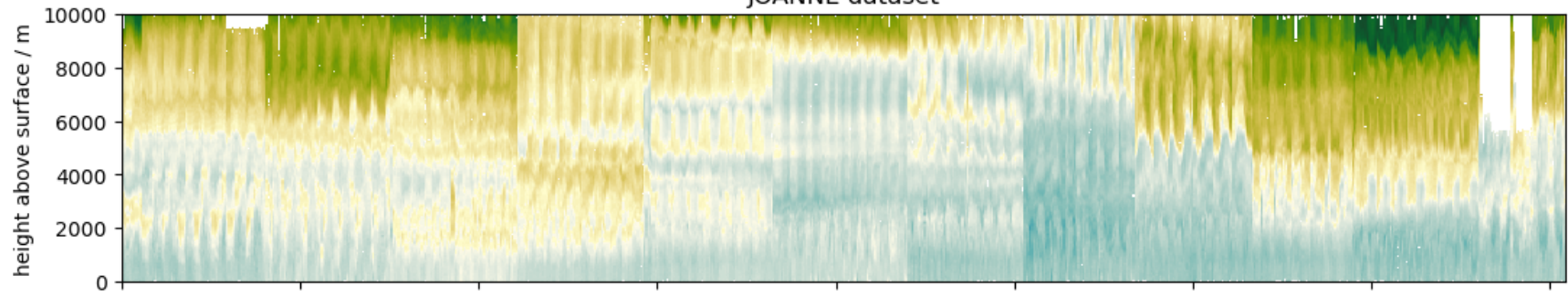
Parallel IO NHR Workshop on 07.05.2024

Select ICON model output at all dropsonde locations during EUREC4A field campaign:

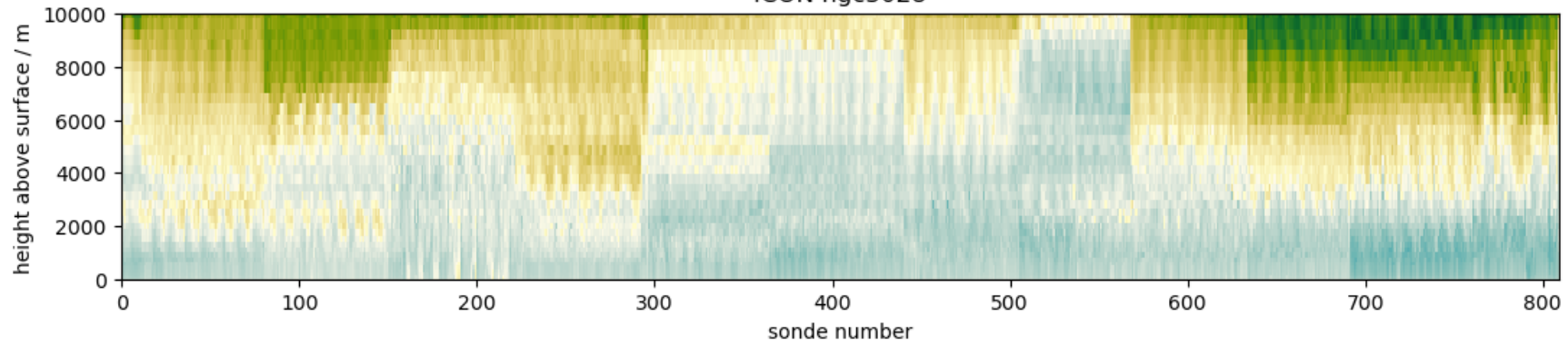
```
1 icon_sondes = (  
2     icon[["ua", "va", "ta", "hus"]]  
3     .sel(time=joanne.launch_time, method="nearest")  
4     .dggs.sel_latlon(joanne.flight_lat, joanne.flight_lon)  
5     .compute()  
6 )
```

with XDGGS (?)

