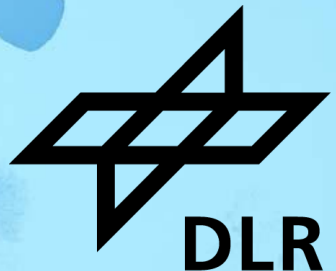


# The Helmholtz Analytics Toolkit (Heat) and its role in the landscape of massively- parallel scientific Python

C. Comito<sup>1</sup>, J. Gutiérrez Hermosillo Muriedas<sup>2</sup>, M. Götz<sup>2</sup>, B. Hagemeyer<sup>1</sup>,  
F. Hoppe<sup>3</sup>, P. Knechtges<sup>3</sup>, K. Krajsek<sup>1</sup>, A. Rüttgers<sup>3</sup>, A. Streit<sup>2</sup>, M. Tarnawa<sup>1</sup>

<sup>1</sup>Jülich Research Centre (FZJ), Institute for Advanced Simulation, Jülich Supercomputing Centre (JSC) <sup>2</sup>Karlsruhe Institute for Technology (KIT),  
Steinbuch Centre for Computing (SCC) <sup>3</sup>German Aerospace Center (DLR), Institute for Software Technology, High-Performance Computing



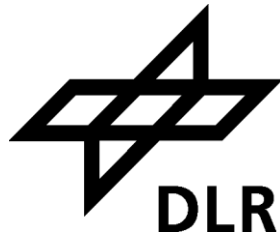
```
for i in range(1, self.n_clusters):
    distances = ht.spatial.distance.cdist(x, centroids, quadratic_expansion=True)
    D2 = distances.min(axis=1)
    D2.resplit_(axis=None)
    prob = D2 / D2.sum()
    random_position = ht.random.rand().item()
    sample = 0
    sum = 0
    for j in range(len(prob)):
        if sum > random_position:
            break
        sum += prob[j].item()
        sample = j
```

# What is Heat and what can I do with it?

# The Helmholtz Analytics Toolkit



- Free and open-source (MIT) Python library for **memory-distributed high-performance data analysis and machine learning** on CPU/GPU-clusters
- allows **scientists and engineers without an extensive background in HPC** to tackle data-intensive problems that go beyond the capabilities of a workstation
- Developed since 2018 at



P. Knechtges  
F. von der Lehr  
A. Rüttgers  
M. Siggel



C. Comito  
B. Hagemeier  
K. Krajsek  
M. Tarnawa



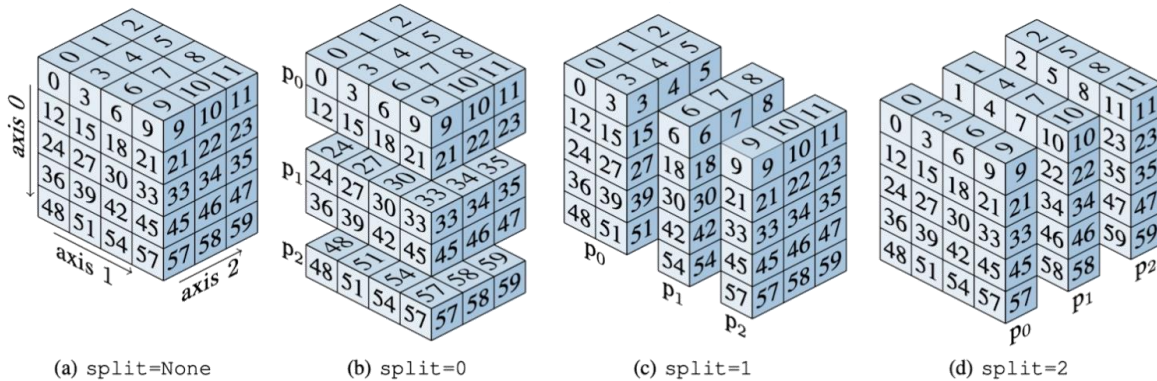
D. Coquelin  
C. Debus  
M. Götz  
J. P. G. Hermosillo Muriedas

+ external *open-source contributors* on github (e.g. GSoC 2022) + *student workers*

# How it works...

Vendor-independent GPU-support  
(Nvidia + AMD) inherited from PyTorch!

- distributed N-dimensional array („DNDarray“) as basic data type



- composed of „local“ torch-Tensors on each MPI-process
- „MPI-glue“ using `mpi4py`
- `numpy`-like API




Clustering  
(k-means, Spectral)

Deep Learning  
(DASO, [3])


Regression  
(QR, LASSO)

Naive Bayes

and more...



- Tensors
- Linear Algebra
- Automatic Differentiation
- NumPy-like interface
- GPU-support (vendor-independent)
- multithreading

