

# HiPSCat: Enabling Storage and Analytics of Large-Scale Catalogs

**Mario Juric**

DiRAC Institute Director | LINCC FW Co-I  
Professor of Astronomy, University of Washington

with Sam Wyatt, Sean McGuire, Melissa DeLucchi, Max West, Doug Branton,  
Neven Caplar, Steven Stetzler, Colin Slater, Jeremy Kubica  
and the LINCC Frameworks Analytics Group



DATA INTENSIVE RESEARCH IN  
ASTROPHYSICS AND COSMOLOGY



LINCC



AST-2003196



# The LINCC Frameworks Project

LSST Interdisciplinary Network For Collaboration And Computing

*To collaboratively develop open computing systems and algorithms needed for large-survey analyses.*

- Data analysis infrastructure (this talk)
- Solar System exploration
- Time domain science
- Extragalactic astronomy

Two LINCC-FW hubs:

- Carnegie Mellon University
- University of Washington



## **The Legacy Survey of Space and Time**

*Deep synoptic optical survey, coming in 2025.*

Repeated imaging of the visible sky to ~24th mag.

10 years of operation.

60 PB of raw data.

40 billion stars, galaxies, asteroids.

30 trillion observations.

# Rubin Observatory, March 15, 2023.

Cerro Pachon, Chile



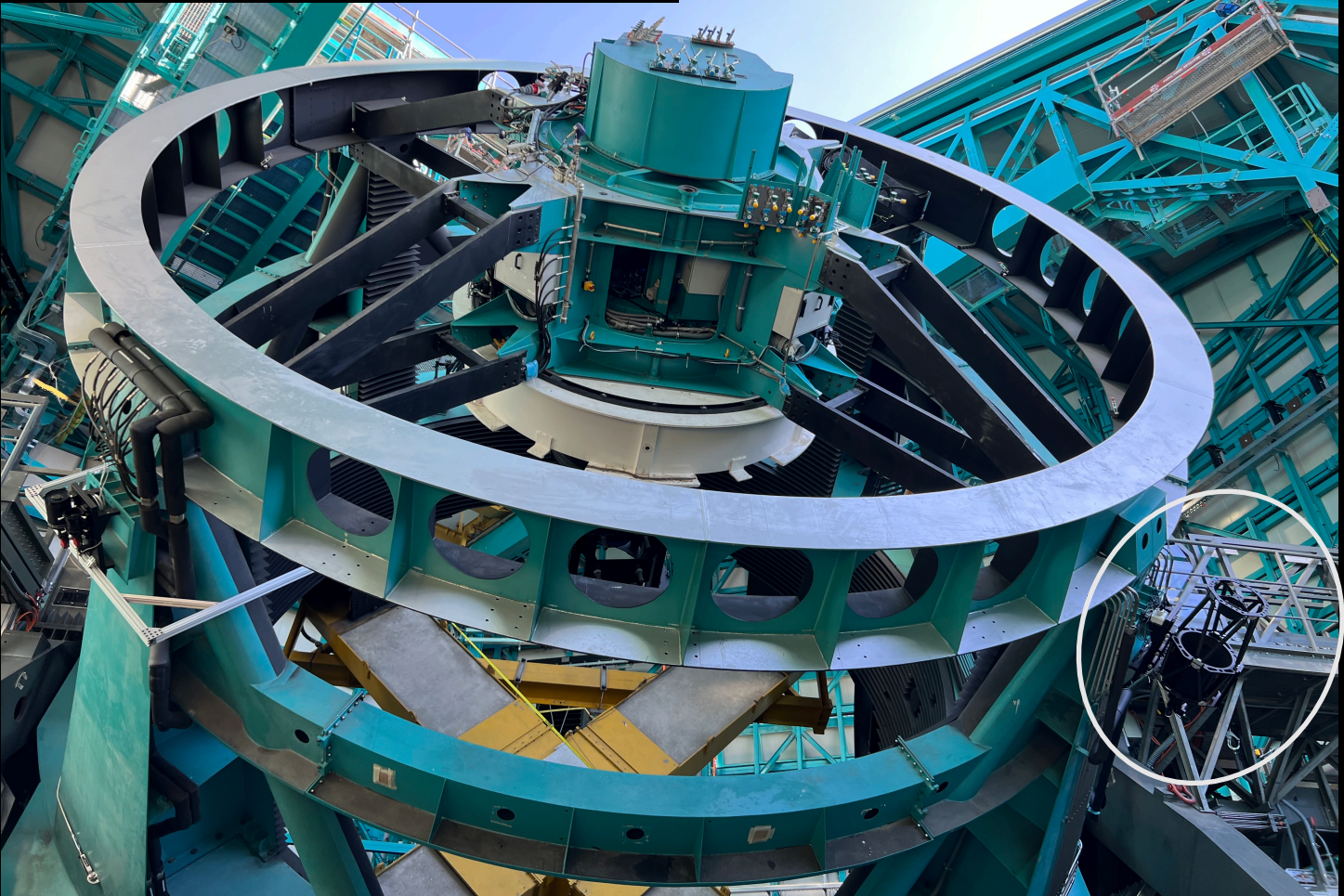






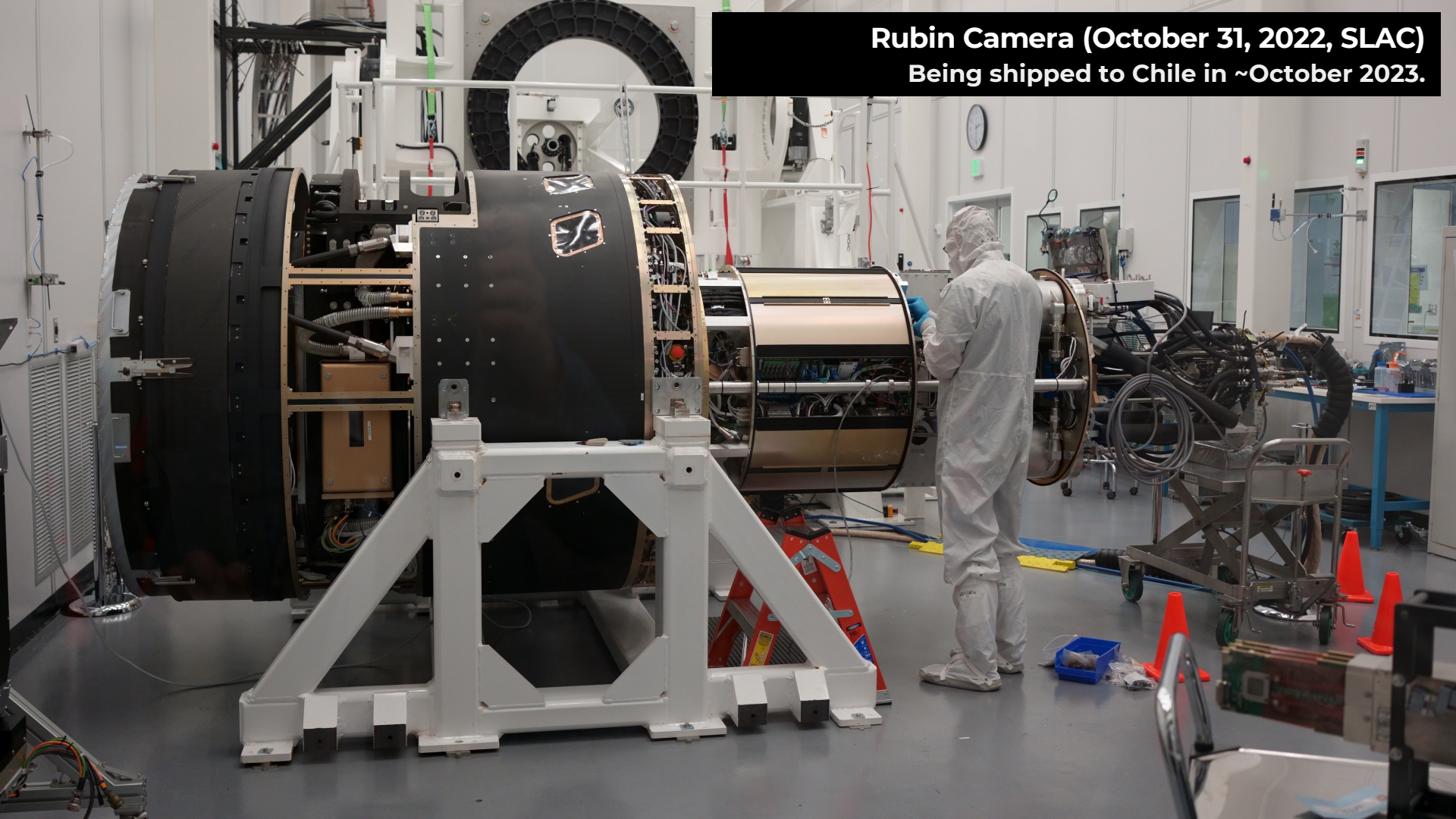
# Rubin Observatory, Telescope Mount Assembly

Ran full night of "observing" (a week ago)



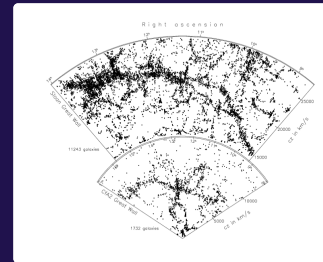
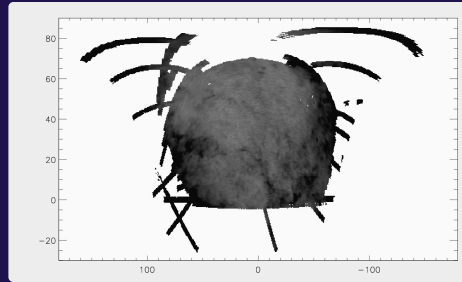
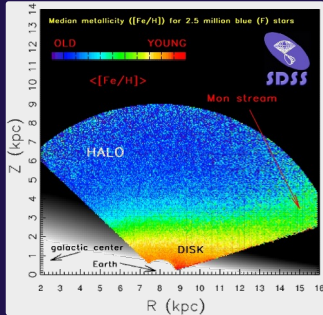
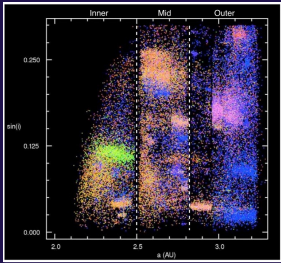


**Rubin Camera (October 31, 2022, SLAC)**  
Being shipped to Chile in ~October 2023.





*One of the major things the community will want to do with Rubin is whole dataset science.*



# Scale of the problem



Rubin Year #1 dataset:

- 10Bn objects
- 100 obsv/object == 1T observations
- 100 bytes/obsv == **100 TB**

Not just a Rubin Problem:

- Gaia, DES, ZTF, WISE, PS, Euclid, Rubin, Roman, SphereX, ...
- Each one of these is Bn+ objects (w. many more measurements)



# Fellow Travelers on the Quest



Sam Wyatt (Product Owner), Sean McGuire, Melissa DeLucchi (Project Manager)  
Max West, Doug Branton, Neven Caplar, Steven Stetzler, Jeremy Kubica

Vandana Desai (IPAC), Troy Raen (IPAC), Dave Shupe (IPAC), Brigitta Sipőcz (IPAC)  
Gregory Dubois-Felsmann (IPAC & Rubin)

Colin Slater (DiRAC & Rubin)

Sharon Shen (STScI)

Susan Mullally (STScI)

Rick White (STScI)

Bernie Shiao (STScI)

Travis Berger (STScI)

Erik Tollerud (STScI)

Josh Peek (STScI)

Tess Jaffe (HESARC)

+ YOU (join the party!)



# Large-dataset Analytics: Partitioned Files

- Relational databases are not ideal for this type of work. Poor UX, too many bottlenecks.
- Industry state-of-the-art is to use distributed analytics tools operating on files.
- Distributed computation achieved through partitioning.

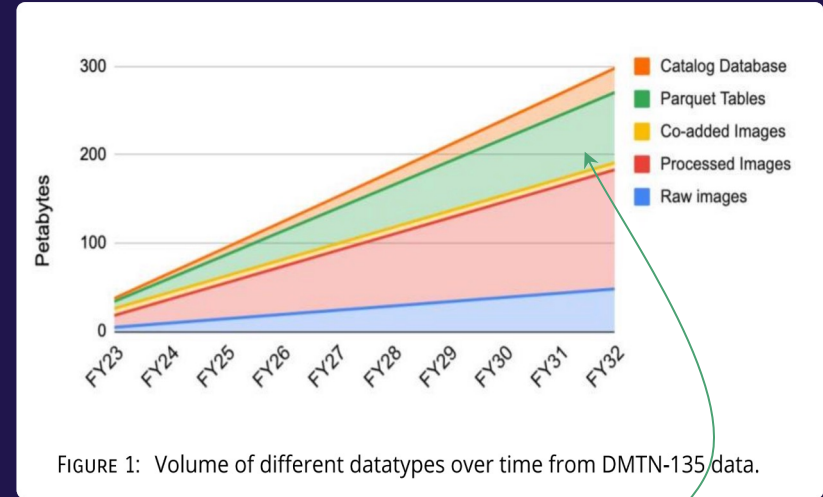
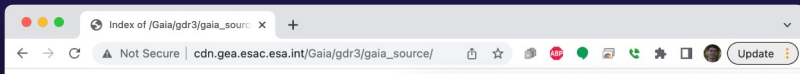


FIGURE 1: Volume of different datatypes over time from DMTN-135 data.

