Active
Storage

Intro

Science
Context

Active
Storage
Concept
Chunking
Workflow
Detail
Implementation
Requirements

Summary

# The ExcaliData Implementation of Active Storage. An opportunity for Lustre?

ExcaliData is an Excalibur Cross Cutting Project
Excalibur is a UK exascale readiness programme
https://excalibur.ac.uk

Bryan Lawrence
(and many more, see next slide)

**National Centre for Atmospheric Science**
NATURAL ENVIRONMENT RESEARCH COUNCIL

**University of Reading**

ExCALIBUR
IO

University of Reading

Active
Storage

Intro

Science
Context

Active
Storage

Concept
Chunking
Workflow
Detail
Implementation
Requirements

Summary

# ExcaliData — Active Storage Authors and Partners

ExcaliData is a big project, active storage is one part of it, where the key participants are:

- Bryan Lawrence (UoR & NCAS)
- Jean-Thomas Acquaviva (DDN)
- Konstantinos Chasapis (DDN)
- Scott Davidson (StackHPC)
- Mark Goddard (StackHPC)
- David Hassell (UoR & NCAS)
- Grenville Lister (UoR & NCAS)
- Valeriu Predoi (UoR & NCAS)
- Matt Pryor (StackHPC)
- Stig Telfer (StackHPC)

**University of Reading**

**National Centre for Atmospheric Science**
NATURAL ENVIRONMENT RESEARCH COUNCIL

**ddn**

**StackHPC**

Funding via

**Met Office**

University of Reading

Co-Design: Not just about HPC Simulation, about analysis too!
There is an opportunity here for big science and big lustre!

University of
Reading

# We want to simulate our world

Active
Storage

University of
Reading

Image: from I. Lafaville, 2006

**Schematic for Global Atmospheric Model**

Horizontal Grid (Latitude-Longitude)

Vertical Grid (Height or Pressure)

Much complexity:

Active Storage

Intro

Science Context

Active Storage
- Concept
- Chunking
- Workflow
- Detail
- Implementation
- Requirements

Summary

University of Reading

Given knowledge of state of
- o(100) variables at every grid point for time $t$,
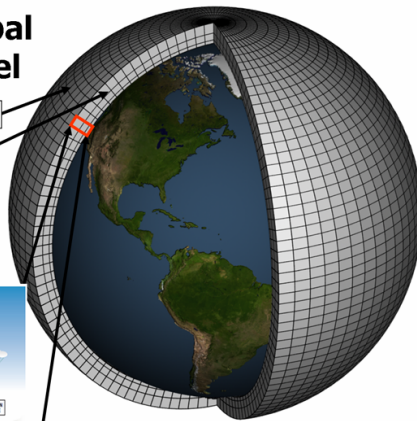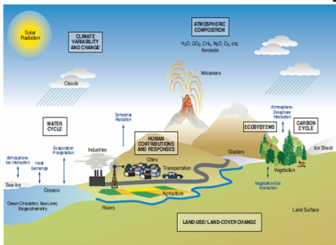**calculate**
at every grid point for every variable, state at $t + \Delta t$.

Many points, integrated for years with timestep of *o(minutes)*!

# The Evolution of Resolution: A better global microscope!



**300km**
**N48**
**(from <2000)**

**130km**
**N96**
**(from ~2002)**

**60km**
**N216**
**(from 2005)**

**25km**
**N512**
**(from 2012)**

**12km**
**N1024**
**(from 2013)**

Zooming in on
**Horizontal Model Resolution**
(These are global models,
this is just an illustration of resolution.)

University of
Reading

# extratropical cyclones

Wind speed and Sea Level Pressure (130 km)

Wind speed and Sea Level Pressure (25 km)

Accumulated precipitation (130 km)

Accumulated precipitation (25 km)

As simulated by the Met Office
https://uip.primavera-h2020.eu/storymaps/extra-tropical-cyclones

