

RFC: Data Analysis Extensions

Jerome Soumagne
Quincey Koziol

Accessing, selecting and retrieving data from an HDF5 container can be a time consuming process, particularly so when data is very large. To enable, ease and accelerate this process, we introduce in this RFC extensions to the library to efficiently query, select and index data.

1 Introduction

When working on large datasets, finding and selecting the interesting pieces of the data can be a cumbersome process. Currently, the HDF5 library enables the application developer to select, read and write data but does not provide any mechanism to select and retrieve pieces without prior knowledge of their content, or without the developer to provide the exact data coordinates that he is willing to access. To satisfy that need, one must be able to issue queries by specifying a data selection criteria. These queries, when applied to the data, can then generate a selection, which contains the coordinates that satisfy the query—this may imply accessing the data and selecting it, depending on the query condition satisfaction. To accelerate and facilitate this process (i.e., so that the data no longer needs to be directly accessed), one can use indexing techniques, which consist of generating an index and using that index to answer the query selection criteria and find the matching elements.

We define in this RFC the components that can enable application developers to create complex and high-performance queries on both metadata and data elements within an HDF5 container and retrieve the results of applying those query operations. Support for these operations can be defined via:

- New *query* objects¹ and API routines, enabling the construction of query requests for execution on HDF5 containers;
- New *index* objects and API routines, which allows the creation of indices on the contents of HDF5 containers, to improve query performance.

2 Query Objects

Query objects are the foundation of the data analysis operations in HDF5 and can be built up from simple components in a programmatic way to create complex operations using *Boolean* operations. The current API is presented below:

```
1 | hid_t H5Qcreate(H5Q_type_t query_type, H5Q_match_op_t match_op, ...);  
2 | herr_t H5Qclose(hid_t query_id);
```

¹Query objects are *in-memory* objects, which therefore do not modify the content of the container.

```

3 | hid_t H5Qcombine(hid_t query1_id, H5Q_combine_op_t combine_op, hid_t query2_id);
4 | herr_t H5Qget_type(hid_t query_id, H5Q_type_t *query_type);
5 | herr_t H5Qget_match_op(hid_t query_id, H5Q_match_op_t *match_op);
6 | herr_t H5Qget_components(hid_t query_id, hid_t *subquery1_id, hid_t *subquery2_id);
7 | herr_t H5Qget_combine_op(hid_t query_id, H5Q_combine_op_t *op_type);
8 | herr_t H5Qencode(hid_t query_id, void *buf, size_t *nalloc);
9 | hid_t H5Qdecode(const void *buf);

```

2.1 Query Creation

The core query API is composed of two routines: `H5Qcreate` and `H5Qcombine`. `H5Qcreate` creates new queries, by specifying an aspect of an HDF5 container, such as:

- `H5Q_TYPE_DATA_ELEM` (data element);
- `H5Q_TYPE_LINK_NAME` (link name);
- `H5Q_TYPE_ATTR_NAME` (attribute name);
- `H5Q_TYPE_ATTR_VALUE` (attribute value);

as well as a match operator, such as:

- `H5Q_MATCH_EQUAL` (`=`);
- `H5Q_MATCH_NOT_EQUAL` (`≠`);
- `H5Q_MATCH_LESS_THAN` (`≤`);
- `H5Q_MATCH_GREATER_THAN` (`≥`);

and a value for the match operator. Created query objects can be serialized and deserialized using the `H5Qencode` and `H5Qdecode` routines², and their content can be retrieved using the corresponding accessor routines. `H5Qcombine` combines two query objects into a new query object, using boolean operators such as:

- `H5Q_COMBINE_AND` (`^`);
- `H5Q_COMBINE_OR` (`∨`).

Queries created with `H5Qcombine` can be used as input to further calls to `H5Qcombine`, creating more complex queries.