u wind @ dropsonde locations in EUREC4A
JOANNE dataset

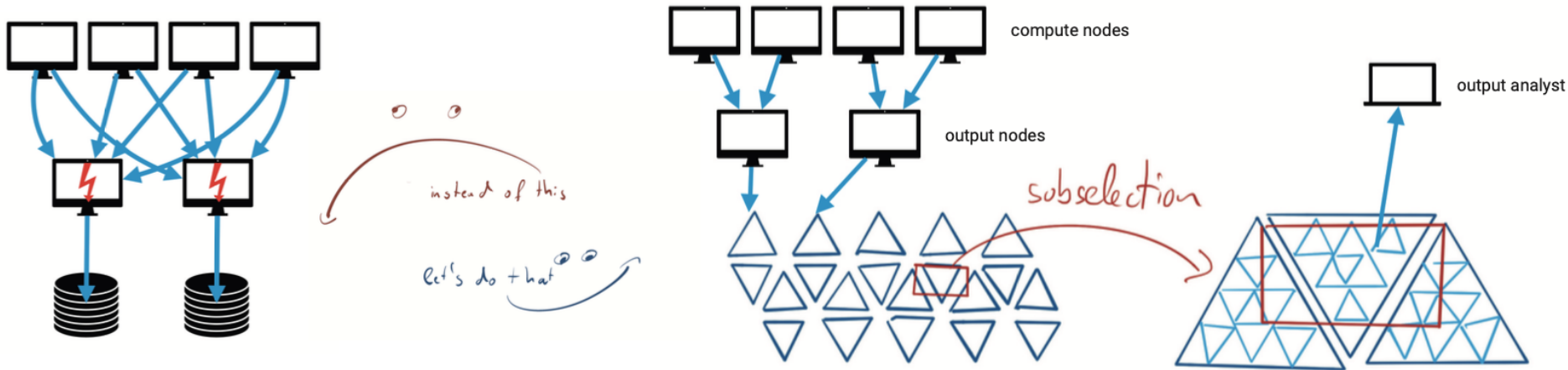


ICON ngc3028

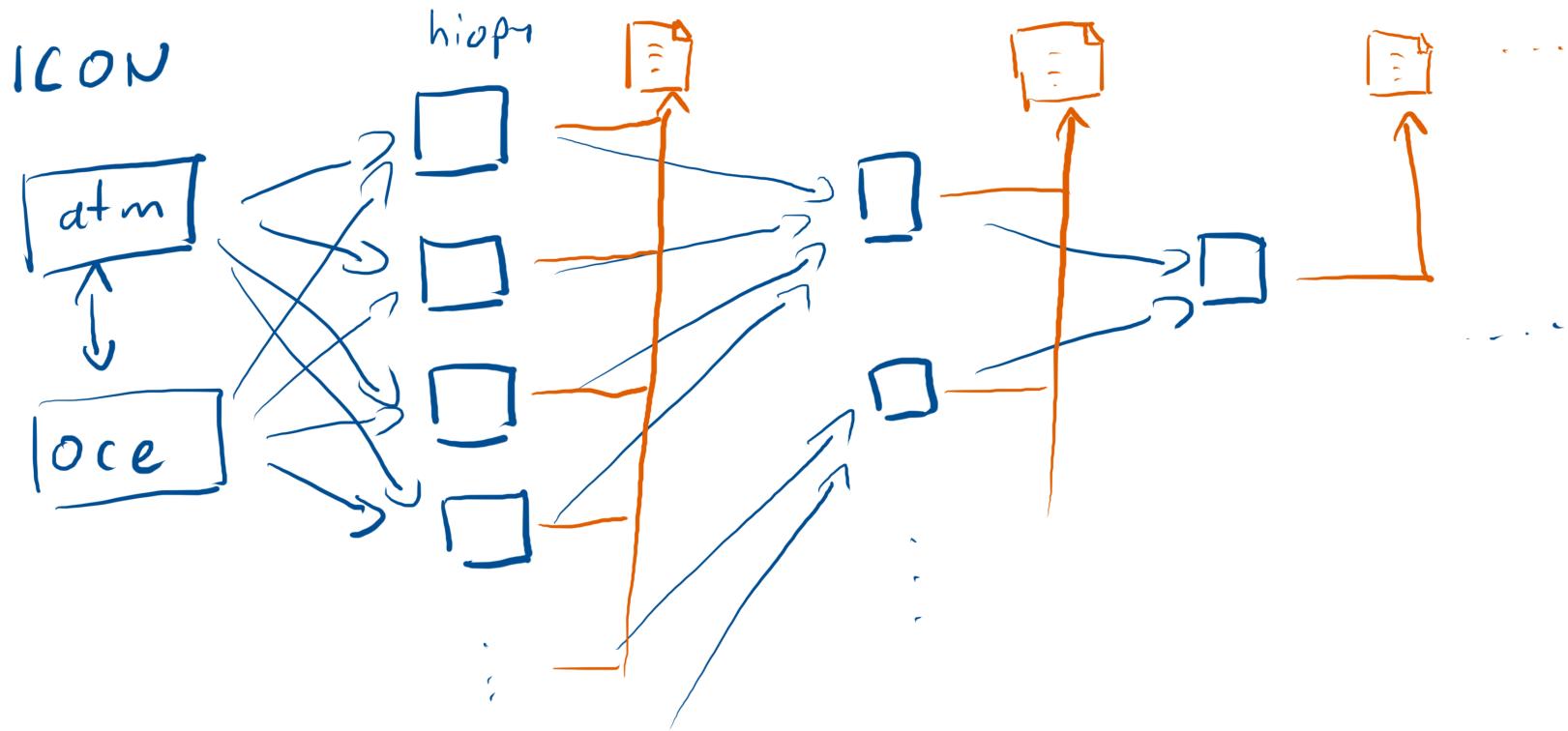


Parallel IO NHR Workshop on 07.05.2024

Parallel IO NHR Workshop on 07.05.2024



```
1 yac = YAC(xml, xsd)
2 sgd = HealPixSubgridDefinition(2**10, nchunks, ichunk)
3 dataset = zarr.open_consolidated(output_folder, mode="r+")
4 comp_id = yac.def_comp("healpix_io")
5 point_id, grid = make_yac_grid(sgd)
6 fields = {varname: Field.create(varname, comp_id, point_id)
7           for varname in varnames}
8 put_fields = ...
9 yac.search()
10 steps = compute_nsteps(yac, fields)
11 for i in range(steps):
12     for varname, field in fields.items():
13         buffer = field.get()
14         put_fields[varname].put(coarsen(buffer))
15         dataset[varname][i,sgd.cell_chunk_slice] = buffer
```



Parallel IO NHR Workshop on 07.05.2024

healpix resolution

$$n_{\text{side}} = 2^{\text{zoom}}$$



10

3

8

...



0

zoom →

temporal resolution

PT30M

30 minutes



...

ISO duration

PT3H

3 hours



...

P1D

1 day



...

time ↓

each box is a complete dataset
but all look (mostly) the same



- output process is coupled to the running model
- writes entire hierarchy at once
- dataset is accessible as soon as the model **starts**