# Climate modelling: an exaflop and an exabyte challenge

Active
Storage

Intro
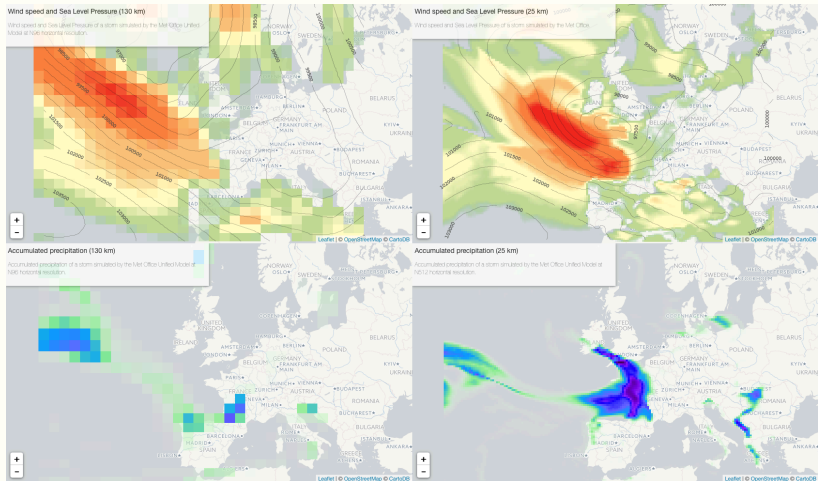
Science
Context

Active
Storage

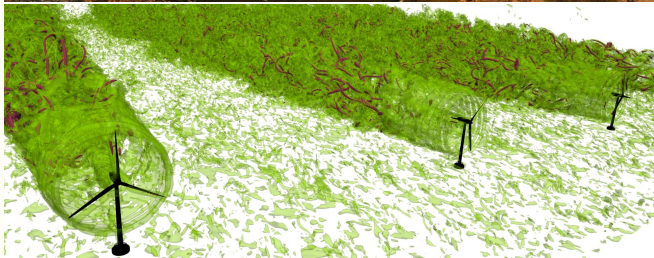Concept
Chunking
Workflow
Detail
Implementation
Requirements

Summary

# Excalidata and Active Storage are cross-cutting

Other big data challenges which should benefit from active storage:



Square Kilometre
Array (SKA),
Simulating turbulence,
and much more.

Picture Credits:

- SKA Project
- Imperial College

Active
Storage

- Exascale science requires access to increasingly large amounts of data
- These data might not be located near to the compute resources
- Moving data over the network, from storage to compute nodes, is costly
- The cost can be thought of as
    - Time or CPU cycles taken to move the data (particularly for read operations, which are always blocking, we can buffer or offload writes)
    - Network bandwidth used (and cost of ensuring the necessary fabric exists)
    - Energy consumed (to move the data)

- Can we avoid some of that movement?
    - We could, if we could do some computation in storage!

University of
Reading

- Nearly all storage systems now have lots of compute, it's necessary for error correction, and managing which bytes go where, but it's basically underused
  - JASMIN Quobyte storage system has (in 2021) 40PB usable storage powered by 131 compute nodes (3460 cores)
- For a long time now, people have been talking about utilising some of that underused compute, by doing "some compute tasks in storage".
- Previous ideas have included shipping functions, VMs or containers
  - If I ship a VM/container, how do I know it is not going to do bad things to the data and/or system?
  - How do I include complex systems in a workflow? How does this VM/container/function fit into my greater workflow?

University of
Reading

- Nearly all storage systems now have lots of compute, it's necessary for error correction, and managing which bytes go where, but it's basically underused
    - JASMIN Quobyte storage system has (in 2021) 40PB usable storage powered by 131 compute nodes (3460 cores)
- For a long time now, people have been talking about utilising some of that underused compute, by doing "some compute tasks in storage".
- Previous ideas have included shipping functions, VMs or containers
    - If I ship a VM/container, how do I know it is not going to do bad things to the data and/or system?
    - How do I include complex systems in a workflow? How does this VM/container/function fit into my greater workflow?
- We limit ourselves to reductions (a la MPI), and we let the Dask workflow tool handle the workflow (including identifying when we can use reductions).

University of Reading

- A task-based parallel computing library for Python.

- Dask partitions work into *computational chunks*.

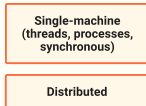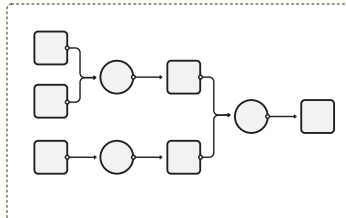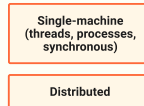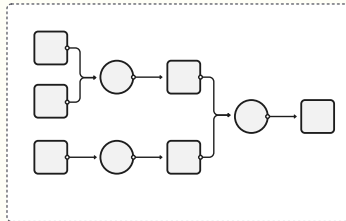**Collections** → **Task Graph** → **Schedulers**
(create task graphs)                (execute task graphs)



- Each *computational chunk* may map to multiple *storage chunks*, each of which will consist of a bunch of blocks in storage.