For example, a single call to `H5Qcreate` could create a query object that would match data elements in any dataset within the container that are equal to the value “17”. Another call to `H5Qcreate` could create a query object that would match link names equal to “Pressure”. Calling `H5Qcombine` with the `^` operator and those two query objects would create a new query object that matched elements equal to “17” in HDF5 datasets with link names equal to “Pressure”. Creating the data analysis extensions to HDF5 using a *programmatic interface* for defining queries avoids defining a text-based query language as a core component of the data analysis interface, and is more in keeping with the design and level of abstraction of the HDF5 API. The HDF5 data model is more complex than traditional database tables and a simpler query model would likely not be able to express the kinds of queries needed to extract the full set of components of an HDF5 container. A text-based query language (or GUI) could certainly be built on top of the query API defined here to provide a more user-friendly (as opposed to *developer-friendly*) query syntax like “Pressure = 17”. However, we regard this as out-of-scope for now.

²Serialization/deserialization of queries were introduced so that queries can be sent through the network.

2.2 Query Execution and Views

Applying a query to an HDF5 container creates an HDF5 *view*. A view is an anonymous HDF5 group in a temporary memory container. A view group contains up to three 1-D dataset objects with references into the contents of the HDF5 container that the query was applied to: one with *object* references, one with *region* references and one with *attribute* references.

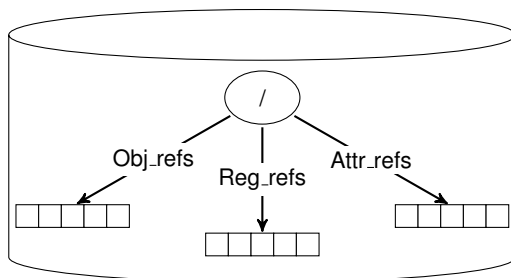


Figure 1 – Temporary memory container with datasets containing query results (i.e., HDF5 references).

A new view is created by the following `H5Qapply` routine, which applies a query to an HDF5 container, group hierarchy, or individual object and returns the object ID of the newly created group. The attributes, objects, and/or data regions referenced within a view’s datasets can be retrieved by further HDF5 dataset (`H5D*`) API calls.

```
|| hid_t H5Qapply(hid_t loc_id, hid_t query_id, unsigned *result, hid_t vcpl_id);
```

Although the created view is stored in memory, it can be persisted by calling `H5Ocopy()` to copy the group (stored in a virtual container) to a persistent container. For coherency, a time stamp may also be attached to it so that its state has a meaning compared to the state of the container that the query was applied to (as the container may have been modified in the meantime). It may be useful in that case to define different states to the view (*dead* or *live*) so that the user knows whether the view is current or not.

For example, starting with the HDF5 container described in Figure 2a, applying a “`link_name = Pressure`” query would result in the view shown in Figure 2b, highlighted in blue.

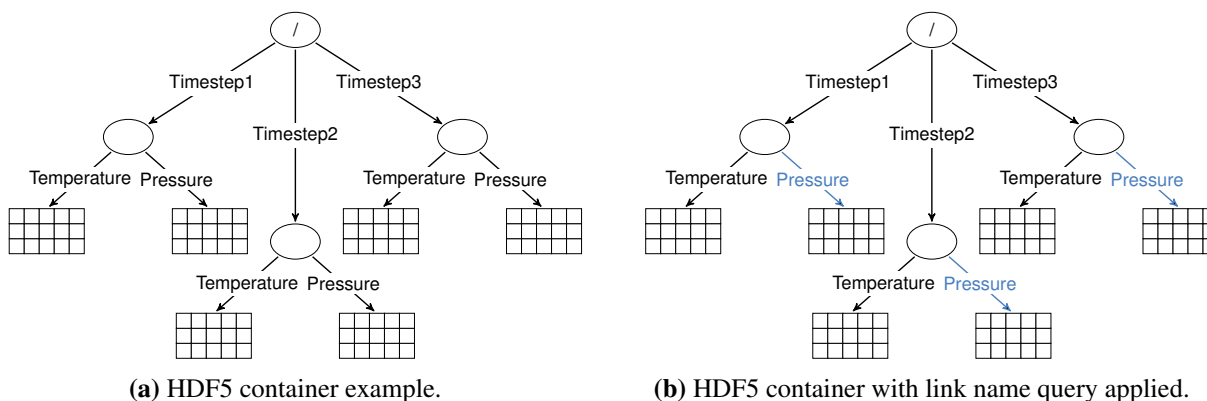


Figure 2 – HDF5 view creation example with simple query.