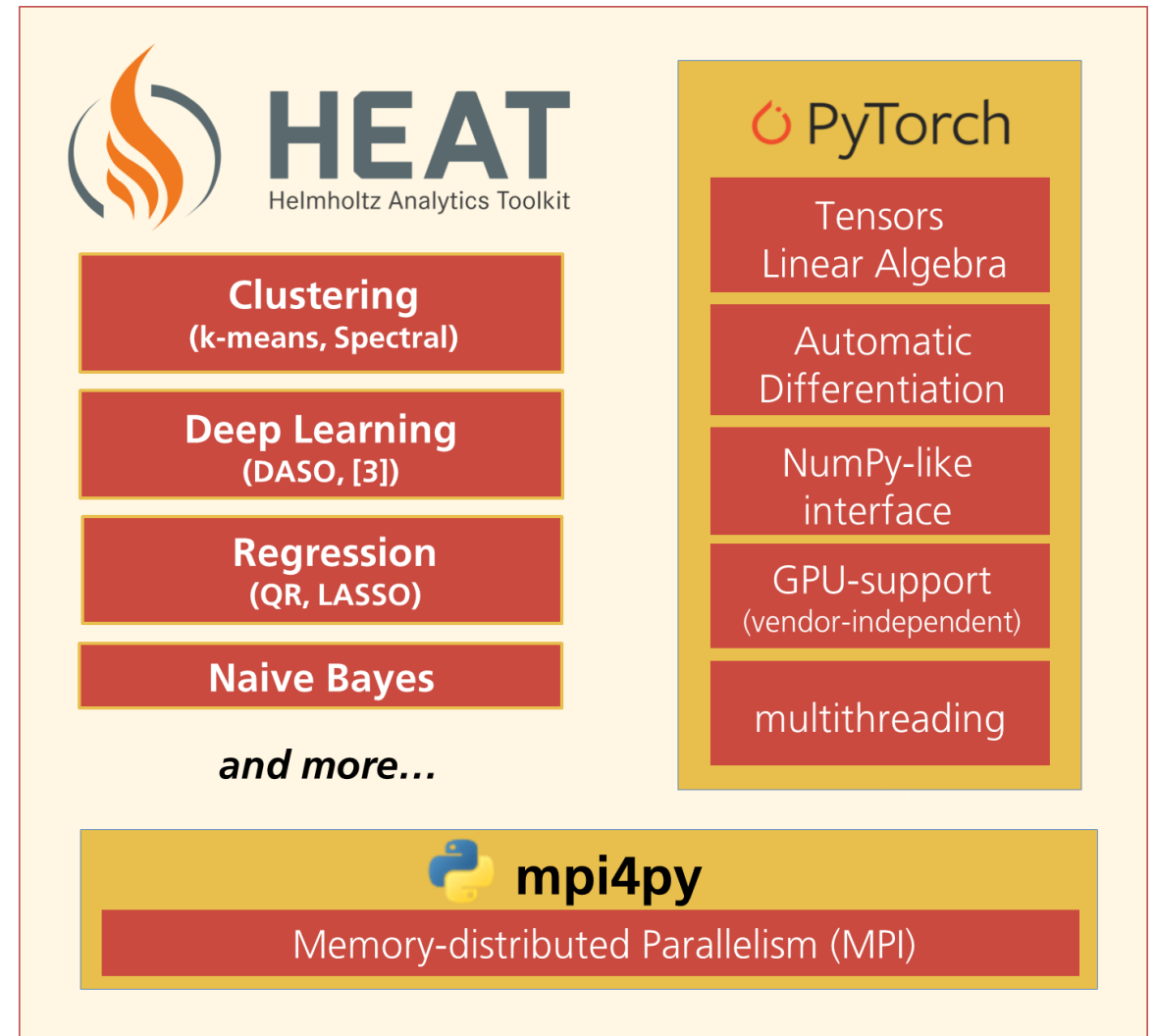


Memory-distributed Parallelism (MPI)

# ...and what you can do with it

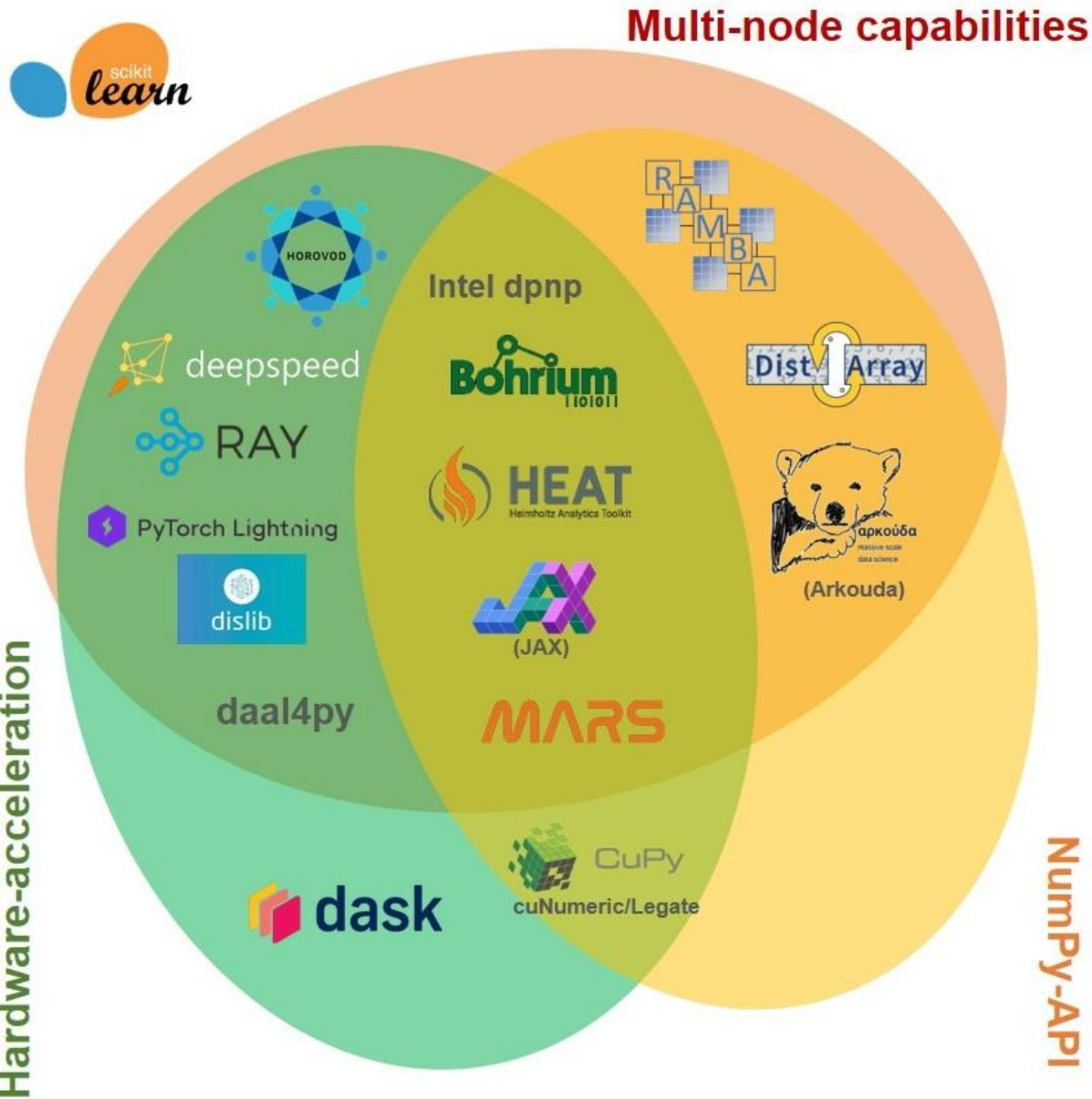
(main-branch as of August 2023)