- Scheduler analyses workflow to create a task graph and allocates task to the available computing elements

**Collections** (create task graphs) → **Task Graph** → **Schedulers** (execute task graphs)

Dask Array
Dask DataFrame
Dask Bag
Dask Delayed
Futures

Single-machine (threads, processes, synchronous)
Distributed

- A task-based parallel computing library for Python.

- Dask partitions work into *computational chunks*.
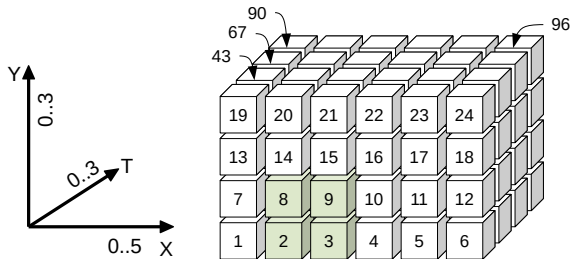
- Each *computational chunk* may map to multiple *storage chunks*, each of which will consist of a bunch of blocks in storage.

- Scheduler analyses workflow to create a task graph and allocates task to the available computing elements

- Objective: Manipulate the task graph to push some tasks into the storage!

**Active Storage**

University of Reading

Active
Storage

Intro

Science
Context

Active
Storage

Concept
Chunking
Workflow
Detail
Implementation
Requirements

Summary

Consider a simple $6 \times 4 \times 4$ grid:



- Data is written one chunk at a time (the size of a chunk is under user control). Deeper down the stack one chunk might be multiple blocks.
- For HDF/NetCDF/Zarr, data is also read into memory one chunk at a time.

Consider a simple $6 \times 4 \times 4$ grid:



- Data is written one chunk at a time (the size of a chunk is under user control). Deeper down the stack one chunk might be multiple blocks.
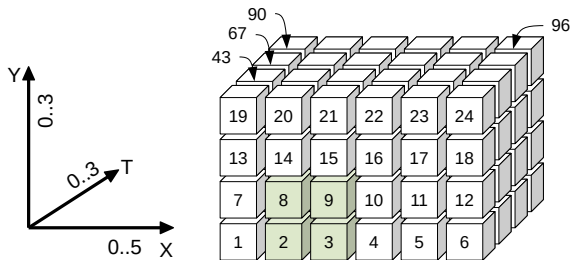- For HDF/NetCDF/Zarr, data is also read into memory one chunk at a time.
- Consider a chunk size of 6 and what happens if we we want to get a map (XY data) at a specific time (T) which corresponds to elements $(2, 3, 8, 9)$ – we need to read two chunks, and extract what we need.

- The default chunking has a preferred access. If you are reading in the direction you wrote, it will be efficient.
- But much of the time we don't do that.
- What about alternatives?

Active
Storage

Intro

Science
Context

Active
Storage

Concept
Chunking
Workflow
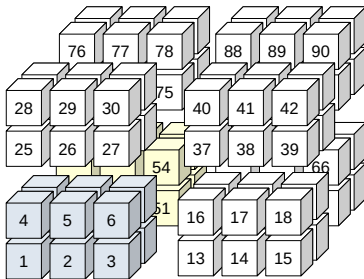Detail
Implementation
Requirements

Summary

# Storage Chunks (cont)

- The default chunking has a preferred access. If you are reading in the direction you wrote, it will be efficient.
- But much of the time we don't do that.
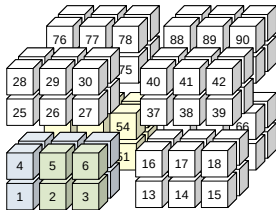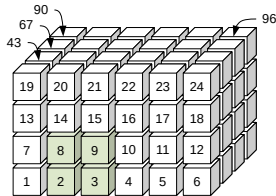- What about alternatives?



- Now nearly all ways of sampling into the cube along different routes than the original will be approximately equally efficient (but chunk size and dimensioning matter . . . a lot)!

Active Storage

- When we read data, we read data one chunk at a time
- If compressed, the chunk is uncompressed, then
- We extract the slice we want from the chunk.
- In the (green) map example, with the two toy chunking strategies shown, EITHER
    - read two chunks, decompress two chunks, extract two slices, build map, OR
    - read one chunk, decompress one chunk, extract one slice build map
- We might have other interesting chunk properties too: missing data, filters etc.