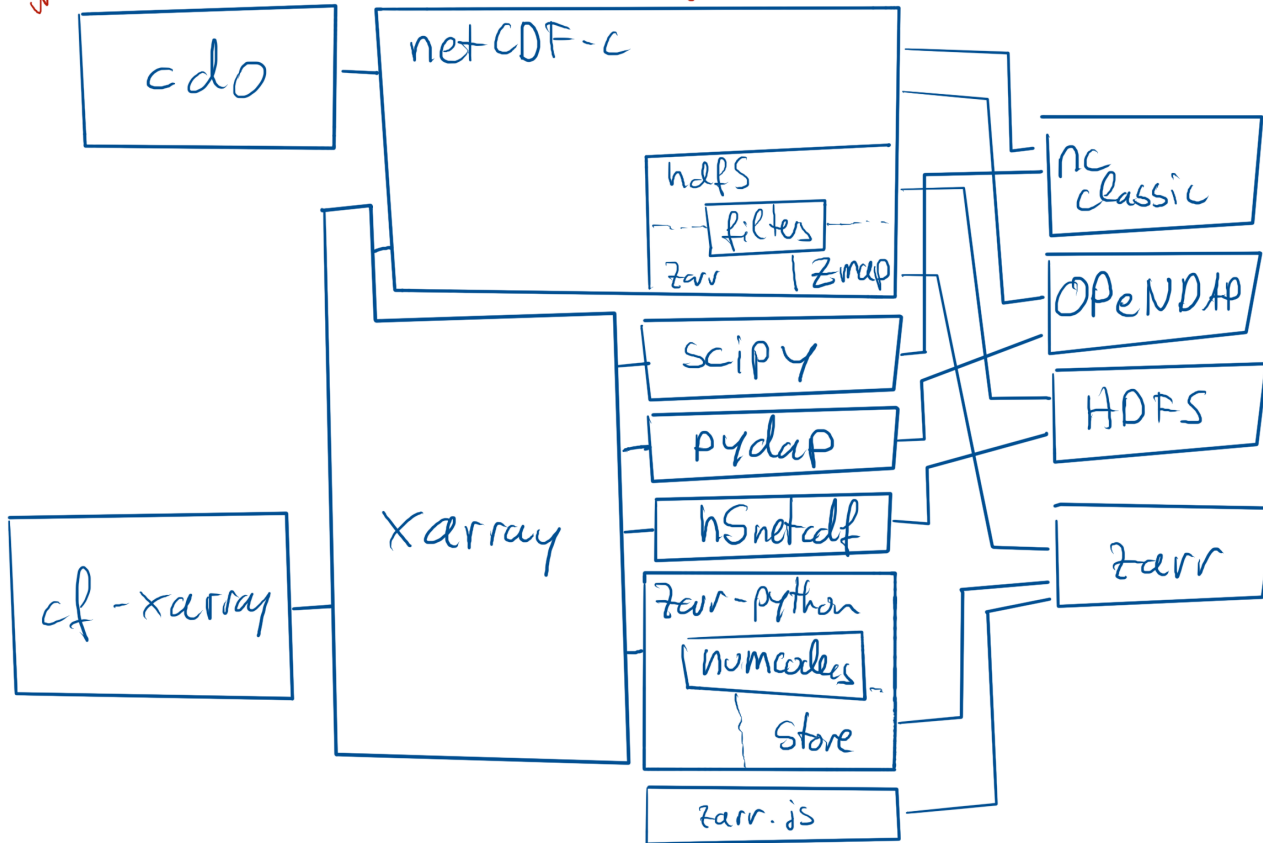
Parallel IO NHR Workshop on 07.05.2024

physical meaning

dataset

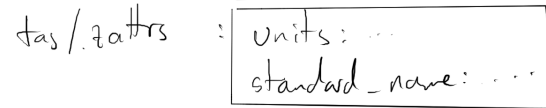
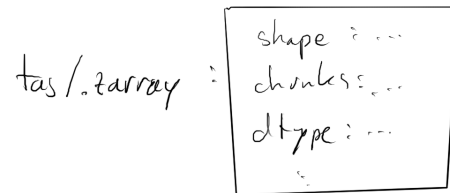
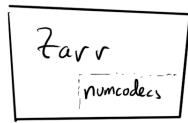
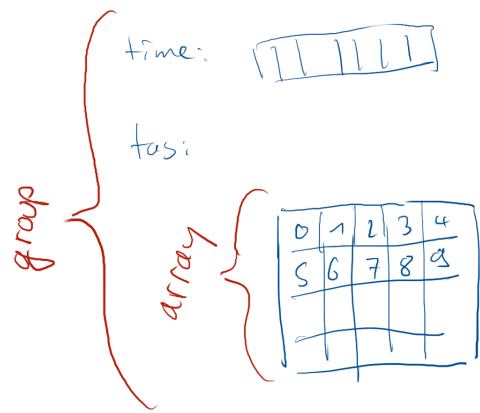
arrays & attributes

storage



no compression
 no chunking
 single file
 ← subsets - over - HTTP

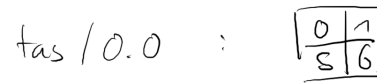
single file (mostly)
 flexible storage backend
 chunking compression



↑

"zmap"

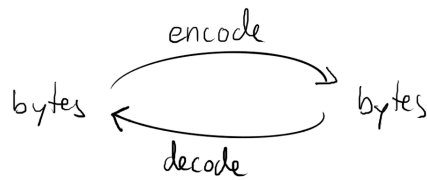
in netCDF-c



...

...

numcodecs / filter / etc...