- **Parallel I/O** via hdf5, csv, netcdf
- **Array creation and manipulation:** get/set item, concatenate, reshape, zeros, ones, arrays with random entries, math functions...
- basic sparse array class
- **Statistics:** min, max, mean, percentile, variance, ...
- **Linear algebra:** transpose, matmul, QR, Lanczos, (approximate) SVD, ...
- **Machine learning:** Clustering (k-means, spectral), Regression (LASSO), Naive Bayes, ...
- **Data parallel NNs** („DASO“)



```
3 # load data set from .h5-file (parallel I/O) on CPU
4 data = ht.load_hdf5('my_data_file.h5', 'my_data_set', split=0, dtype=ht.floa
5
6 # cluster the data set by K-Means
7 clustering = ht.cluster.KMeans(n_clusters=10, init="kmeans++")
8 clustering.fit(data)
9
10 # print the centers found
11 centers = clustering.cluster_centers_
```

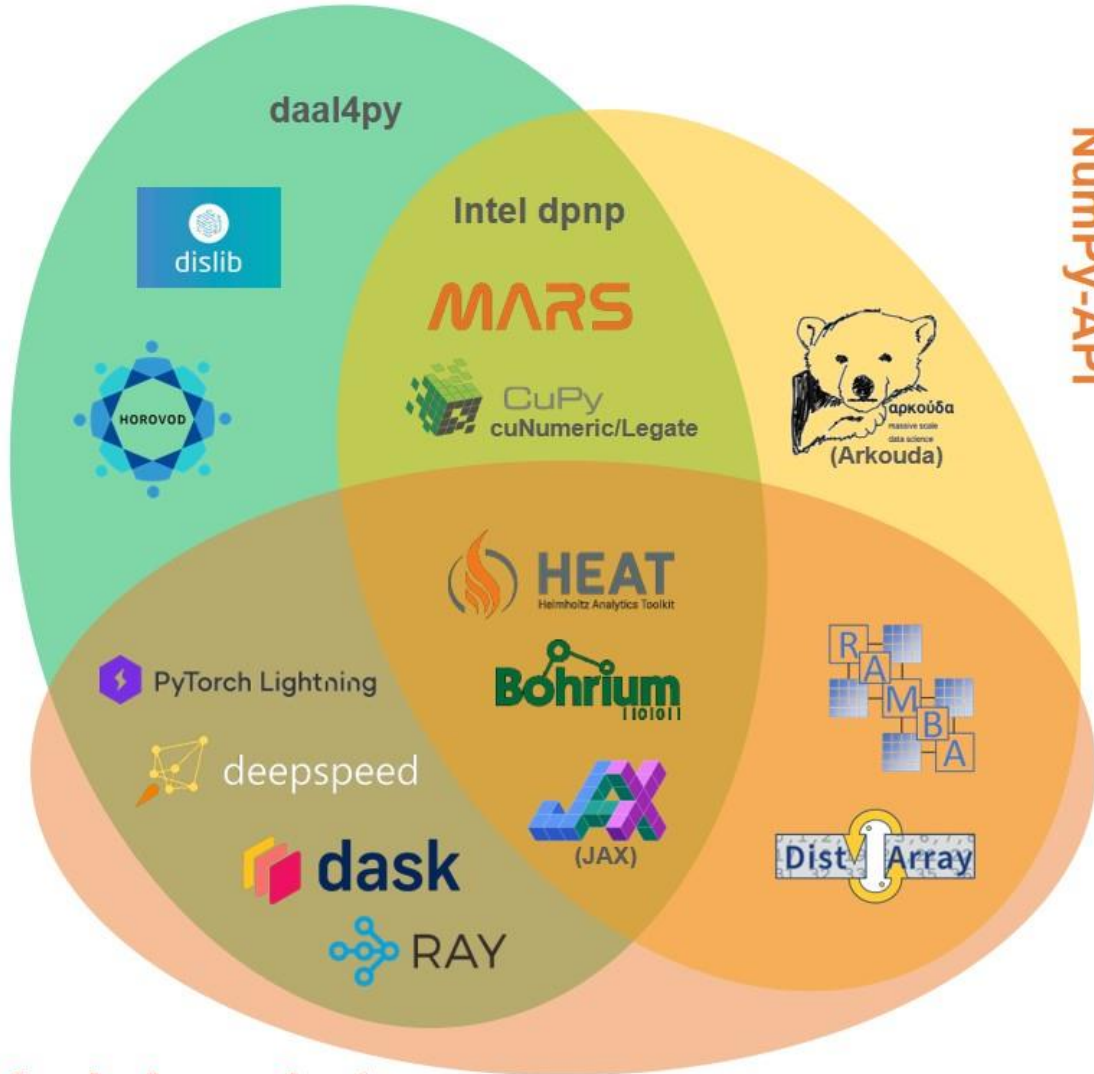
# Why another library for distributed array computing and ML in Python?



Heat...

- ...allows for **memory-distributed**, i.e. **multi-node**, parallelism and is not limited to embarrassingly parallel applications
- ...supports **hardware acceleration** (GPU)
- ...has a **simple NumPy-/scikit-learn-like API** that allows for easy adaption of legacy code