University of Reading

Active
Storage

Intro

Science
Context

Active
Storage

Concept
Chunking
Workflow
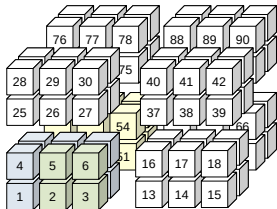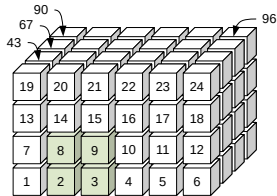Detail
Implementation
Requirements

Summary

- When we read data, we read data one chunk at a time
- If compressed, the chunk is uncompressed, then
- We extract the slice we want from the chunk.
- In the (green) map example, with the two toy chunking strategies shown, EITHER
    - read two chunks, decompress two chunks, extract two slices, build map, OR
    - read one chunk, decompress one chunk, extract one slice build map
- We might have other interesting chunk properties too: missing data, filters etc.
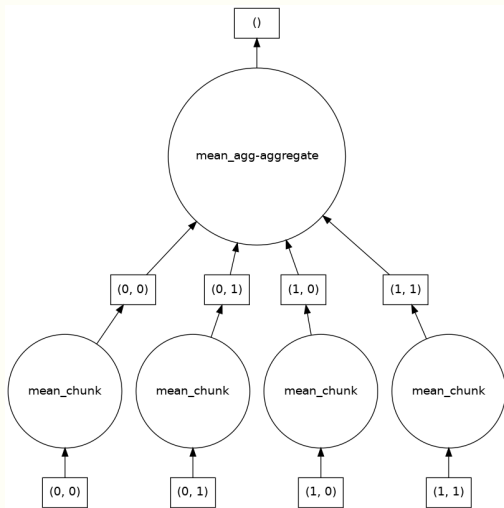- Now extend our thinking from "extracting data for a map" to "doing some calculation" . . .

Active
Storage

Intro

Science
Context

Active
Storage

Concept
Chunking
**Workflow**
Detail
Implementation
Requirements

Summary

# A simple reduction workflow: mean



Taking a mean using four computational chunks:

- four lots of data re read from storage
- four means are taken
- means are aggregated result is calculated from the partial means
- requires reading all the data from the storage system, moving it into the compute node(s).

  (see also Blelloch algorithm.)

Active
Storage

Intro

Science
Context

Active
Storage

Concept
Chunking
**Workflow**
Detail
Implementation
Requirements

Summary

(This time with three, and note direction of flow has changed)

Active
Storage

Intro

Science
Context

Active
Storage
Concept
Chunking
**Workflow**
Detail
Implementation
Requirements

Summary

# An active storage reduction workflow: mean

(This time with three, and note direction of flow has changed)



For, say a 3GB array, instead of reading 1 GB to each processor, we are reading a few bytes to each processor — much less data movement!

University of
Reading

Active
Storage

Intro

Science
Context

Active
Storage
Concept
Chunking
Workflow
Detail
Implementation
Requirements

Summary

University of
Reading

# Some gory details (1): the interface

Each computational chunk is operating on many storage chunks, but we need to know quite a lot about each storage chunk to do the calculations - knowledge held by the application, but not the storage, unless we tell it.

Active
Storage

## Vocabulary

**Storage Service (OSS)**: An (object) service which provides access to the storage system, takes data and writes it, or takes a read request and returns a series of blocks.

**Storage Target (OST)**: An (object) storage target used by an OSS or a kernel to actually write to one or more disks (e.g in a RAID system).

OSS receives and serves entire chunks:



This will work, OSS receive a series of blocks which constitute a chunk.
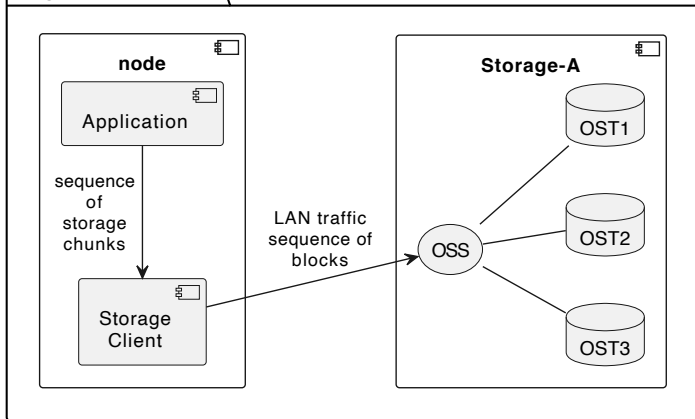
University of
Reading

Active
Storage

## Vocabulary

Storage Service (OSS):
An (object) service which
provides access to the
storage system, takes
data and writes it, or
takes a read request and
returns a series of
blocks.
Storage Target (OST): An
(object) storage target
used by an OSS or a
kernel to actually write to
one or more disks (e.g in
a RAID system).