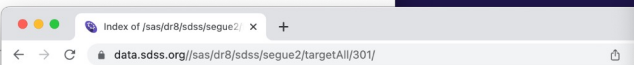
```
import pandas as pd
pd.read_parquet('example_pa.parquet', engine='pyarrow')
```



# How to partition? Historically, we haven't generally given this much thought...



File Name	File Size
GaiaSource_000000-003111.csv.gz	05-May-2022
GaiaSource_003112-005263.csv.gz	05-May-2022
GaiaSource_005264-006601.csv.gz	05-May-2022
GaiaSource_006602-007952.csv.gz	05-May-2022
GaiaSource_007953-010234.csv.gz	05-May-2022
GaiaSource_010235-012597.csv.gz	05-May-2022
GaiaSource_012598-014045.csv.gz	05-May-2022
GaiaSource_014046-015169.csv.gz	05-May-2022
GaiaSource_015170-016240.csv.gz	05-May-2022
GaiaSource_016241-017018.csv.gz	05-May-2022
GaiaSource_017019-017658.csv.gz	05-May-2022
GaiaSource_017659-018028.csv.gz	05-May-2022
GaiaSource_018029-018477.csv.gz	05-May-2022
GaiaSource_018478-019161.csv.gz	05-May-2022
GaiaSource_019162-019657.csv.gz	05-May-2022
GaiaSource_019658-020091.csv.gz	05-May-2022
GaiaSource_020092-020493.csv.gz	05-May-2022
GaiaSource_020494-020767.csv.gz	05-May-2022
GaiaSource_020768-020984.csv.gz	05-May-2022
GaiaSource_020985-021233.csv.gz	05-May-2022
GaiaSource_021234-021411.csv.gz	05-May-2022
GaiaSource_021412-021665.csv.gz	05-May-2022
GaiaSource_021666-021919.csv.gz	05-May-2022
GaiaSource_021920-022158.csv.gz	05-May-2022
GaiaSource_022159-022410.csv.gz	05-May-2022
GaiaSource_022411-022698.csv.gz	05-May-2022
GaiaSource_022699-022881.csv.gz	05-May-2022
GaiaSource_022882-023058.csv.gz	05-May-2022
GaiaSource_023059-023264.csv.gz	05-May-2022
GaiaSource_023265-023450.csv.gz	05-May-2022
GaiaSource_023451-023649.csv.gz	05-May-2022
GaiaSource_023650-023910.csv.gz	05-May-2022
GaiaSource_023911-024205.csv.gz	05-May-2022
GaiaSource_024206-024526.csv.gz	05-May-2022
GaiaSource_024527-025166.csv.gz	05-May-2022
GaiaSource_025167-025691.csv.gz	05-May-2022
GaiaSource_025692-026057.csv.gz	05-May-2022
GaiaSource_026058-026390.csv.gz	05-May-2022
GaiaSource_026391-026648.csv.gz	05-May-2022
GaiaSource_026649-027106.csv.gz	05-May-2022
GaiaSource_027107-027517.csv.gz	05-May-2022
GaiaSource_027518-027832.csv.gz	05-May-2022
GaiaSource_027833-028076.csv.gz	05-May-2022
GaiaSource_028077-028318.csv.gz	05-May-2022



File Name	File Size
Parent directory/	-
1000/	-
1006/	-
1009/	-
1010/	-
1011/	-
1013/	-
1022/	-
1024/	-
1033/	-
1035/	-
1037/	-
1040/	-
1043/	-
1045/	-
1055/	-
1056/	-
1057/	-
109/	-
1119/	-
1120/	-
1122/	-
1133/	-
1140/	-
1142/	-
1231/	-
1233/	-
1239/	-
1241/	-
125/	-
1302/	-
1329/	-
1331/	-

## ZTF ALERT ARCHIVE

What is included?

Below you will find compressed tar archives of ZTF event alerts (observations detected in image differences). Each tar file contains alerts collected in the given night (UTC-based), with each alert stored in a separate file in the AVRO format. To get you started, we offer a repository with low basic utilities for reading AVRO-serialized data, as well as an example Jupyter notebook. The schema fields are described here.

Why this service?

We are providing this archive as simple alternative to public event brokers. Full-featured event brokers that provide real-time access to these alerts include MARS, Lasair, ANTARES, and ALERCE.

Known caveats

- The data provided on this site is generated automatically. The files provided contain a full, unfiltered, 5-sigma alert stream. Depending on your science case, you may wish to improve the purity of your sample by filtering the data on the included attributes such as the signal-to-noise ratio or the real-bogus score.
- Users interested in un-subtracted archival photometry should consider the ZTF Data Releases, accessible at IRSA.
- A subset of events obtained through Caltech time are made public here in "program33" tarballs; as of this writing these are additional observations of the current TESS sector.

Name	Last modified	Size
ztf_public_20220810.tar.gz	10 hours ago	8.3G
ztf_public_20220809.tar.gz	1 day ago	8.3G
ztf_public_20220808.tar.gz	2 days ago	1.7G
ztf_public_20220807.tar.gz	3 days ago	10G
ztf_public_20220806.tar.gz	4 days ago	11G
ztf_public_20220805.tar.gz	5 days ago	3.6G
ztf_public_20220804.tar.gz	5 days ago	6.3G
ztf_public_20220803.tar.gz	7 days ago	11G
ztf_public_20220802.tar.gz	8 days ago	17G
ztf_public_20220801.tar.gz	9 days ago	3.4G
ztf_public_20220731.tar.gz	10 days ago	2.8G
ztf_public_20220730.tar.gz	11 days ago	14G

# There's value in thinking this through, and standardizing

- Users know what to expect and how to handle the dataset
- High quality, shared, analytics tools can be written
- Multi-dataset analytics can supported
- Pre-staging/ETL may be avoided
  
- Bulk export files == bulk analytics files
- Easier to generate and support for providers
- Can share code and infrastructure (e.g. mirroring, caching)

*Think of all the wonderful tools and ecosystems that sprung up around HiPS!*





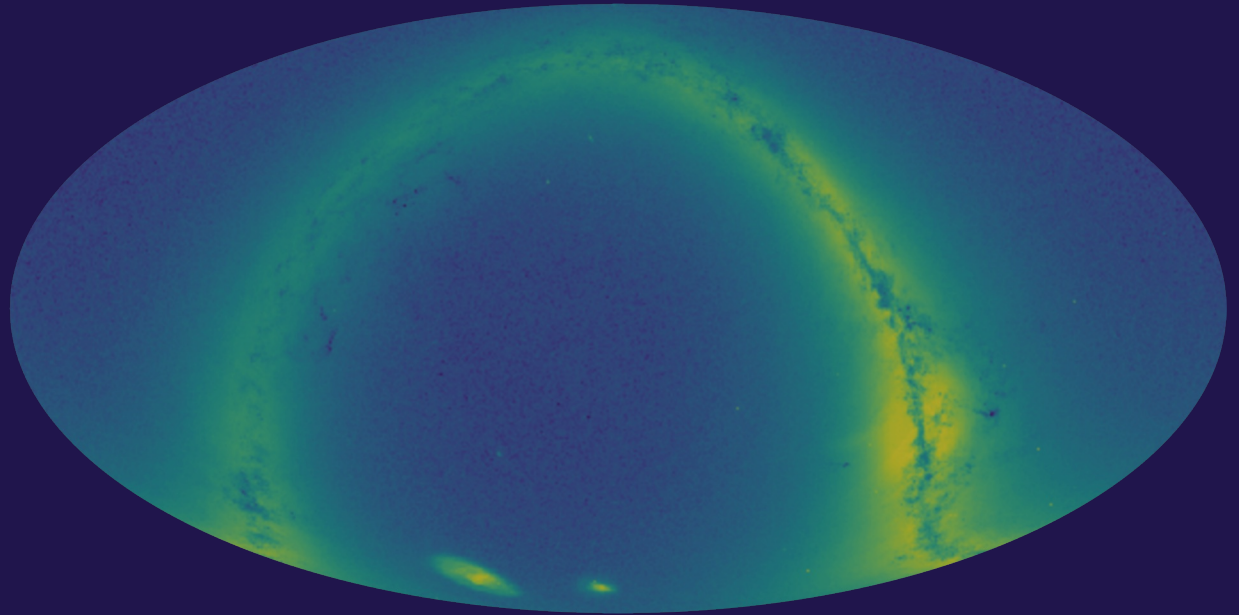
***Imagine a dataset with approximately  
1,811,709,771 sources...***

***... how could we store it?***

# Gaia DR3

The 1.8Bn sources  
released in Gaia DR3

A single ASCII file would  
be about ~680GB in size,  
(gzip compressed).





# 1. Partitioning: HEALPix

Partition the sky into NSIDE=1 (order=0)  
HEALPix tiles, map tiles to files.

Example:

```
Norder0-Npix0.tsv.gz  
Norder0-Npix1.tsv.gz  
Norder0-Npix2.tsv.gz  
Norder0-Npix3.tsv.gz  
Norder0-Npix4.tsv.gz  
Norder0-Npix5.tsv.gz  
Norder0-Npix6.tsv.gz  
Norder0-Npix7.tsv.gz  
Norder0-Npix8.tsv.gz  
Norder0-Npix9.tsv.gz  
Norder0-Npix10.tsv.gz  
Norder0-Npix11.tsv.gz
```



# Problem: Severely unbalanced file sizes

Pixel 4 (Galactic pole) ~ 20GB

Pixel 10 (Galactic center) ~ 400GB.

Simple file-based parallelization fails.