Hardware-acceleration



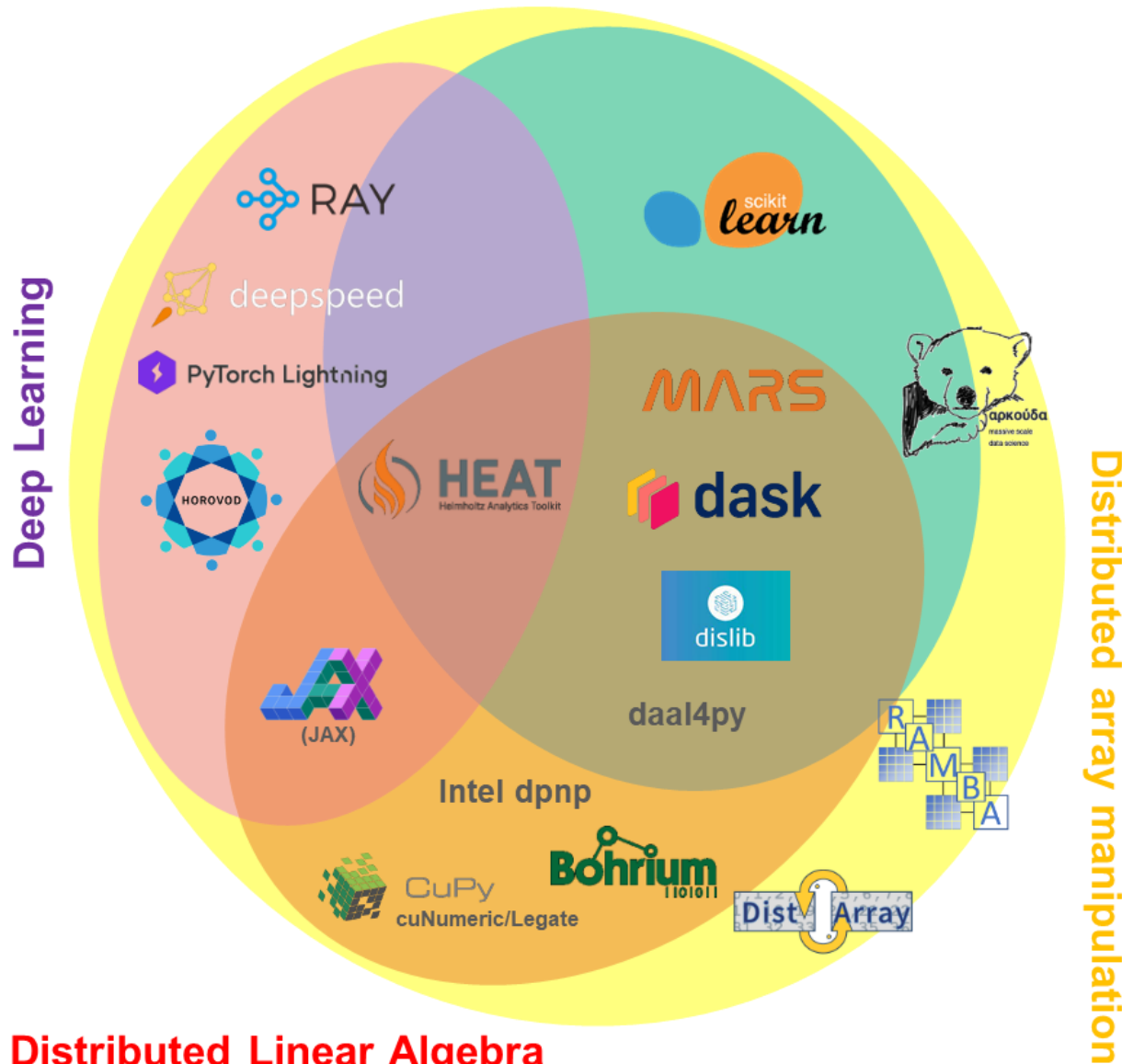
**Vendor independent**  
(hardware, software)

Heat...

- ...allows for **memory-distributed**, i.e. **multi-node**, parallelism and is not limited to embarrassingly parallel applications
- ...supports **hardware acceleration** (GPU)
- ...has a **simple NumPy-/scikit-learn-like API** that allows for easy adaption of legacy code
- ...is **vendor-/platform-independent**
- ...is **interoperable** (PyTorch)



## Machine Learning/Data Analytics



### Distributed Linear Algebra

Heat...

- ...allows for **memory-distributed**, i.e. **multi-node**, parallelism and is not limited to embarrassingly parallel applications
- ...supports **hardware acceleration** (GPU)
- ...has a **simple NumPy-/scikit-learn-like API** that allows for easy adaption of legacy code
- ...is **vendor-/platform-independent**
- ...is **interoperable** (PyTorch)
- ...is **general-purpose**
- ...is developed **by scientists for scientists**.

```
distances = ht.spatial.distance.cdist(x, centroids, quadrati
D2 = distances.min(axis=1)
D2.resplit_(axis=None)
prob = D2 / D2.sum()
random_position = ht.random.rand().item()
sample = 0
sum = 0
for j in range(len(prob)):
    if sum > random_position:
        break
    sum += prob[j].item()
sample = j
proc = 0
for p in range(x.comm.size):
    if displ[p] > sample:
        break
proc = p
```

In a bit more detail...

# Avoid memory becoming a bottleneck (a use-case from DLR)

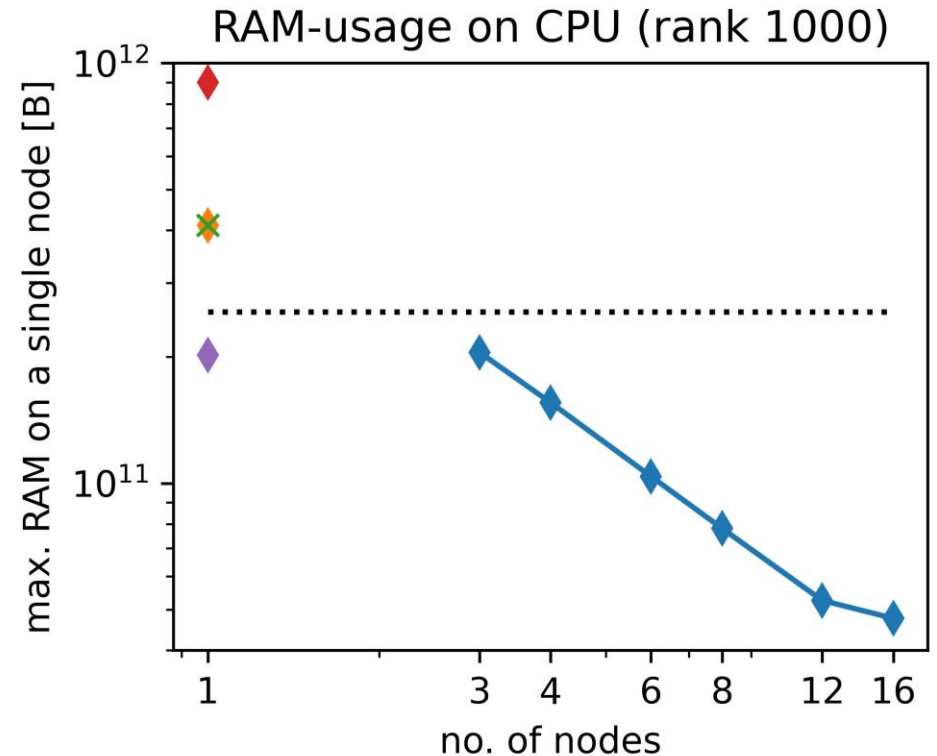


## Use-case:

analysis of a data set of 85487 RGB-images with 1024x192 pixels

- 589824 x 85487 matrix (~200GB)
- PCA (here: *truncated PCA with rank 1000*) and other typical tasks (e.g. clustering, matmul) are not data-parallel
- embarrassingly parallel batch-processing not possible!
- memory becomes a hard limit easily...