lz4

zstd

blosc

tlib

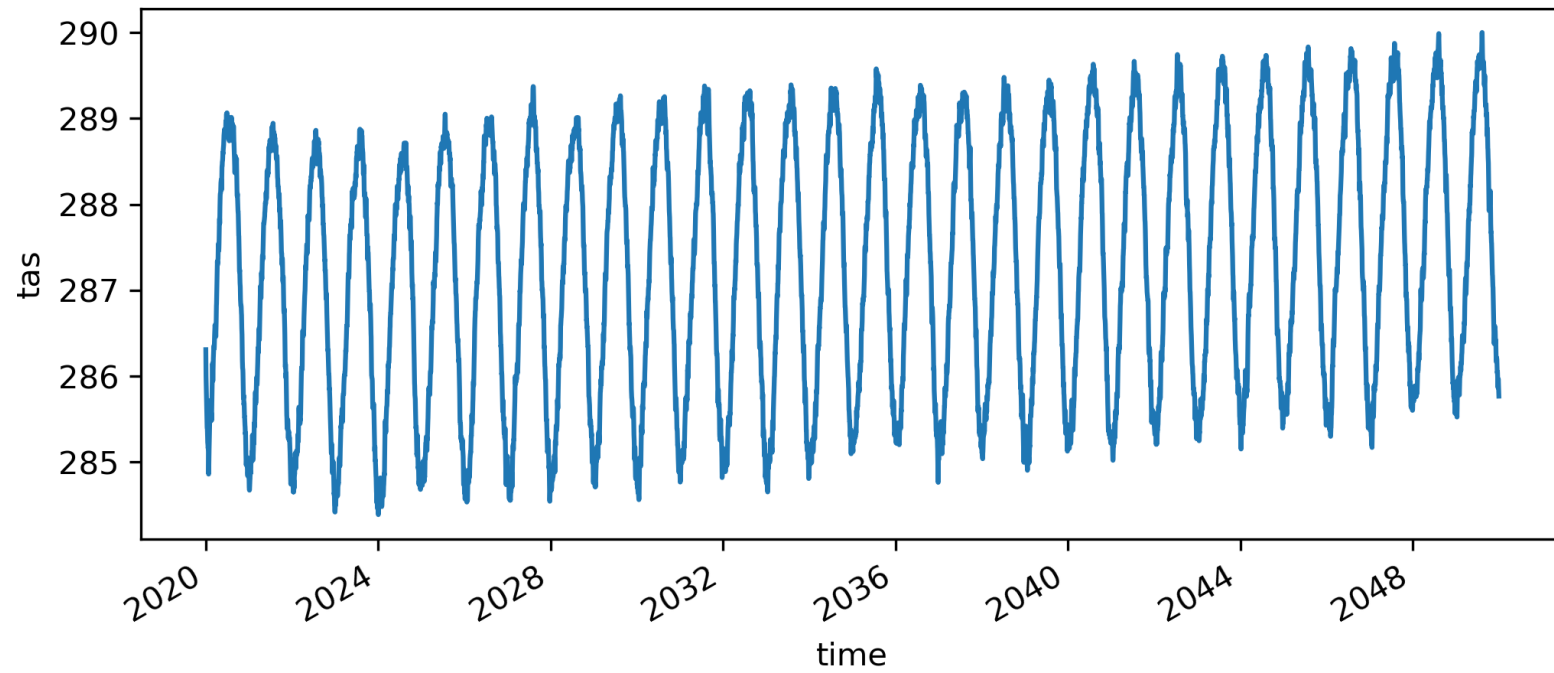
...

- Filesystem
- GCS, S3, SWIFT
- MySQL
- SQLite
- Zipfile
- reference
- Memory
- SLURM
- IPFS / IPLD
- Sharded
- ...

- can **improve** performance (less I/O)
- **reduces** storage cost
- requires good chunking

Parallel IO NHR Workshop on 07.05.2024

```
1 cat.ICON.ngc4008(time="P1D", zoom="0").to_dask().tas.mean("cell").p
```

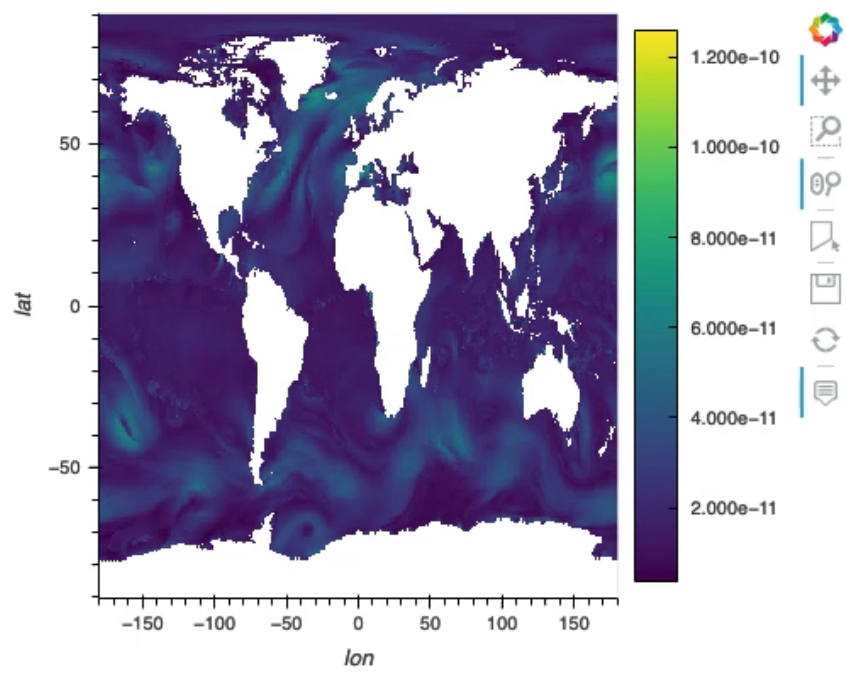


(100ms, 250MB, single thread)