## Example

Norder0-Npix0.tsv.gz

Norder0-Npix1.tsv.gz

Norder0-Npix2.tsv.gz

Norder0-Npix3.tsv.gz

Norder0-Npix4.tsv.gz

Norder0-Npix5.tsv.gz

Norder0-Npix6.tsv.gz

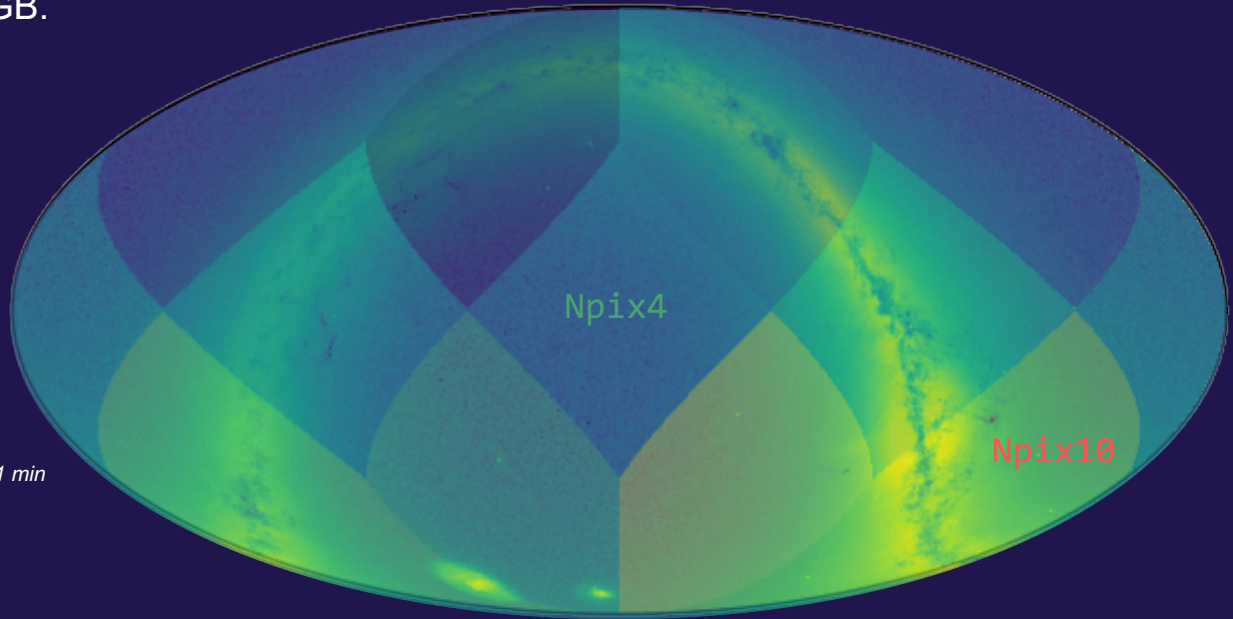
Norder0-Npix7.tsv.gz

Norder0-Npix8.tsv.gz

Norder0-Npix9.tsv.gz

Norder0-Npix10.tsv.gz

Norder0-Npix11.tsv.gz

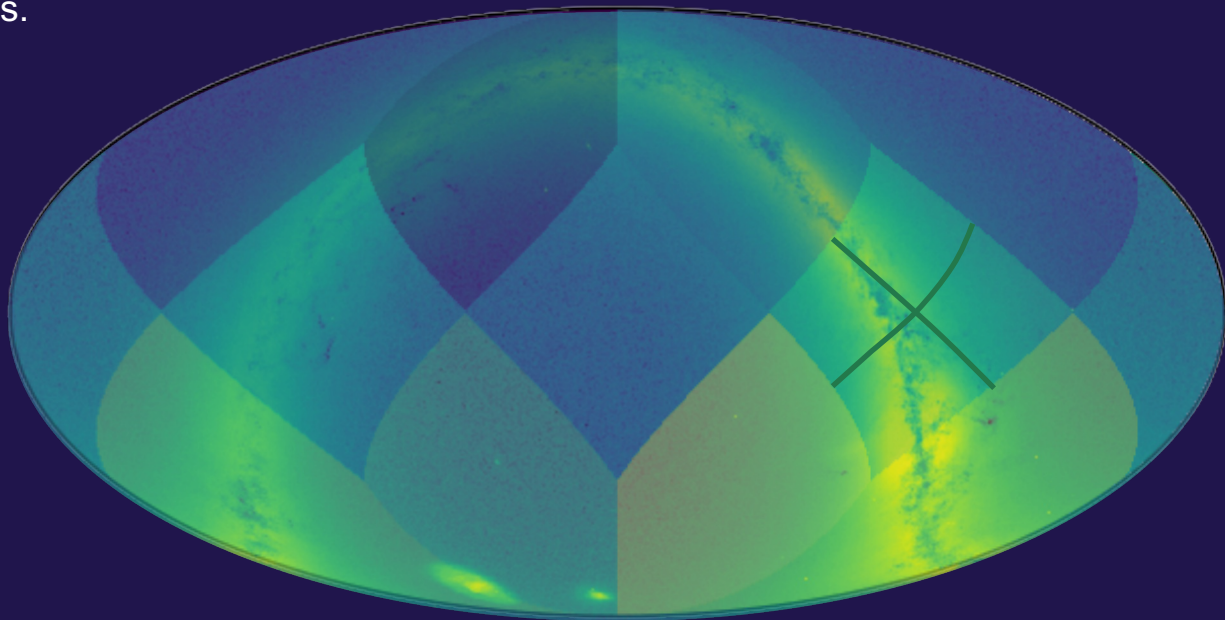


# Solution: Partition Hierarchically

If too many sources fall into a pixel,  
split it into four higher order pixels.

## Example

Norder0-Npix0.tsv.gz  
...  
Norder0-Npix7.tsv.gz  
...  
Norder0-Npix11.tsv.gz



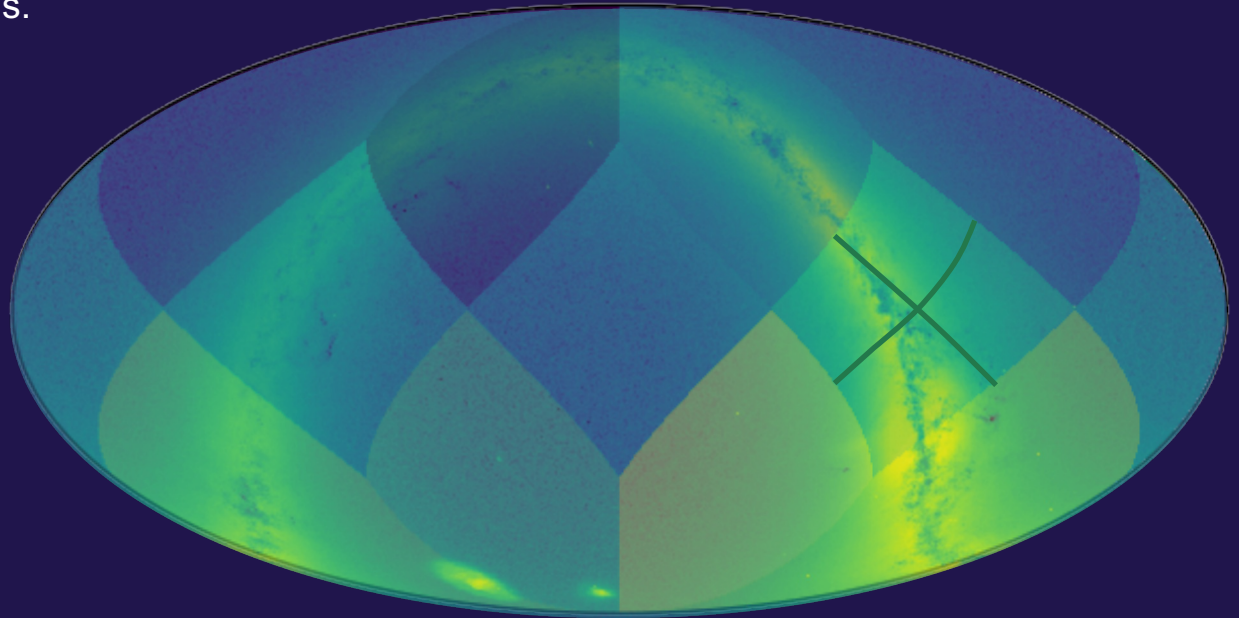


# 1. Partition Hierarchically

If too many sources fall into a pixel,  
split it into four higher order pixels.

## Example

```
Norder0-Npix0.tsv.gz  
...  
Norder1-Npix28.tsv.gz  
Norder1-Npix29.tsv.gz  
Norder1-Npix30.tsv.gz  
Norder1-Npix31.tsv.gz  
...  
Norder0-Npix11.tsv.gz
```



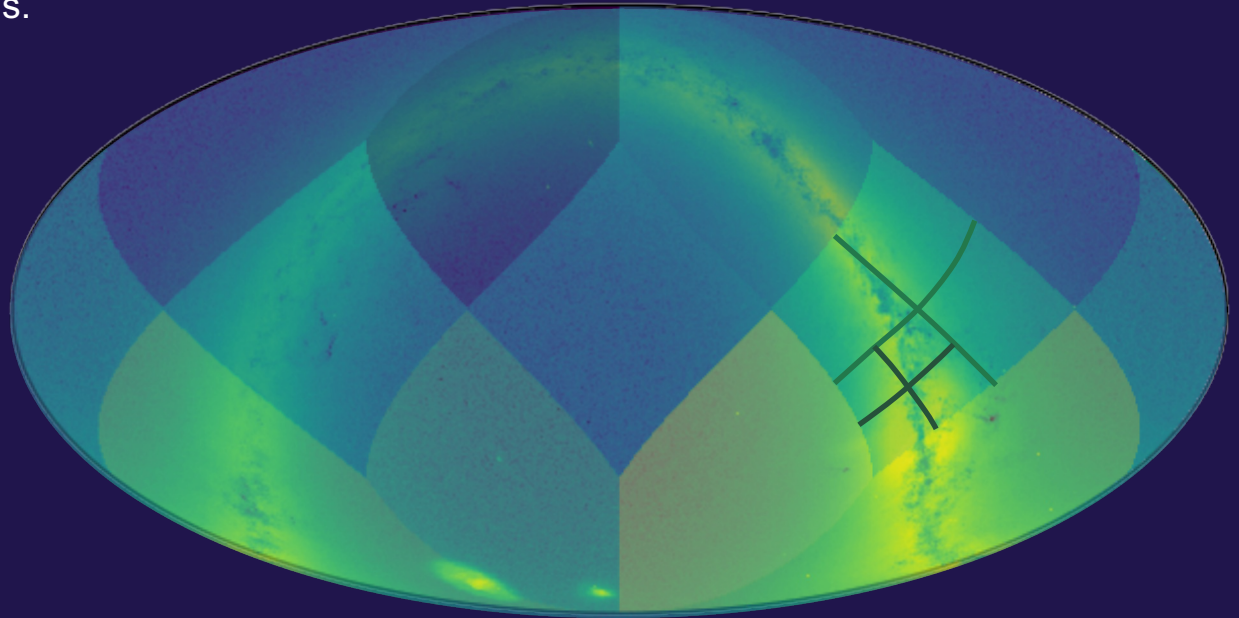
# 1. Partition Hierarchically

If too many sources fall into a pixel,  
split it into four higher order pixels.

Repeat.

## Example

```
Norder0-Npix0.tsv.gz  
...  
Norder1-Npix28.tsv.gz  
Norder1-Npix29.tsv.gz  
Norder1-Npix30.tsv.gz  
Norder1-Npix31.tsv.gz  
...  
Norder0-Npix11.tsv.gz
```





# 1. Partition Hierarchically

If too many sources fall into a pixel,  
split it into four higher order pixels.

Repeat.

## Example

Norder0-Npix0.tsv.gz

...

Norder1-Npix28.tsv.gz

Norder1-Npix29.tsv.gz

Norder1-Npix30.tsv.gz

Norder2-Npix112.tsv.gz

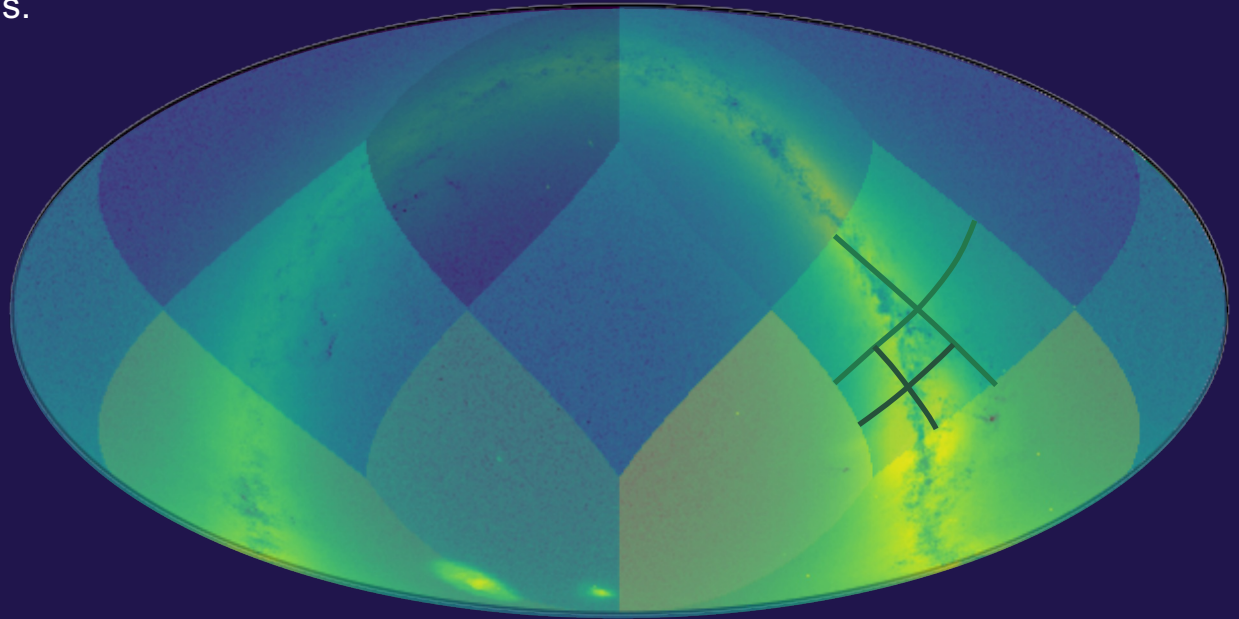
Norder2-Npix113.tsv.gz

Norder2-Npix114.tsv.gz

Norder2-Npix115.tsv.gz

...

Norder0-Npix11.tsv.gz



# 1. Partition Hierarchically

If too many sources fall into a pixel,  
split it into four higher order pixels.

Repeat until each file size  
is beneath some  
pre-defined threshold.