Experiments conducted on up to 16 ‚medium‘ compute nodes (Heat) and a single ‚bigmem‘ node (sklearn, numpy) of DLRs cluster CARO



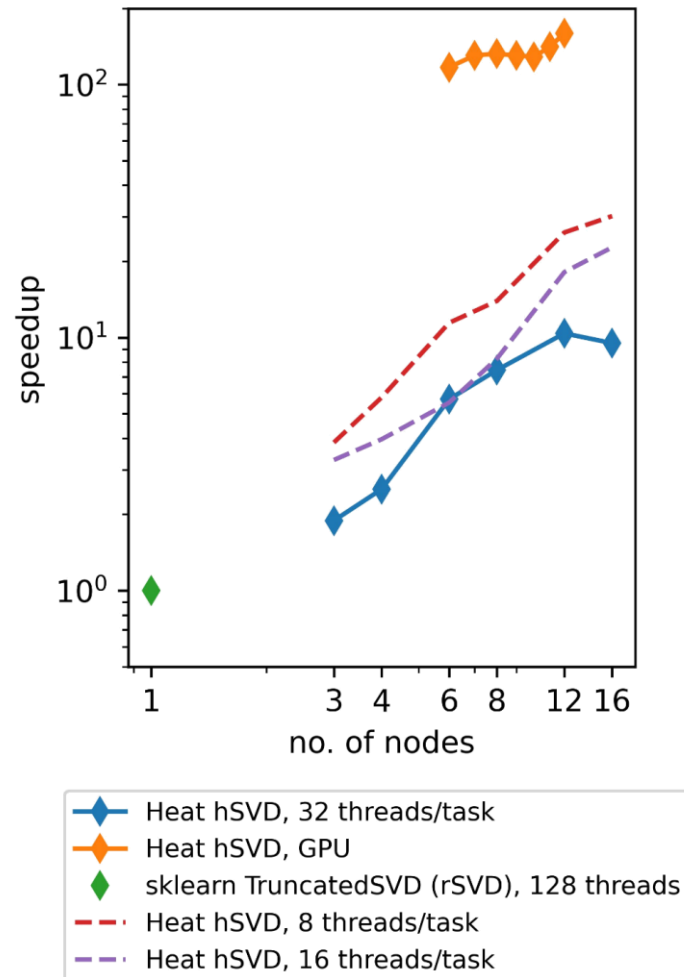
- ◆ Heat hSVD, 32 threads/task
- ◆ sklearn TruncatedSVD (rSVD), 128 threads
- × sklearn TruncatedSVD (ARPACK), 128 threads
- ◆ numpy full SVD (40% of columns only), 128 threads
- ◆ size of the array as hdf5-file

# Exploit your existing hardware: strong scaling

(use-case from DLR continued)



Computing times (Trunc. PCA, rank 50)



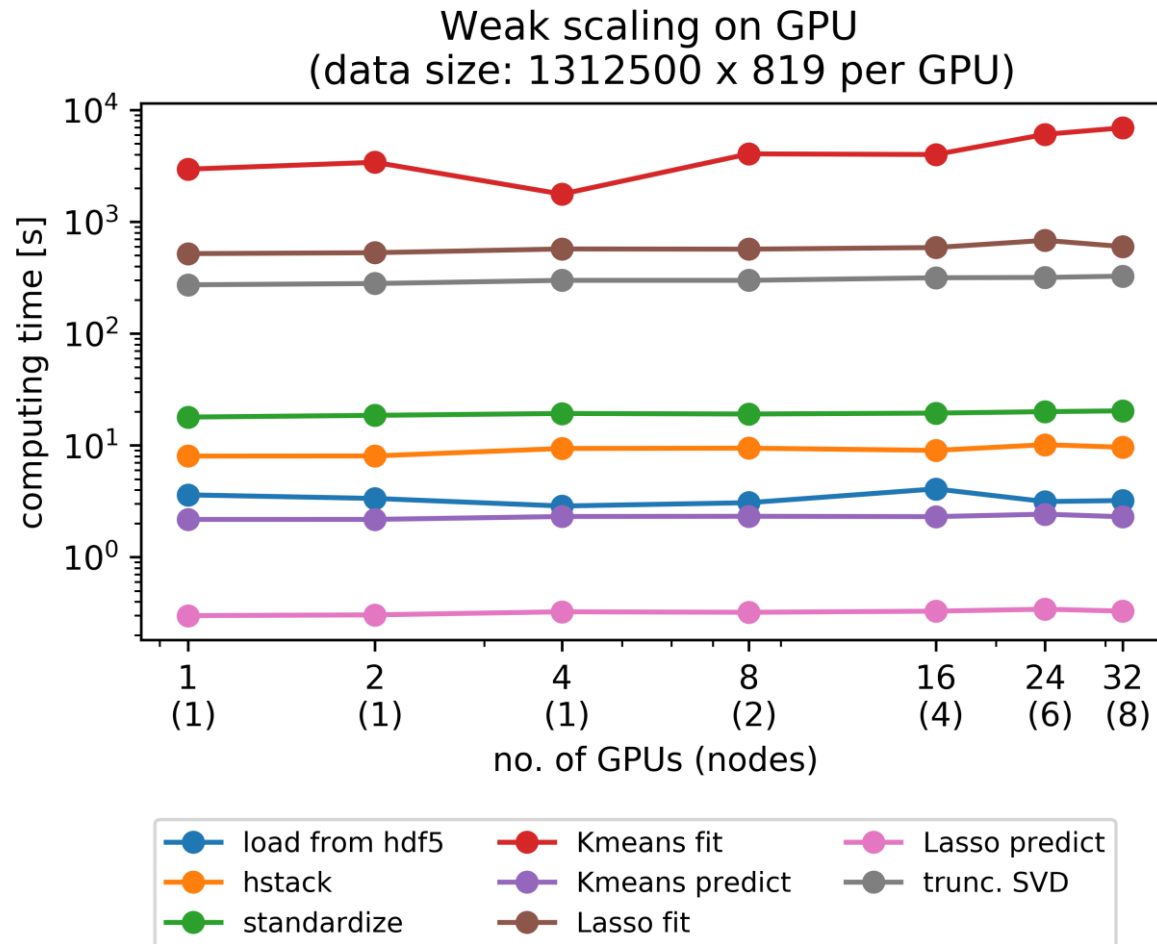
- **CPU:** hybrid parallelism (MPI+OpenMP) may achieve significant speedup compared to single node/single process
- **GPUs:** multi-GPU/multi-node setting is sometimes unavoidable due to memory limitations
- Heat so far runs on:
  - Intel and AMD CPUs
  - Nvidia GPUs
  - AMD GPUs (helmholtz.codebase runners)
  - Apple CPU+GPU (developers notebook, experimental)