The client is responsible for striping a chunk across
OSTs:



This will not be worth doing, a chunk is broken up across
OSTs, and needs to be re-assembled in the client before
reduction.

University of
Reading

Active
Storage

Intro

Science
Context

Active
Storage
Concept
Chunking
Workflow
Detail
Implementation
Requirements

Summary

University of
Reading

Reductionist: https://github.com/stackhpc/reductionist-rs

University of
Reading

Request:

```
POST /operations/{operation}
Headers:
  Authorization: Basic =auth_token=
  Content-Type: application/json
Body:
{
"source": "https://s3.server.address/,
"bucket": "my-bucket",
"object": "path/to/object",
"dtype": "int32",
"byte_order": "little",
"offset": 0,
"size": 128,
"shape": [20, 5],
"order": "C",
"selection": [[0, 19, 2], [1, 3, 1]],
"compression": {"id": "zlib"},
"filters": [{"id": "shuffle", "element_size"
    : 4}],
"missing": {"missing_value": 42}
}
```

Response:

```
HTTP/1.1 200 OK
Headers:
  Content-Length: 4
  Content-Type: application/octet-stream
  x-activestorage-dtype: int32
  x-activestorage-byte-order: little
  x-activestorage-shape: []
  x-activestorage-count: 1
Body:
42
```

Active
Storage

Intro

Science
Context

Active
Storage
  Concept
  Chunking
  Workflow
  Detail
  Implementation
  Requirements

Summary

University of
Reading
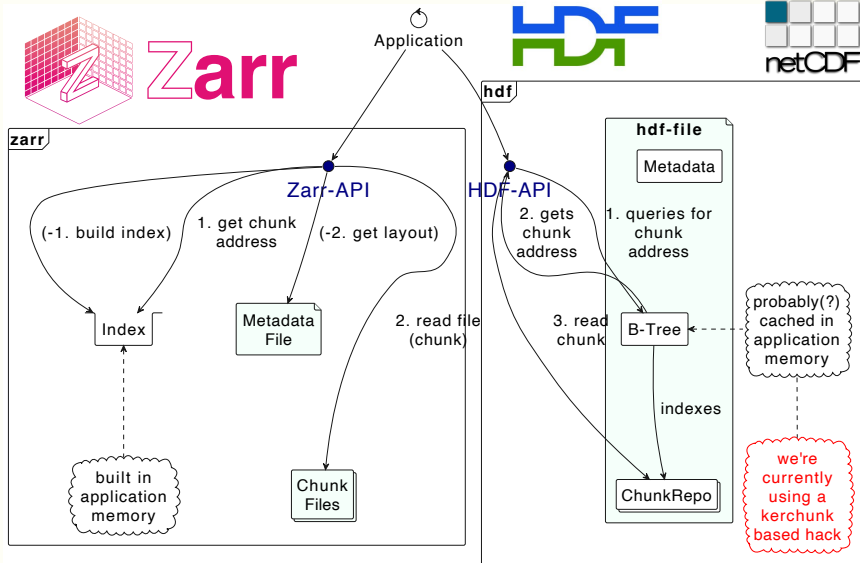
# Implementing Active Storage — Reductionist HTTP API

Request:

```
POST /operations/{operation}
Headers:
  Authorization: Basic =auth_token=
  Content-Type: application/json
Body:
{
"source": "https://s3.server.address/,
"bucket": "my-bucket",
"object": "path/to/object",
"dtype": "int32",
"byte_order": "little",
"offset": 0,
"size": 128,
"shape": [20, 5],
"order": "C",
"selection": [[0, 19, 2], [1, 3, 1]],
"compression": {"id": "zlib"},
"filters": [{"id": "shuffle", "element_size"
     : 4}],
"missing": {"missing_value": 42}
}
```

Response:

```
HTTP/1.1 200 OK
Headers:
  Content-Length: 4
  Content-Type: application/octet-stream
  x-activestorage-dtype: int32
  x-activestorage-byte-order: little
  x-activestorage-shape: []
  x-activestorage-count: 1
Body:
42
```

How do we find the offset for the
chunk operation?