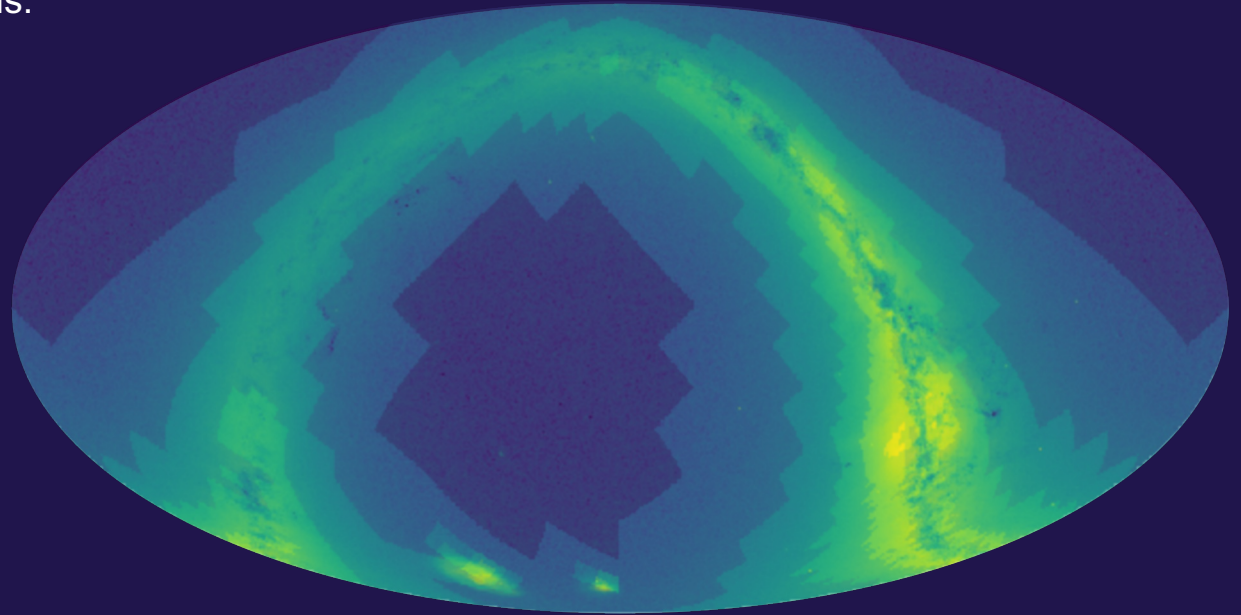


Figure: an overlay of  
Gaia counts and the  
partitioning map, taking  
MAXOBJECTS=1e6

order:



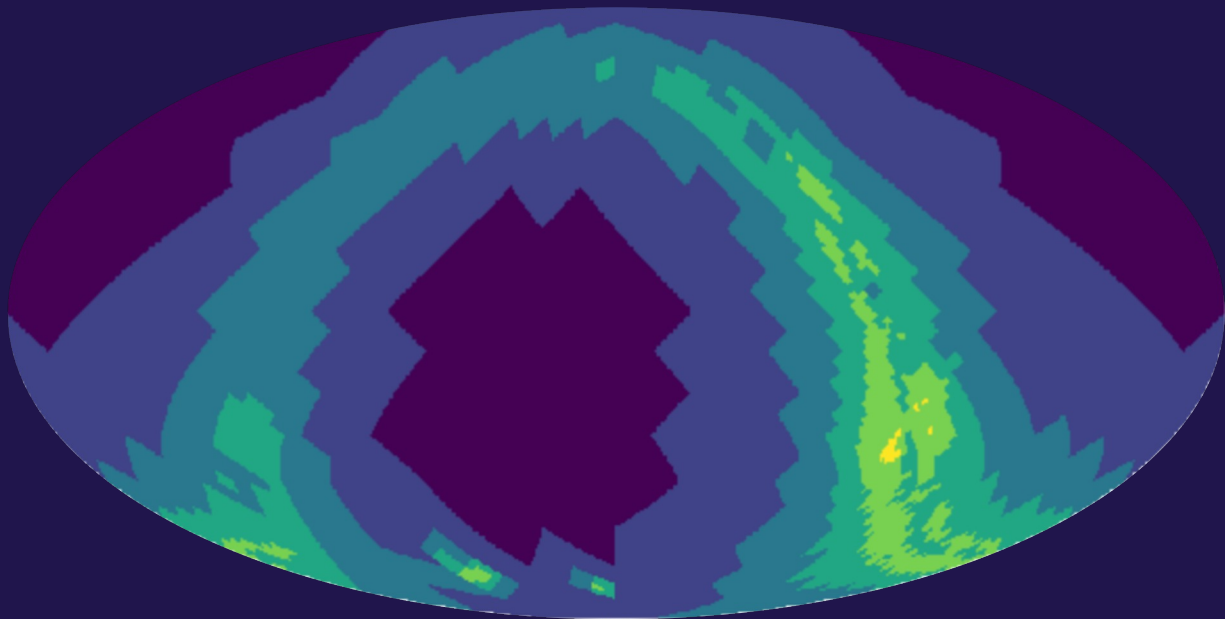
3993 partitions for Gaia DR3, with 1M  
object/partition threshold

## 2. On-disk organization

Holding everything in a single directory is unwieldy (at best).

A directory structure encoding the hierarchy would be helpful.

```
Norder0-Npix0.tsv.gz  
...  
Norder1-Npix28.tsv.gz  
Norder1-Npix29.tsv.gz  
Norder1-Npix30.tsv.gz  
Norder2-Npix112.tsv.gz  
Norder2-Npix113.tsv.gz  
Norder2-Npix114.tsv.gz  
Norder2-Npix115.tsv.gz  
...  
Norder0-Npix11.tsv.gz
```

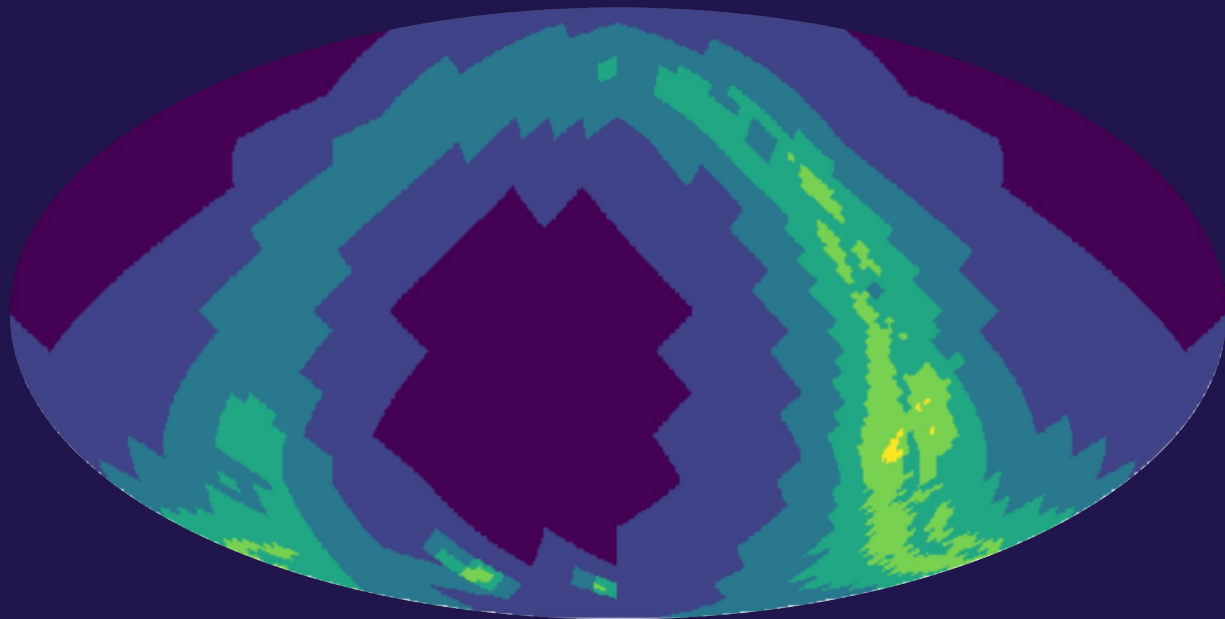




## 2. On-disk organization: HiPS-like Directories

Fortunately, we have a precedent in VO on how to organize hierarchically partitioned HEALPix data: [HiPS](#).

```
Norder0/Dir0/Npix0.tsv.gz  
...  
Norder1/Dir0/Npix28.tsv.gz  
Norder1/Dir0/Npix29.tsv.gz  
Norder1/Dir0/Npix30.tsv.gz  
Norder2/Dir0/Npix112.tsv.gz  
Norder2/Dir0/Npix113.tsv.gz  
Norder2/Dir0/Npix114.tsv.gz  
Norder2/Dir0/Npix115.tsv.gz  
...  
Norder0/Dir0/Npix11.tsv.gz
```

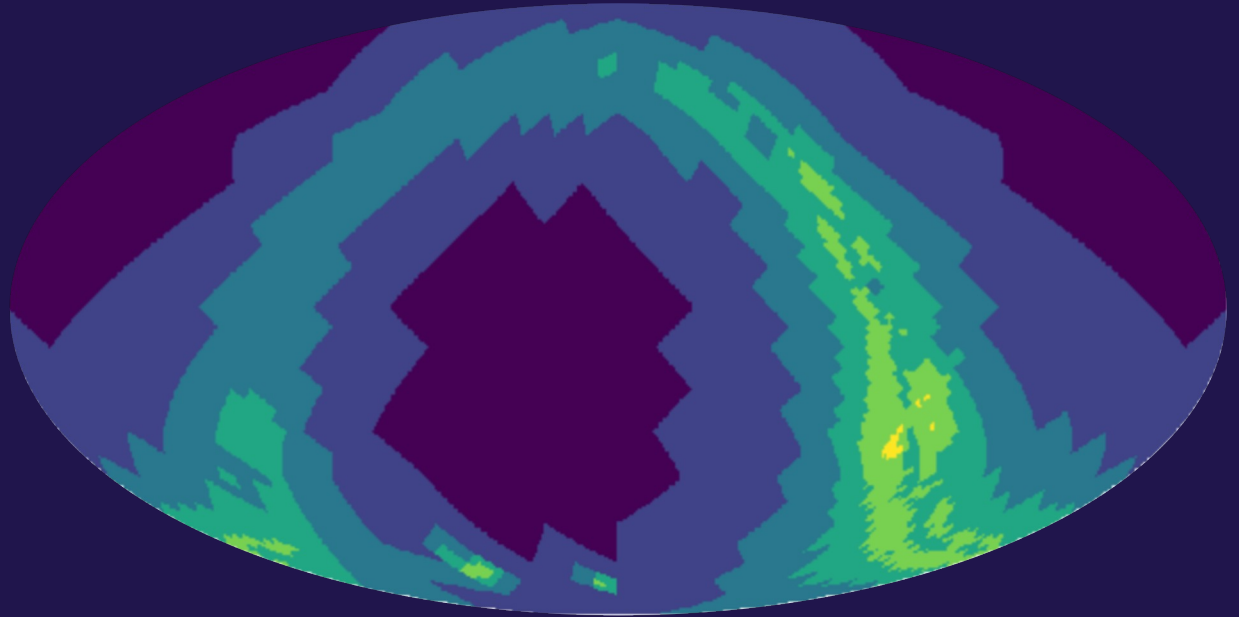


# Note: Hierarchical Progressive Surveys

While we use a HiPS-like directory structure, we only store data at the leaf nodes.