## Experiments conducted on:

- **sklearn:** one ,bigmem' node of CARO
- **Heat CPU:** up to 16 ,medium' compute nodes of CARO
- **Heat GPU:** up to 12 nodes (with 4 Nvidia A100 80GB each) of DLRs cluster Terrabyte



# Scale your applications: weak scaling



## Test data set:

*ATLAS Top Tagging Open Data Set* (42 mio. entries with 819 features each ~ 129GB (“train” part) )

[ATLAS collaboration (2022). ATLAS Top Tagging Open Data Set. CERN Open Data Portal. DOI:10.7483/OPENDATA.ATLAS.FG5F.96GA]

Experiments conducted on up to 8 GPU-nodes of DLRs cluster Terrabyte (with 4 Nvidia A100 80GB each)



(Load with split=0, hstack along axis=1, Kmeans with 4 clusters, Lasso with regularization parameter 0.1 and tolerance 0.01, truncated SVD by hSVD with rel. tolerance 0.01)

# Benefits of the simple API



- Straight-forward adaptation of legacy NumPy-/scikit-learn-workflows:

```
import heat as ht

data = ht.load("your_hdf5_data_set.h5",
dataset="name_of_the_data", split=1, device="gpu")
clustering = ht.cluster.KMeans(n_clusters=10,
init="kmeans++")
clustering.fit(data)
print(clustering.cluster_centers_)
```

- Simple specification of devices (similar to PyTorch) allows to **use GPUs with almost no additional effort**

# Benefits of the simple API



- Except for prescribing a split-direction, **parallelization does not enter your python-script** explicitly
  - no explicit chunking of the data required
  - easy to run on HPC-systems using SLURM:

```
# on GPU
```

```
srun --nodes=8 --ntasks-per-node=4 --gres:gpu=4 python  
script.py
```

```
# on CPU (incl. multithreading)
```

```
srun --nodes=8 --ntasks-per-node=8 --cpus-per-task=20  
python script.py
```

# Comparison with Dask (CPU)

129GB data set\* on 8 CPU-nodes



	HEAT	DASK
<i>load from hdf5</i>	3 s	3 min
<i>stack data</i>	3 s	18 ms
<i>standardize data</i>	2 s	5 min
<i>K-means fit (4 clusters)...</i>	40 min	12 h
<i>...and predict</i>	1 s	7 min
<i>PCA (trunc. to rank 100)</i>	31 s	10 min
<i>randn...</i>	9 s	3 ms
<i>...and save to hdf5</i>	41 s	48 s
<b>MaxRSS for all this:</b>	<b>14 GB</b>	<b>460 GB</b>

Experiments conducted on up to 8 CPU-nodes of DLRs cluster Terrabyte (with 2 Intel Xeon Platinum 8380 40C 270W 2.3GHz each)

8 tasks per node, 20 threads per task



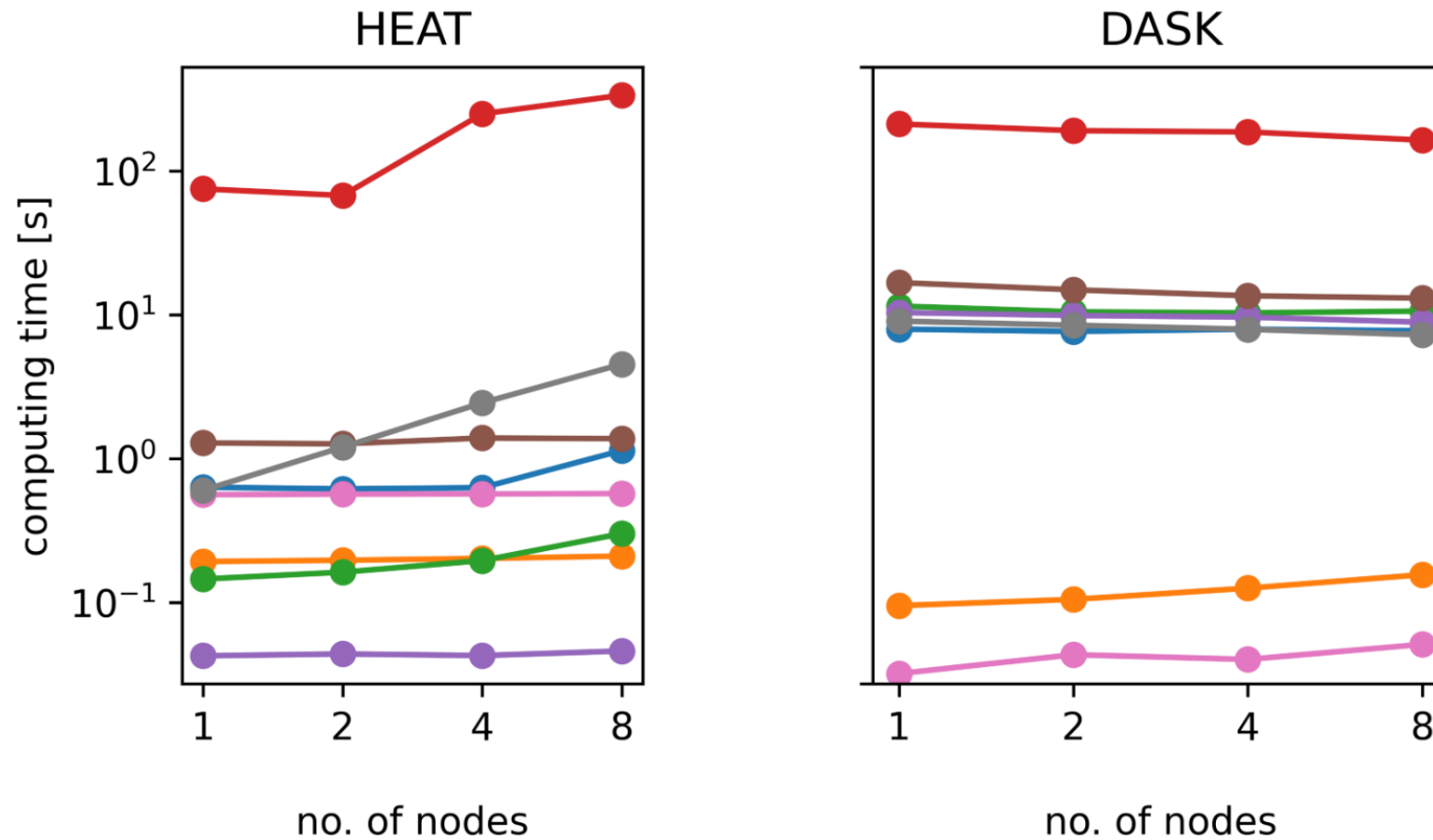
**Disclaimer:** These experiments have been run by a „Heat-Expert“, but a „Dask-Non-Expert“ – code available upon request.