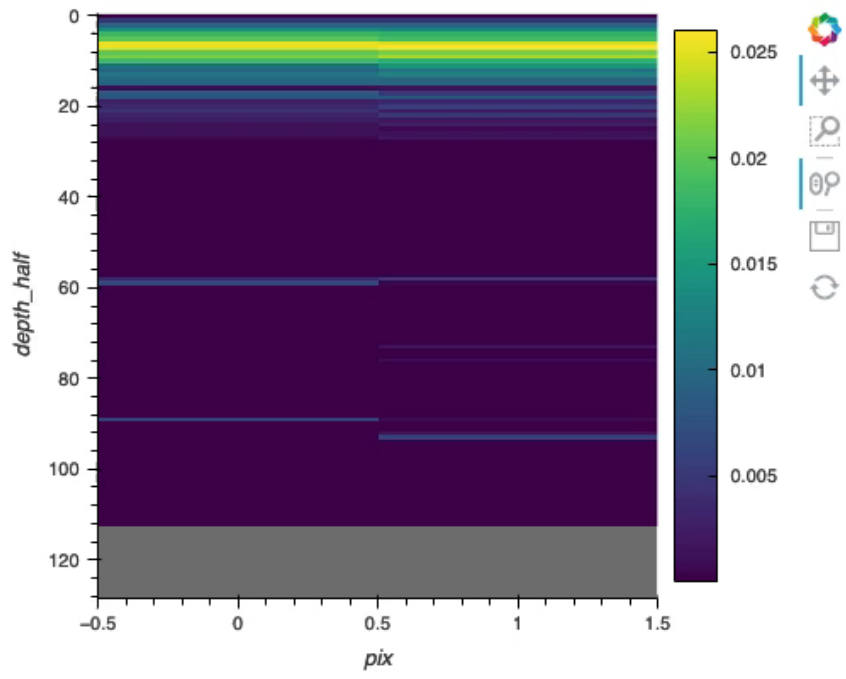
Time: 0



a_tracer_v_to



a_tracer_v_to



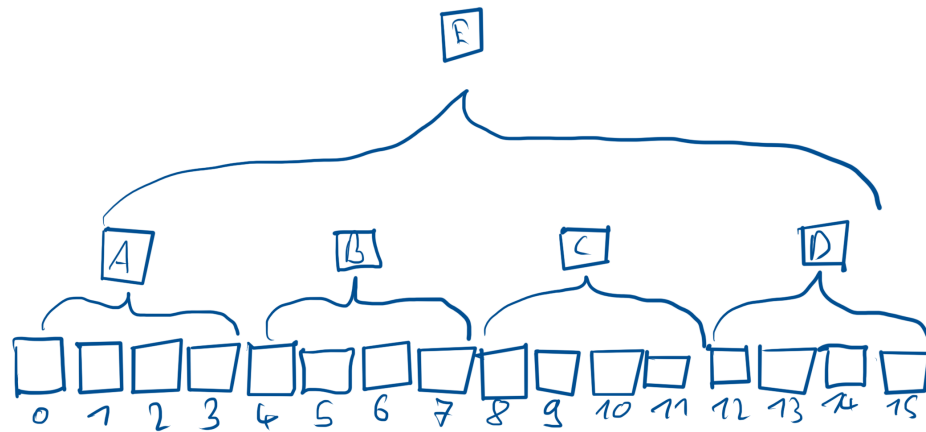
Parallel IO NHR Workshop on 07.05.2024

Output tested on multiple $\mathcal{O}(\text{PB})$ -scale model runs, 100+ users:

- remarkably little issues raised
- very positive general feedback
- enabled diagnostics which seemed impossible before

- Large datasets simplify access
- Chunking, hierarchies and compression play best together
- Healpix supports this ↗
- [hiopy](#) can build such datasets
- It's about data structures, not so much the tooling

Parallel IO NHR Workshop on 07.05.2024



mean = $\frac{E - 0 - D + 12}{12}$

↳ loads instead of 12 → 3 times less reads

1 point

4 points

16 points

27 points

31% increase