I.e., there are no lower-resolution files at lower orders.

```
Norder0/Dir0/Npix0.tsv.gz
...
Norder1/Dir0/Npix28.tsv.gz
Norder1/Dir0/Npix29.tsv.gz
Norder1/Dir0/Npix30.tsv.gz
Norder2/Dir0/Npix112.tsv.gz
Norder2/Dir0/Npix113.tsv.gz
Norder2/Dir0/Npix114.tsv.gz
Norder2/Dir0/Npix115.tsv.gz
...
Norder0/Dir0/Npix11.tsv.gz
```



e.g., there's no Norder1/Dir0/Npix31.csv

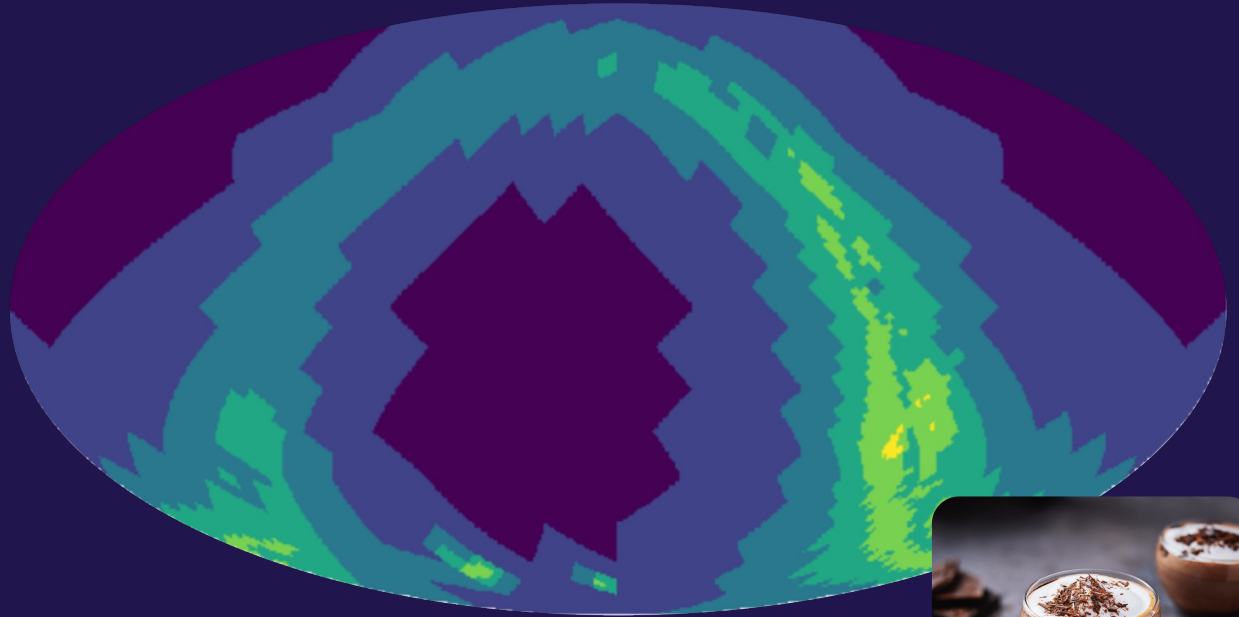
# Note: Hierarchical Progressive Surveys

While we use a HiPS-like directory structure, *we only store data at the leaf nodes.*

I.e., there are no lower-resolution files at lower orders.

```
Norder0/Dir0/Npix0.tsv.gz
...
Norder1/Dir0/Npix28.tsv.gz
Norder1/Dir0/Npix29.tsv.gz
Norder1/Dir0/Npix30.tsv.gz
Norder2/Dir0/Npix112.tsv.gz
Norder2/Dir0/Npix113.tsv.gz
Norder2/Dir0/Npix114.tsv.gz
Norder2/Dir0/Npix115.tsv.gz
...
Norder0/Dir0/Npix11.tsv.gz
```

This is really a MOC  
(multi-order coverage map; <https://ivoa.net/documents/MOC/>)



HiPSCat → MOCHA? (e.g. Multi-order Catalogs in HEALpix for Astronomy? Bacronym suggestions welcome!)





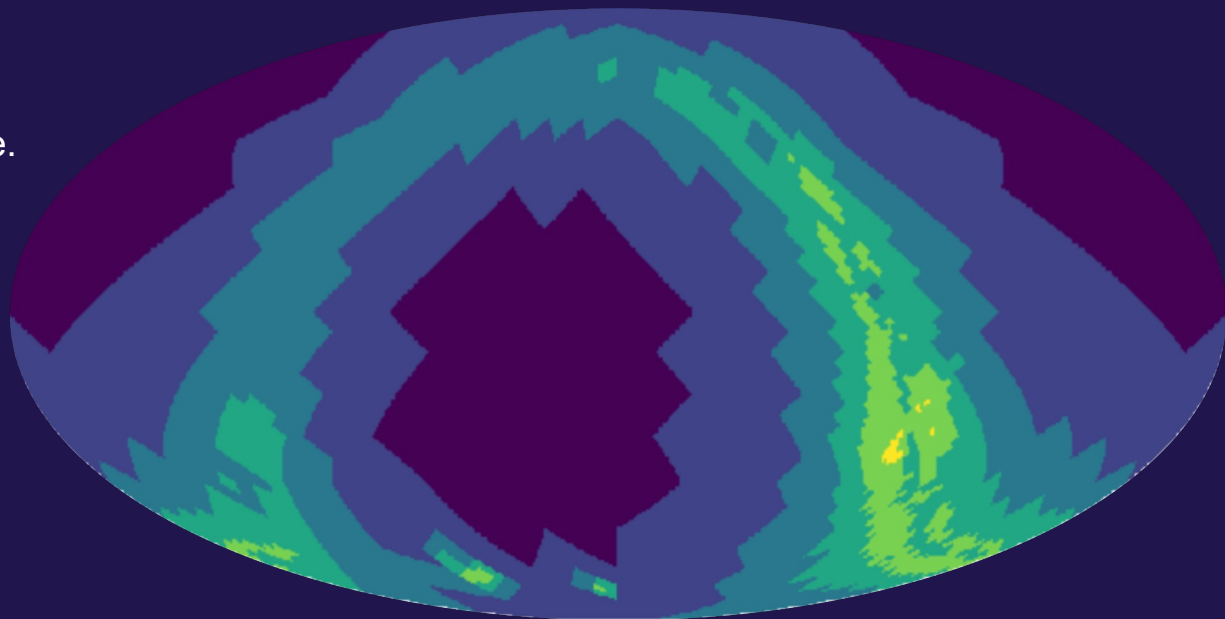
# 3. Serialization

TSV is not ideal for large catalog storage and analytics.

Time-consuming to parse and (de)compress. Also not seekable.

FITS? HDF5?

```
Norder0/Dir0/Npix0.tsv.gz
...
Norder1/Dir0/Npix28.tsv.gz
Norder1/Dir0/Npix29.tsv.gz
Norder1/Dir0/Npix30.tsv.gz
Norder2/Dir0/Npix112.tsv.gz
Norder2/Dir0/Npix113.tsv.gz
Norder2/Dir0/Npix114.tsv.gz
Norder2/Dir0/Npix115.tsv.gz
...
Norder0/Dir0/Npix11.tsv.gz
```



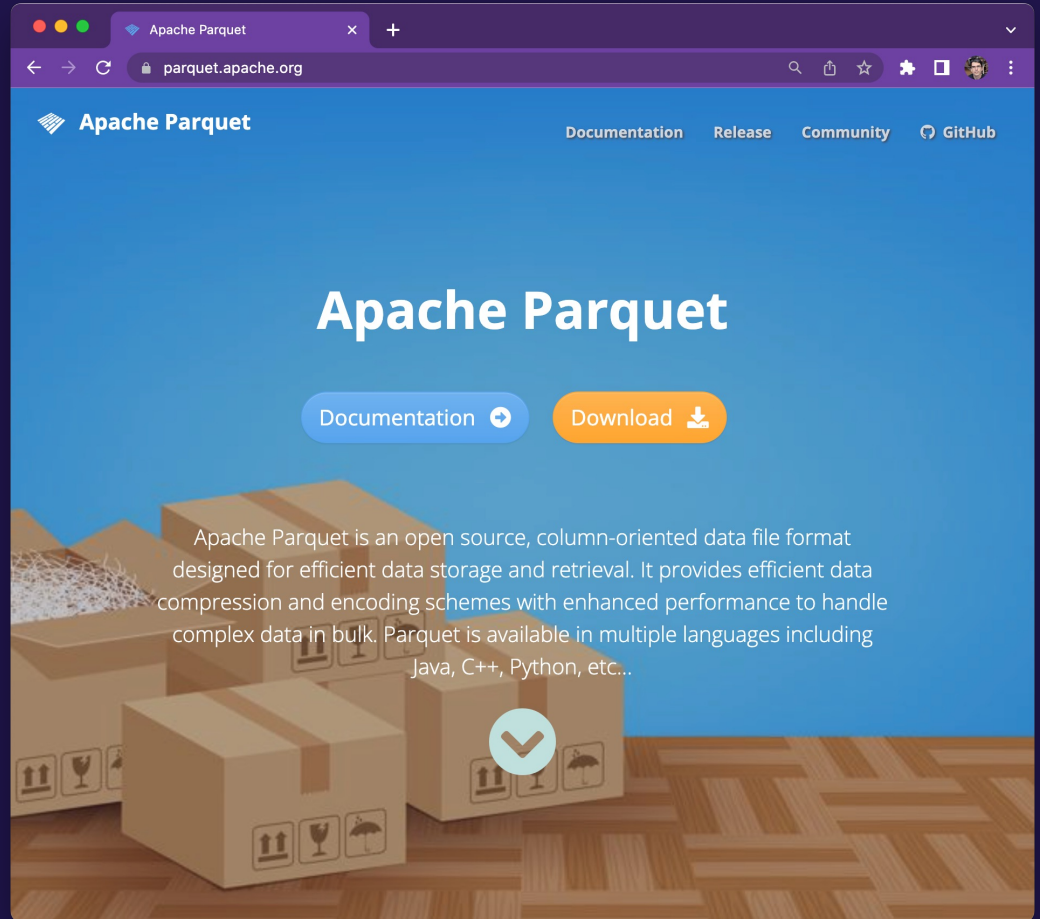
# 3. Serialization: Parquet

TSV is not ideal for large catalog storage and analytics.

We use Parquet.

## Key features:

- ✓ Designed for storage of large tables
- ✓ Columnar
- ✓ Efficient (binary)
- ✓ Transparent compression
- ✓ Data Integrity (checksums)
- ✓ Partitioning
- ✓ Broad multi-language support
- ✓ Broad tool support
- ✓ Strong industry backing
- ✓ Open source

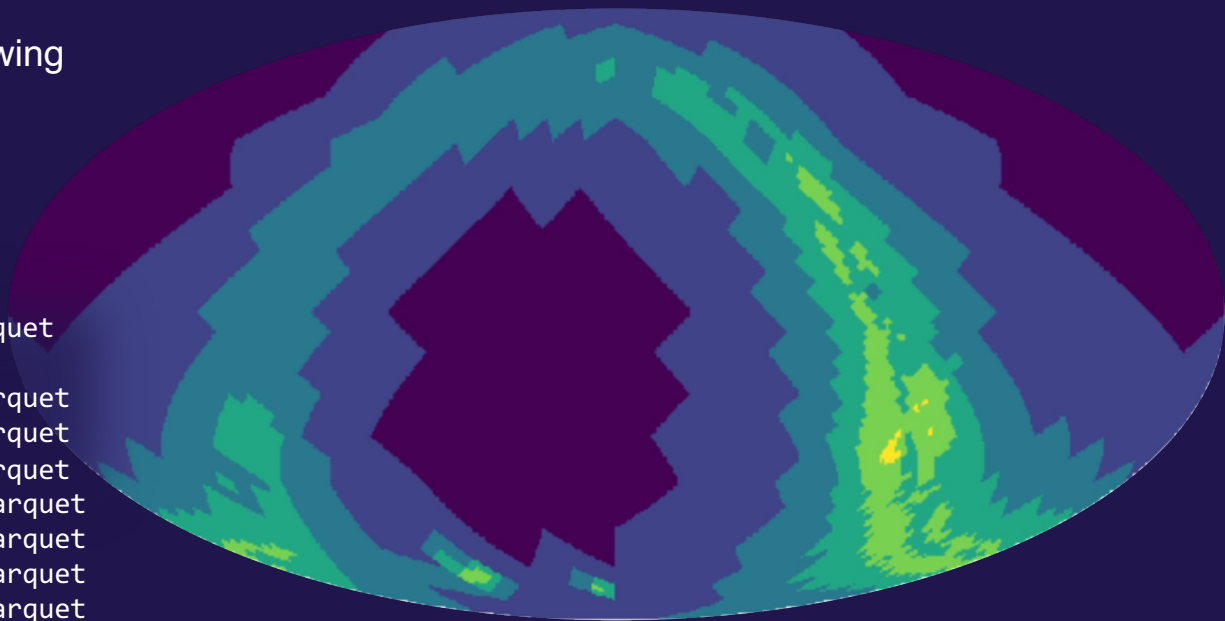


# 3. Serialization: Parquet

Parquet readers natively support reading partitioned datasets if they're stored in directories following `<key>=<value>` naming format.

We make that small tweak...