\* „train“ part of *ATLAS Top Tagging Open Data Set*



# Comparison with Dask (CPU)

Weak scaling of runtime on a small data set\* (7.6GB)



Experiments conducted on up to 8 CPU-nodes of DLRs cluster Terrabyte (with 2 Intel Xeon Platinum 8380 40C 270W 2.3GHz each)

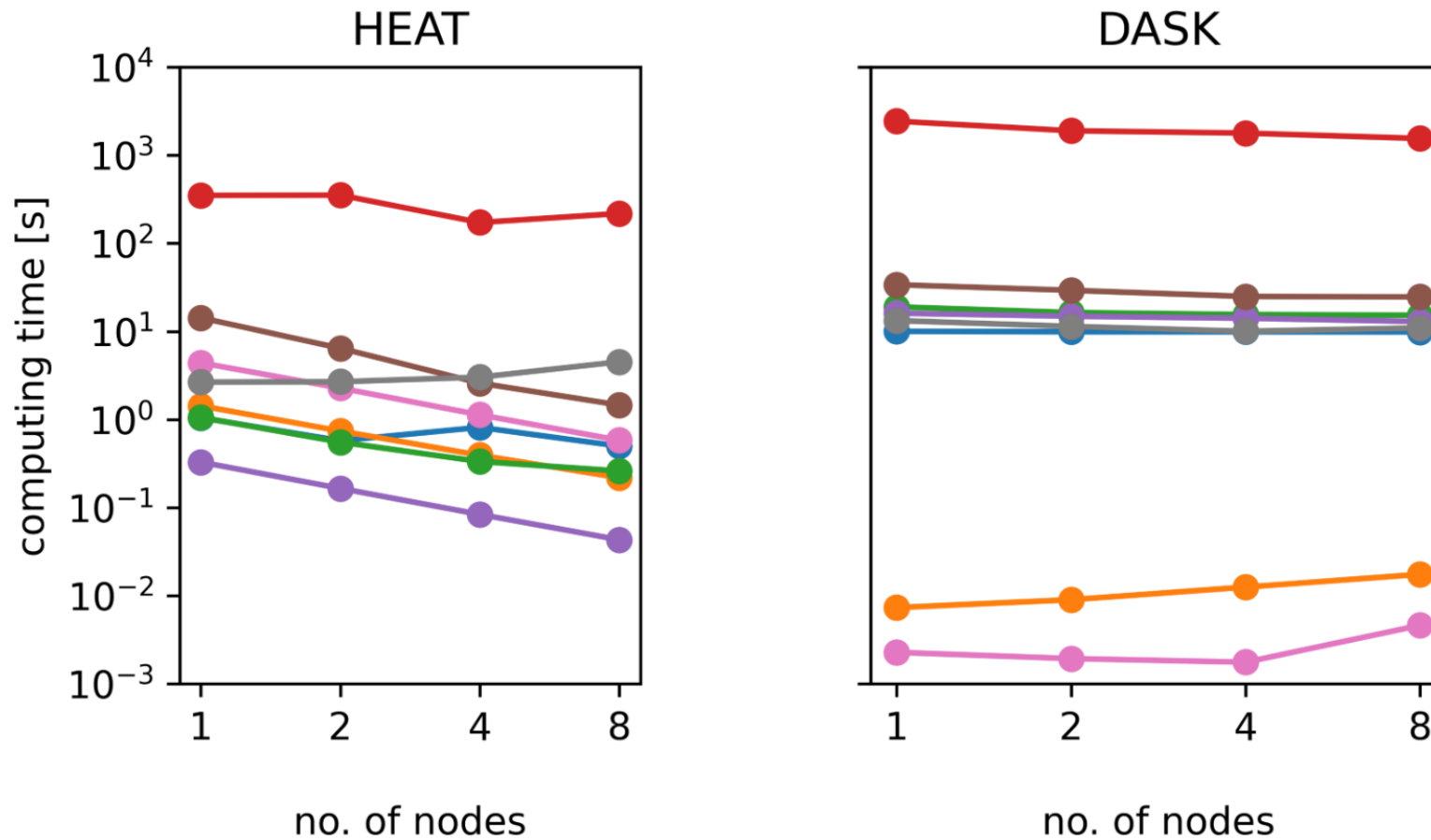
8 tasks per node, 20 threads per task

\* „test“ part of *ATLAS Top Tagging Open Data Set*



# Comparison with Dask (CPU)

Strong scaling of runtime on a small data set\* (7.6GB)