Active
Storage

Active
Storage

**Single Server Pattern**

**node**

Application

sequence
of
storage
chunks

Storage
Client

traffic
sequence of
blocks

**Storage-A**

OST1

OSS

OST2

OST3

- The big question is how to get the request through the storage stack?

University of
Reading

- The big question is how to get the request through the storage stack?
- Solution: Hack `ioctl`.
- (Implemented first in Fuse for proof of concept, then in DDN Infinia)

# HDF5 files, chunks and striping

Metadata

Dataset array data

HDF5 File

Metadata is mixed with raw data in HDF5 file

2K metadata block; may be partially filled

HDF5 File

Metadata blocks of different lengths



Multiple OSS Pattern

- Reading: chunk locations are known to HDF-library, and can be extracted by application and active storage (a la kerchunk), but
- Writing: File system doesn't know how to align them with OSS when writing them.

Much easier to handle chunk alignment with Zarr!

University of Reading

- Can address a byte range within the file (a chunk) and fill a buffer with those bytes
    - (return error if buffer can't support decompressed chunk.)
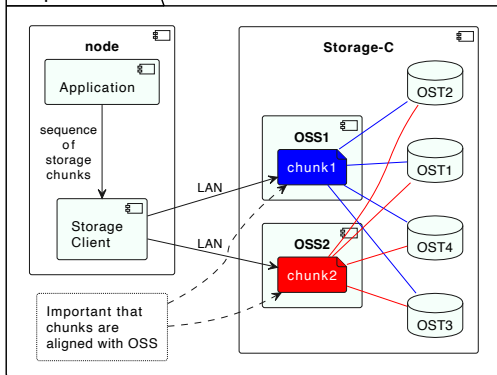    - (respect float description and return error if unrecognised.)
- Support for specific filter and decompression algorithms (those supported by Zarr and netcdf).
    - (Return error for unrecognised filter and/or decompression.)
- Support for the defined operations: mean, sum, max, min, count (more could be added provided they were reductions); applied to a slice within the buffered array and respecting missing data.
    - (Return error for unrecognised operation.)
    - (Return error if slice operation is out of bounds or misunderstood.)
    - (Return error for unsupported missing description.)

(We have discussed a streaming version of this as well, but this is the simplest version.)

University of
Reading

- Science Application (untouched!): any script written in Python, utilising

- CF-Python (climate forecasting domain specific analysis library), which itself uses

- Dask (provides flexible task-based parallelism from threads to nodes).

In the current implementation we have modified CF-Python to patch the Dask graph when it recognises active storage to use the

- PyActiveStorage middleware (handles computational offload to storage and returns results).

to return the first computational reduction.
(Rather than what would otherwise be at the bottom of the graph: a read and then the computation).

University of Reading

- Science Application (untouched!): any script written in Python, utilising

- CF-Python (climate forecast... ...n specific

- ...allelism from

In th... ...ified
CF-P... ...ecognises
activ...

It is possible to handle it all in Dask without modification further up the stack.
We have demonstrated that in a fork of Dask, but we wont try to push our patches upstream until there are some active storage servers deployed for real.

- ...omputational

...
to re... ...ion.
(Rath... ...an what would otherwise be at the bottom of the graph: a read and then the computation).

Active Storage

## Active Storage

- We've built and tested code to push filters, decompression, and reductions into storage.
    - S3 implementation is feature complete. Performance testing planned for early next year.
    - POSIX implementation in DDN Infinia shows promise.
- Existing version is domain specific (utilising CF-Python) but methodology to make discipline independent (via minor modification to Dask) exists.

University of Reading

Active
Storage

Intro

Science
Context

Active
Storage

Concept
Chunking
Workflow
Detail
Implementation
Requirements

Summary

University of
Reading

## Active Storage

- We've built and tested code to push filters, decompression, and reductions into storage.
    - S3 implementation is feature complete. Performance testing planned for early next year.
    - POSIX implementation in DDN Infinia shows promise.
- Existing version is domain specific (utilising CF-Python) but methodology to make discipline independent (via minor modification to Dask) exists.

## Whither Lustre?

- How do we get chunks aligned on OSS?
- Do we need striping if we are putting chunks on different OSS?
- How do we add computational facilities to Lustre?
    - Control?: Overload `ioctl`? Something else?
    - Software: Need implementation of reductions, filters and decompression.
    - Hardware: Buffer size, caching, configuration?