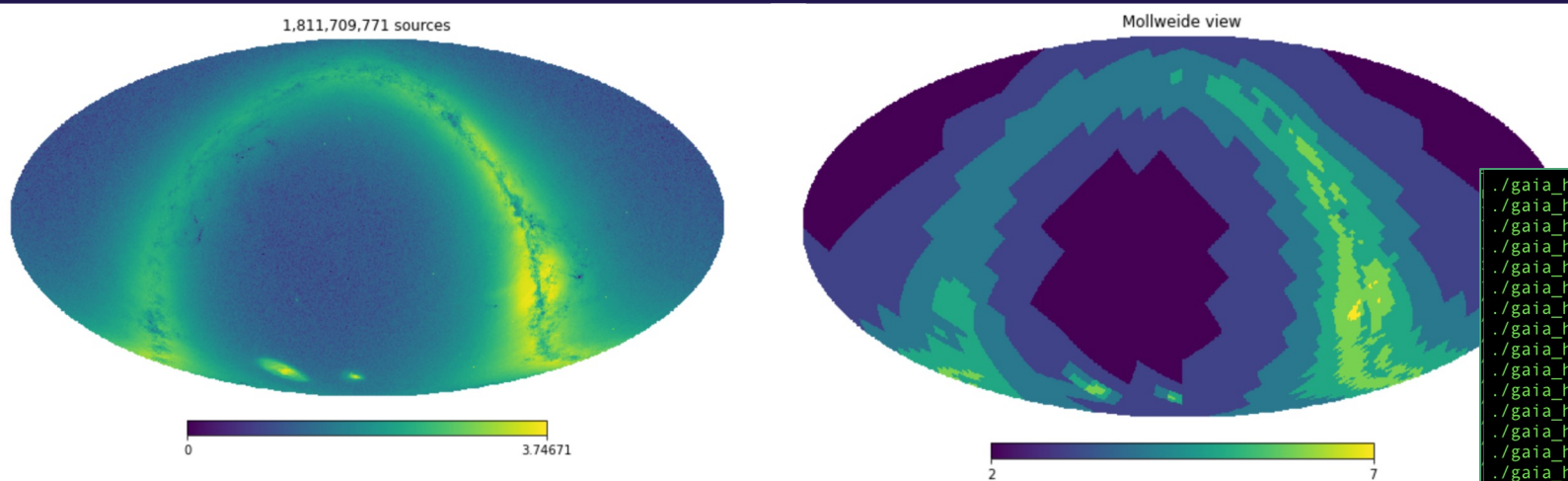
```
Norder=0/Dir=0/Npix=0/catalog.parquet
...
Norder=1/Dir=0/Npix=28/catalog.parquet
Norder=1/Dir=0/Npix=29/catalog.parquet
Norder=1/Dir=0/Npix=30/catalog.parquet
Norder=2/Dir=0/Npix=112/catalog.parquet
Norder=2/Dir=0/Npix=113/catalog.parquet
Norder=2/Dir=0/Npix=114/catalog.parquet
Norder=2/Dir=0/Npix=115/catalog.parquet
...
Norder=0/Dir=0/Npix=11/catalog.parquet
```





(\*) We may need to change the name to avoid confusion with true HiPS catalogs

# All together: HiPSCat\*



Layout on “disk”:

```
./gaia_hipscat/Norder5
./gaia_hipscat/Norder5/Npix0
./gaia_hipscat/Norder5/Npix1
./gaia_hipscat/Norder5/Npix10
./gaia_hipscat/Norder5/Npix100
./gaia_hipscat/Norder5/Npix1000
./gaia_hipscat/Norder6/Npix10000
./gaia_hipscat/Norder6/Npix10001
./gaia_hipscat/Norder6/Npix10002
./gaia_hipscat/Norder6/Npix10003
./gaia_hipscat/Norder5/Npix10004
./gaia_hipscat/Norder5/Npix10005
./gaia_hipscat/Norder5/Npix10006
./gaia_hipscat/Norder5/Npix10007
./gaia_hipscat/Norder5/Npix10008
./gaia_hipscat/Norder5/Npix10009
./gaia_hipscat/Norder5/Npix1001
./gaia_hipscat/Norder5/Npix10010
./gaia_hipscat/Norder5/Npix10011
./gaia_hipscat/Norder7/Npix10012
./gaia_hipscat/Norder7/Npix10013
./gaia_hipscat/Norder7/Npix10014
./gaia_hipscat/Norder7/Npix10015
./gaia_hipscat/Norder5/Npix10016
./gaia_hipscat/Norder5/Npix10017
./gaia_hipscat/Norder5/Npix10018
./gaia_hipscat/Norder5/Npix10019
```

Gaia DR2 Catalog Counts (log scale)

Visualization of file storage (color = healpix level)  
3933 partitions of similar size (128-256 MB)

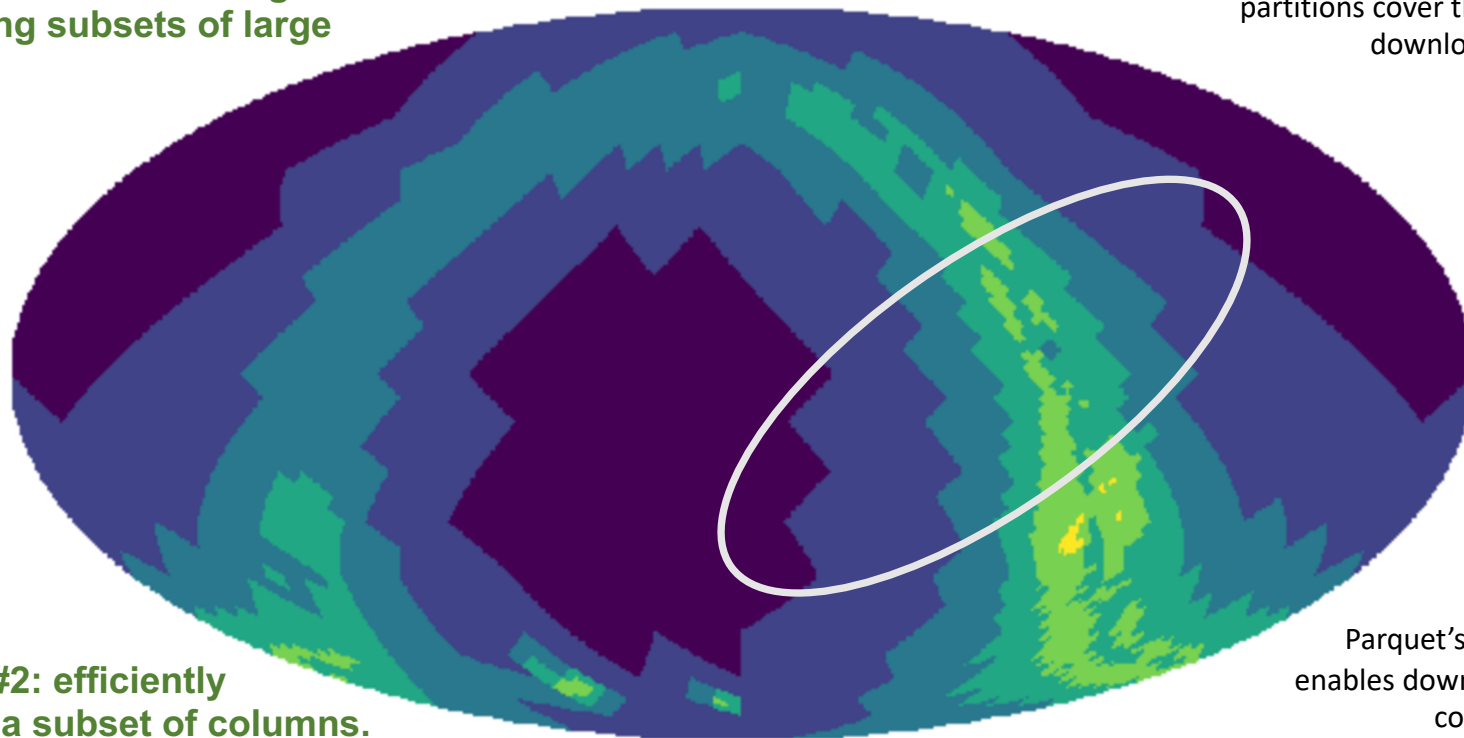
Inspired by the widely used IVOA HiPS standard. Standing on the shoulders of giants (Pierre, Mark, Thomas et al; REC-HIPS-1.0-20170519).

*What can we do with this?*

# Download spatial subsets of a catalog

**Use case #1: downloading overlapping subsets of large catalogs.**

Given a region of the sky, it's straightforward to find which partitions cover the region (and download those files)



**Use case #2: efficiently download a subset of columns.**

Parquet's columnar layout enables downloading only the columns of interest

order = 2  
pixel size size = 14.7deg

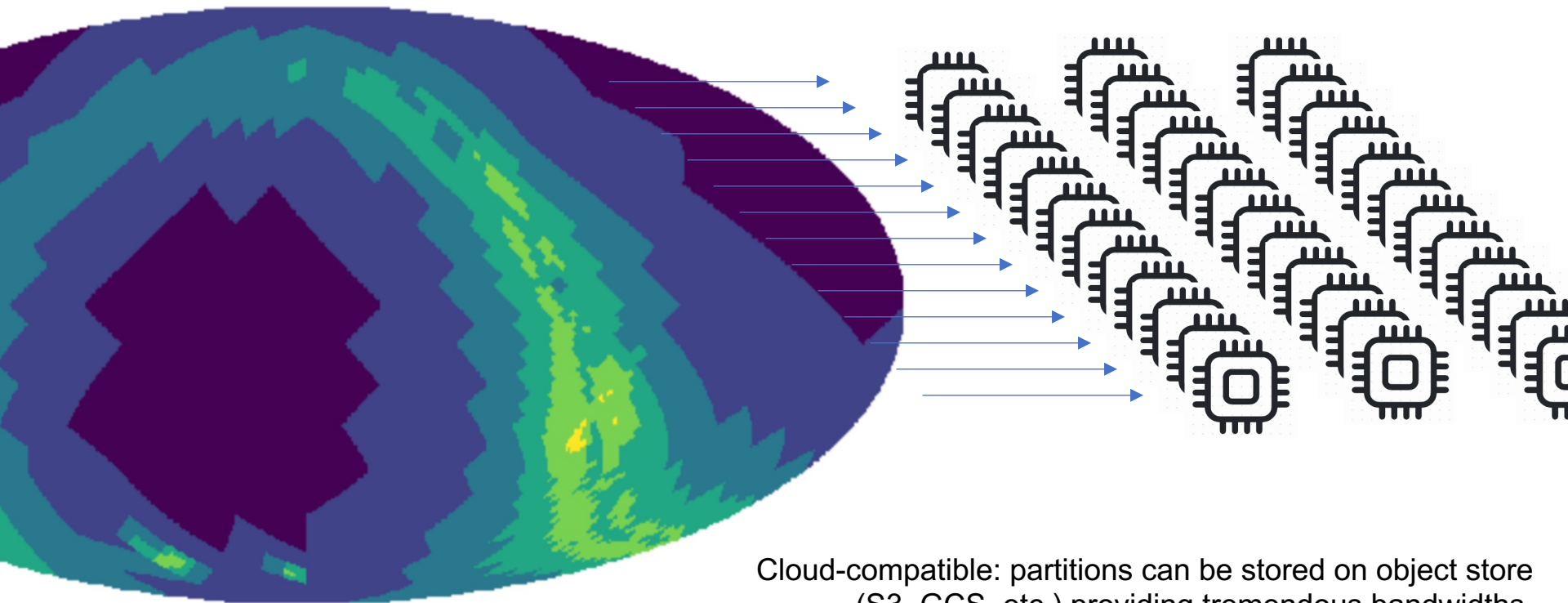


order = 7  
pixel size size = 0.46deg

# Straightforward Parallel Whole-Catalog Computation

Use case #3: complex searches, feature computation, spatial processing (clustering)

Enables very simple parallel computation schemes: per-file parallelization.

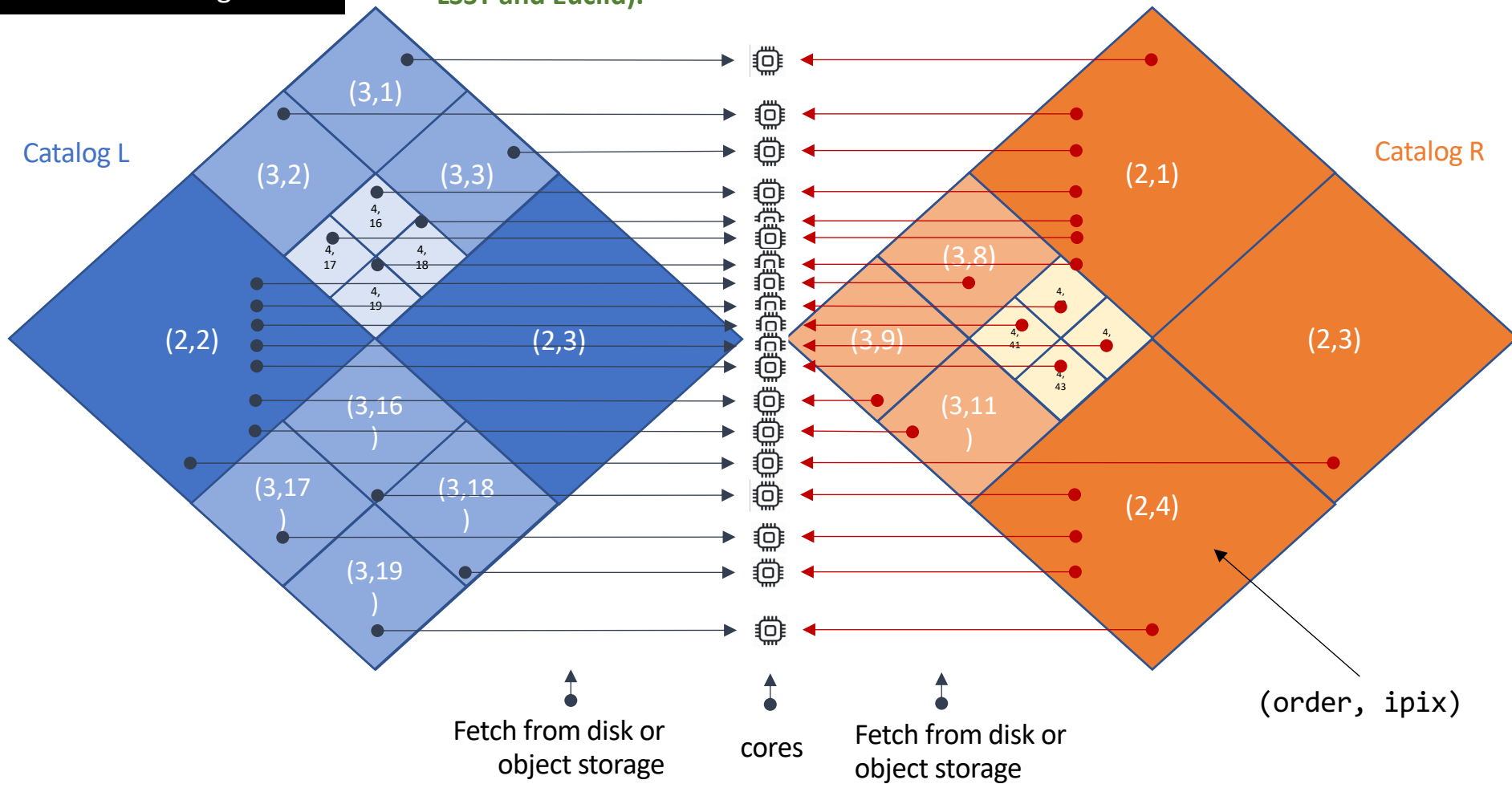


Cloud-compatible: partitions can be stored on object store (S3, GCS, etc.) providing tremendous bandwidths.