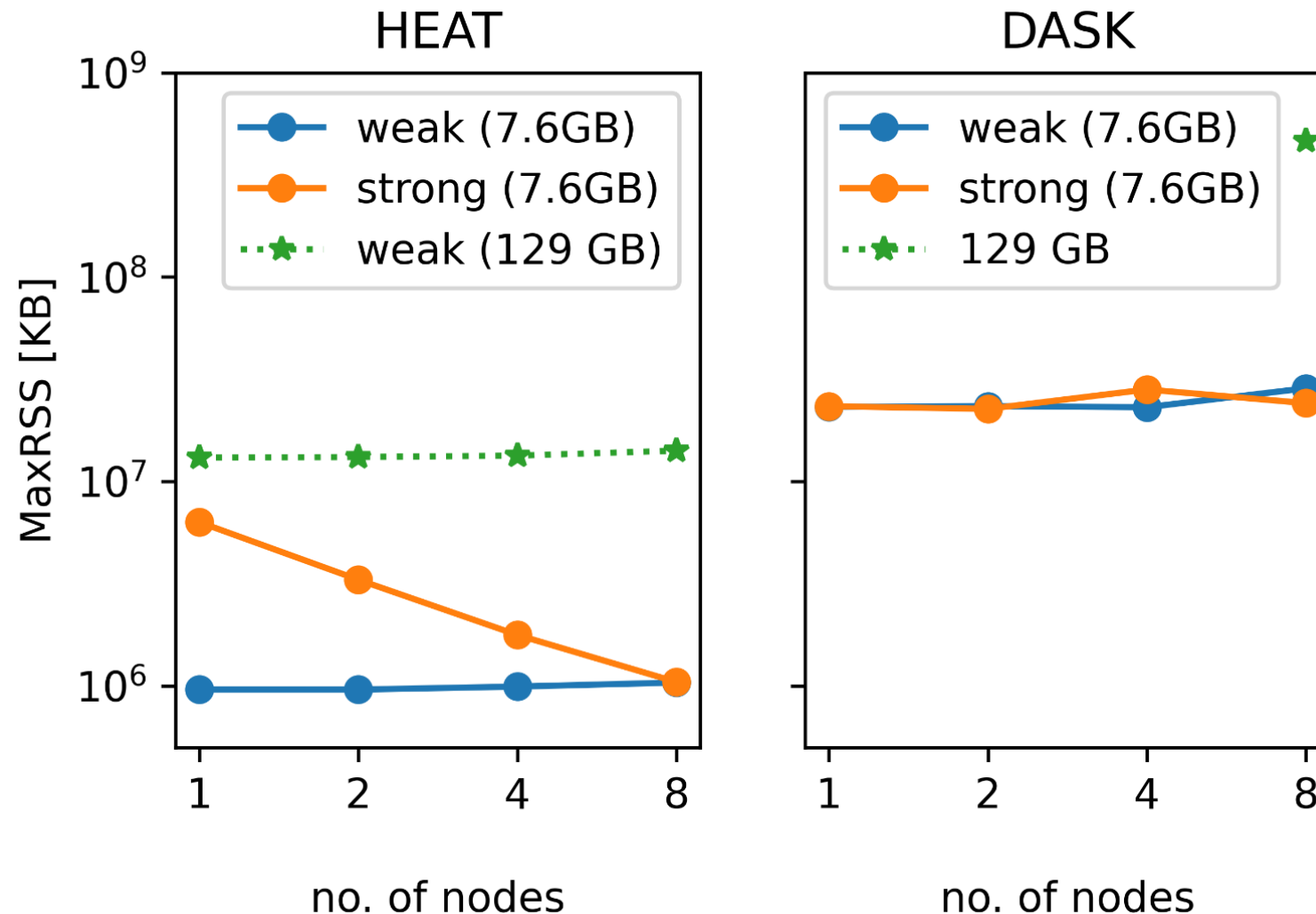
Experiments conducted on up to 8 CPU-nodes of DLRs cluster Terrabyte (with 2 Intel Xeon Platinum 8380 40C 270W 2.3GHz each)

8 tasks per node, 20 threads per task

\* „test“ part of *ATLAS Top Tagging Open Data Set*

# Comparison with Dask (CPU)

Weak and strong scaling of used RAM on a small data set\* (7.6GB)



Experiments conducted on up to 8 CPU-nodes of DLRs cluster Terrabyte (with 2 Intel Xeon Platinum 8380 40C 270W 2.3GHz each)

8 tasks per node, 20 threads per task

\* „test“ part of ATLAS Top Tagging Open Data Set

```
sz = chunk_map[en, axis]
if arr0.comm.rank == en:
    lcl_slice = [slice(None)] * arr0.ndim
    lcl_slice[axis] = slice(ttl, ttl + sz)
    t_arr0 = t_arr0[lcl_slice].clone()
    ttl += sz.item()
```

```
if len(t_arr0.shape) < len(t_arr1.shape):
    t_arr0.unsqueeze_(axis)
```

```
if s1 is None:
```

```
    arb_slice = [None] * len(arr0.shape)
```

```
    arb_slice[axis] = slice(0, arr0.shape[axis])
    chunk_map[arb_slice] -= lshape_map[tuple([0]
```

# Future development and challenges

# Future development and challenges



- **Software engineering:** keep/increase *software quality* by, e.g., ...
  - continuous benchmarking of runtime, memory consumption, and energy consumption („green HPC“)
  - keep/extend interoperability w.r.t. the evolving ecosystem
- **Extend functionality,** e.g.,
  - implementation of support for distributed `xarray`'s has just started (→ EOC-community)
  - implementation of an SVD/PCA/DMD-module (→ ESA)
  - preprocessing module, FFT, ... → ***or maybe your suggestion?***
- **Community interaction,** e.g., ...
  - extend documentation by tutorials/use cases etc.
  - user-support by workshops and/or individual collaboration/help:  
***If you want to try out/use Heat for your research, we are glad to help you; we do not expect you to contribute to Heat.***

*Feel free to try out,  
to suggest features,  
or to contribute!*



visit us on **github**:



# Impressum



- Thema: **The Helmholtz Analytics Toolkit (Heat)** and its role in the landscape of massively-parallel scientific Python
- Datum: 2023-08
- Autor: Fabian Hoppe (+ the Heat development team)
- Institut: Institut für Softwaretechnologie, High-Performance Computing
- Bildcredits: Alle Bilder „DLR (CC BY-NC-ND 3.0)“, sofern nicht anders angegeben