Efficient, parallel, joins and crossmatching

Use case #4: distributed analysis on data from two catalogs (example: LSST and Euclid).



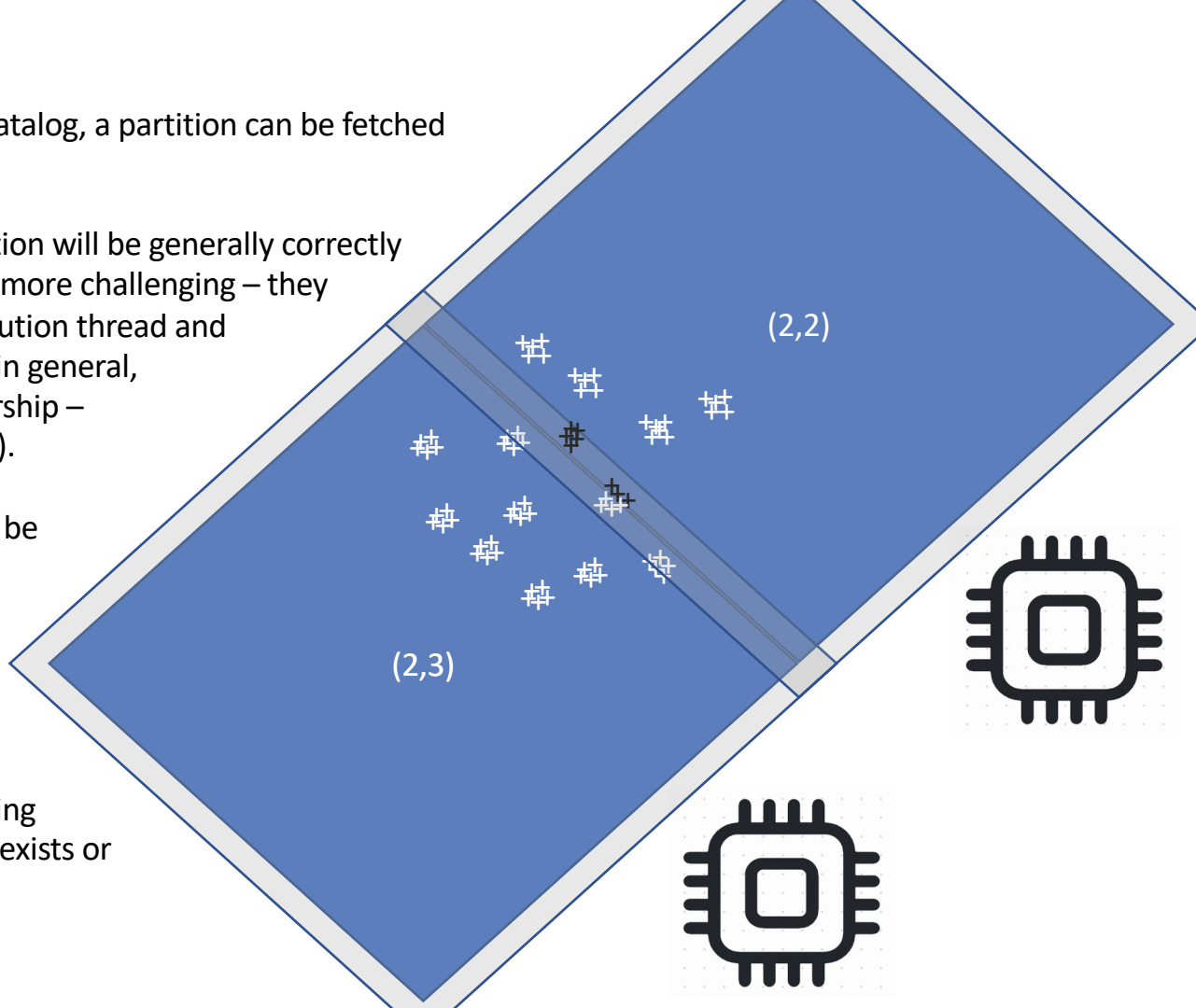
# Clustering

**Clustering algorithms:** Given a source catalog, a partition can be fetched and a clustering algorithm run on it.

Clusters close to the center of the partition will be generally correctly found. Clusters near the edges become more challenging – they can be found by more than one execution thread and will need to be deduplicated/resolved (in general, they won't even be identical in membership – distributed clustering is a *hard* problem).

Simple 1-pass clustering algorithms can be made to work, and complex algorithms can be implemented in two passes where the 2<sup>nd</sup> pass only requires  $O(\sim 1\%)$  of the data to exchange.

The use case here is creation of object catalog from source catalogs via clustering (of interest where no reference catalog exists or cannot be utilized for whatever reason)





# **Analytics Tools**

**How do we expose these capabilities to the user and enable science?**

# Any Tool that Understands Parquet can Read HiPSCat



- Parquet underpins much of modern large-dataset analytics ...
- ... and HiPSCat is a valid, Hive-partitioned, Parquet dataset.
  
- We can immediately use existing tools.
- Spark, Dask, Ray, Pandas\*, Hive, Hadoop, ....
  
- For all of these (except for Pandas) analysis execution is inherently distributed, with an API that hides the complexity from the user.
  
- Still, we would like an "astronomy-aware" layer...





# LSDB: Python Analytics for HiPSCat

- LSDB: Large Survey Database
- Enable Pandas-like analysis on trillions of observations with thousands of cores
- Build on existing tools: Dask (looking at Ray).
- Full HiPSCat awareness: spatial queries, cross-matching, timeseries, multi-dataset joining.
- Very much in pre-alpha/prototype phase; expect usable alphas in the next few weeks

```
img = gaia
    .query("pm > 10")
    .crossmatch(ztf)
    .join(ztf_sources)
    .for_each(varstar_classify)
    .query("pRRLy > 0.95")
    .skymap()

hp.mollview(img)
```

LSDB target APIs: The API center science. Multi-processing, autoscaling, fail-over, etc. are all implicit. Good user experience.

Wyatt et al. (2023)

<https://github.com/astronomy-commons/lldb>



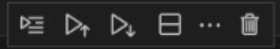
## *Quick Demo*

( see it in full on YouTube at <https://dirac.us/opj> )

```
import numpy as np
import healpy as hp
from lsdb import registry as reg
from dask.distributed import Client
```

[1] ✓ 4.2s

Python



```
#viewing imported hipscats
HIPS_DIR = '/data3/epyc/projects3/ivoa_demo/'
lsdb_registry = reg.Registry(location=HIPS_DIR)
```

[ ]

Python

```
#loading a hipscat
```

[ ]

Python

```
#basic catalog visualization
```

[ ]

Python

```
#print the gaia columns
```

[ ]

Python

( see it in full on YouTube at <https://dirac.us/opj> )

# Status

- Ad-hoc collaboration of scientists/engineers from LINCC, Rubin, MAST, IRSA, HESARC, LINeA (Brazil). *JOIN US!*
- Still heavily prototyping (the format isn't yet static)
- Expecting to stabilize the format a bit by ~June.
- Building tools in parallel (dataset import, end-user analytics)
- Will convert a large number of datasets, run science user tests.
- Gather feedback, iterate until we have a solid solution.
- Draft for winter Interop (?)

*Mailing list:* <https://groups.google.com/g/hipscat-wg>  
*Repositories:* <https://github.com/astronomy-commons>  
*Meetings:* 10am PT, June 2<sup>nd</sup> (then every other wk)



# Much work remains



- Format

- Better integrate with VO standards
- Supporting variable neighbor margins (for cross-matching)?
- Supporting efficient joins on tables of moving objects?
- Temporal partitioning?
- Supporting spectra?
- Future transaction support?

- Tooling

- Core HiPSCat libraries
- Developing LSDB
- Spark on HiPSCat
- Rust, Java, C/C++ libraries
- ...

- Science

- Importing a variety of catalogs
- Science case tests
- Deploy, collect user feedback
- ...

*Mailing list:* <https://groups.google.com/g/hipscat-wg>  
*Repositories:* <https://github.com/astronomy-commons>  
*Meetings:* 10am PT, June 2<sup>nd</sup> (then every other week)



*Collaboratively advancing data-intensive astronomy.*

A UNIVERSE UNDERSTOOD THROUGH  
DATA-INTENSIVE DISCOVERY

Thank You !

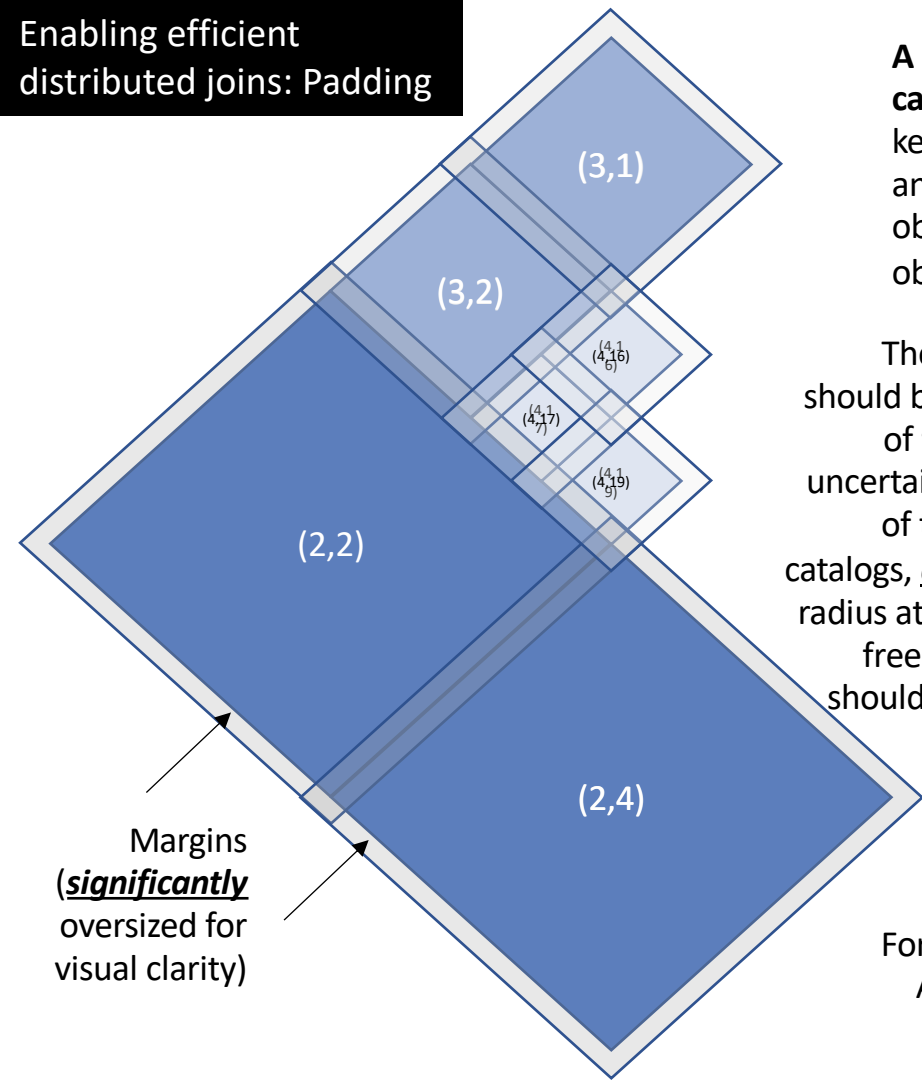
Contact: [mjuric@uw.edu](mailto:mjuric@uw.edu)

UNIVERSITY of WASHINGTON

# Backups

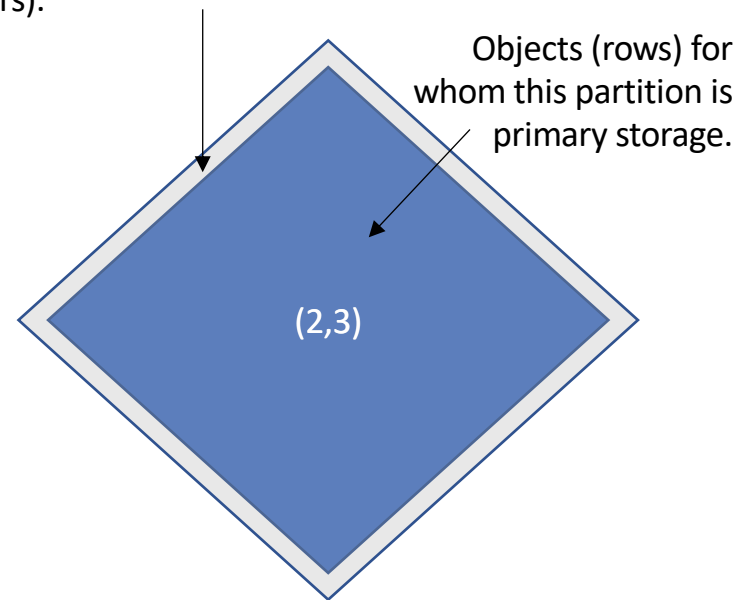


# Enabling efficient distributed joins: Padding



A “padding” or “overlap” or “neighbor margin” or “neighbor cache”. A copy of the rows whose  $(ra, dec)$  fall into this area will be kept with this partition. This enables distributed, shuffle-free, joins and crossmatches of catalogs where  $(ra, dec)$  of same objects/sources match only approximately (any catalog with observational errors).

The margin width should be on the order of  $\sim$ few times the uncertainty in  $(ra, dec)$  of the rows in the catalogs, or, close to the radius at which shuffle-free cross-matches should be supported.



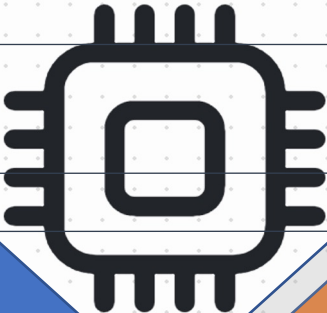
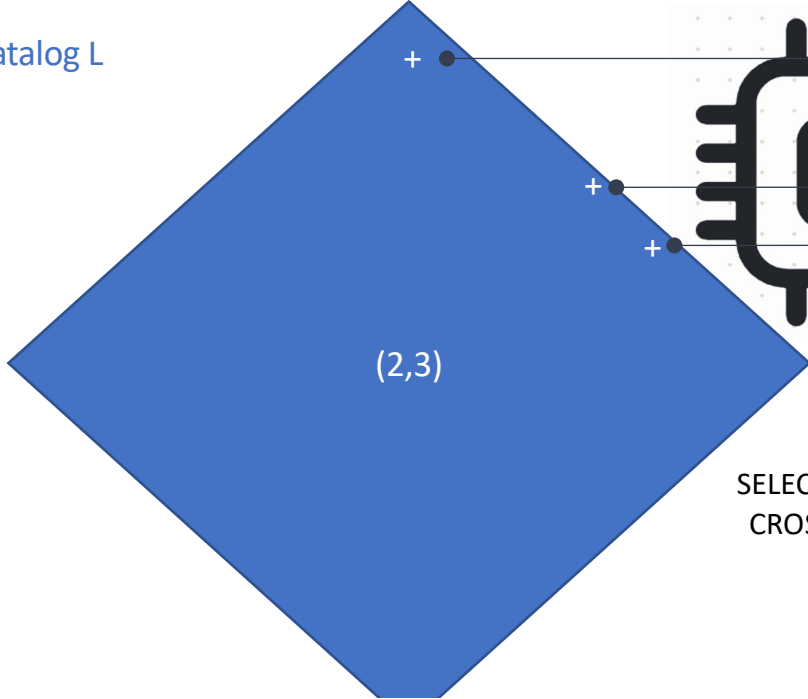
For a typical optical large dataset this may be **on order of  $\sim 10$ arcsec**. Assuming a (very) conservative medium partition size of  $\sim 0.46$ deg (1650 arcsec), **the extra storage overhead is roughly  $\sim 2.4\%$** .



# How distributed crossmatching works

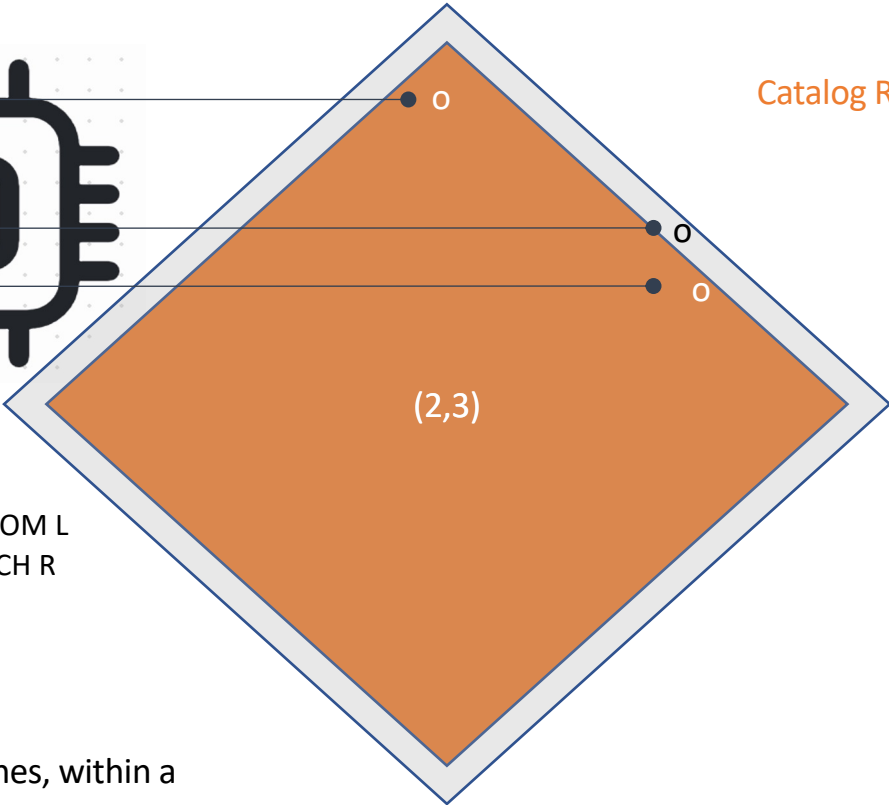
**Crossmatching algo:** Fetch the coordinate data from both partitions. Make sure to also download the padding of the partition to the right. For each row on the left, find nearest neighbor(s) in the table to the right. In some cases, the NN can be right across the partition boundary, and thanks to the padding, it will be found.

Catalog L



SELECT ... FROM L  
CROSSMATCH R

Catalog R

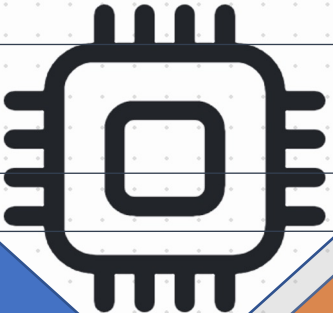
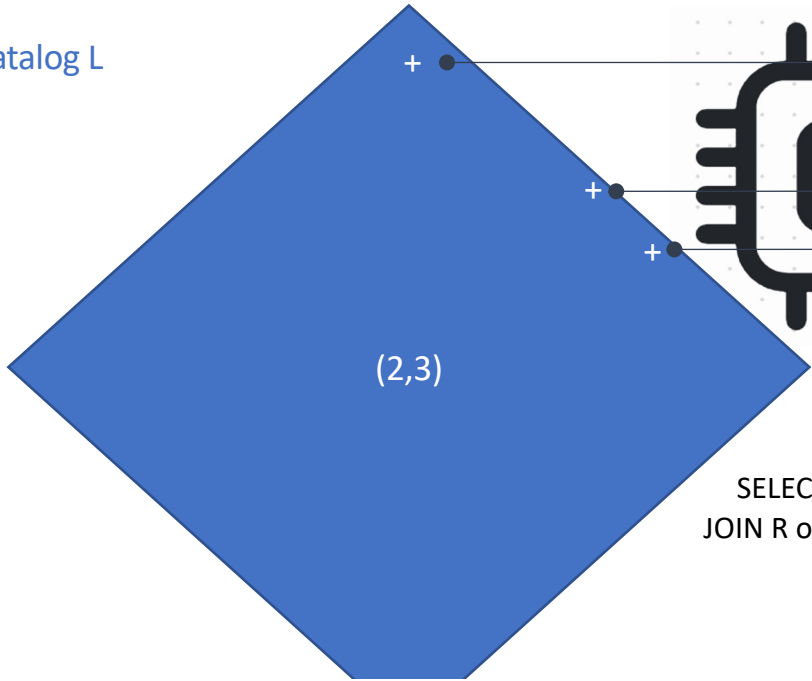


This allows for correct, parallel, N-nearest neighbor cross-matches, within a radius  $r$ , where  $r < \text{width of the padding margin}$ .

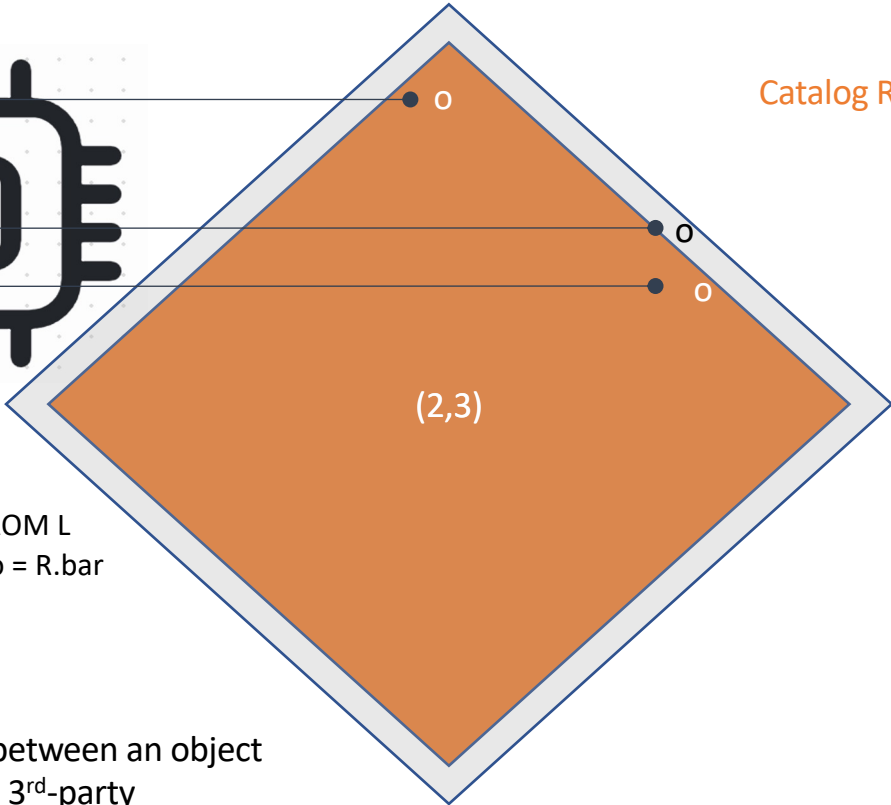
# How distributed joins work

Joining on **co-partitioned columns** just works™. Co-partitioned columns are columns for whose values it's guaranteed they will *all* be in the same partitions (plus margins) on the left and right side of the join. For such columns, a global "SELECT ... FROM L JOIN R ON L.foo = R.foo" type query can be executed on a per-partition basis.

Catalog L



Catalog R



SELECT ... FROM L  
JOIN R on L.foo = R.bar

This allows for correct, parallel, joins. Typical examples: a) a join between an object and a source catalog, on object ID; b) a 1:1 join of a catalog and a 3<sup>rd</sup>-party produced catalog of added-value columns (e.g., object classifications).