Alternatively, applying a “data_element = 17” query would result in the view shown in Figure 3a, highlighted in blue. Finally, applying the combined queries “(link_name = Pressure) \wedge (data_element = 17)” would result in the view shown in Figure 3b, highlighted in blue.

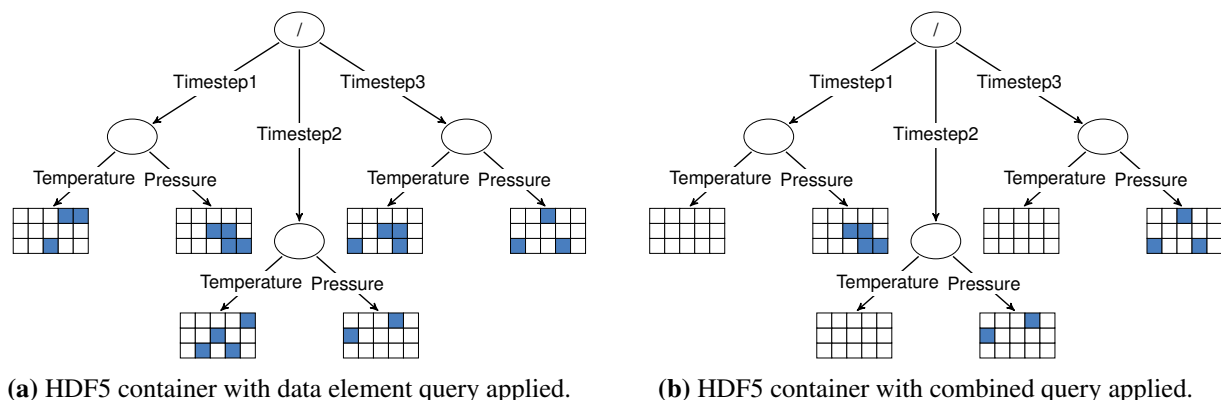


Figure 3 – HDF5 view creation example with combined query.

Since views contain a set of HDF5 references (object, dataset region or attribute references) to components of the underlying container, they retain the context of the original container. For example, the view containing the results in Figure 3b will contain three dataset region references, which can be retrieved from the view object and probed for the dataset and selection containing the elements that match the query with the existing `H5Rdereference` and `H5Rget_region` API calls. Note that selections returned from a region reference retain the underlying dataset’s dimensionality and coordinates—they are not *flattened* into a 1-D series of elements. The selection returned from a region reference can also be applied to a different dataset in the container, allowing a query on pressure values to be used to extract temperature values, for example.

Table 1 describes the result types for atomic queries and combining queries of different types. Query results of *None* type are rejected when `H5Qcombine` is called, causing it to return failure³.

3 Index Objects

Index objects are designed to accelerate creation of views from query operations. For example, if the “(link_name = Pressure) \wedge (data_element = 17)” query (previously described) is going to be either frequently executed or executed on a container that contains significantly large datasets, indices could be created in that container, which would speed up the creation of views when querying for link names and for data element values. Indices created for accelerating the `link_name = Pressure` or `data_element = 17` queries would also improve view creation for the more complex `(link_name = Pressure) \wedge (data_element = 17)` query. Although creating indices for metadata components of queries, such as link or attribute names, is possible, we first focus on index creation for dataset elements, as they represent the largest volume of data in typical HPC application usage of HDF5. Queries with metadata components execute properly, but are not able to be accelerated with an index currently.

³Query results of *None* type may be implemented with another result type in the future, once experience with the query framework is acquired and a meaningful grammar for those results are defined.

Table 1 – Query combinations and associated result reference type.

Query	Result Type
H5Q_TYPE_DATA_ELEM	Dataset Region
H5Q_TYPE_ATTR_VALUE	Attribute
H5Q_TYPE_ATTR_NAME	Attribute
H5Q_TYPE_LINK_NAME	Object
Dataset Element \wedge Dataset Element	Dataset Region
Dataset Element \wedge Attribute	Dataset Region
Dataset Element \wedge Object	Dataset Region
Attribute \wedge Attribute	Attribute
Attribute \wedge Object	Object
Object \wedge Object	Object
Dataset Element \vee Dataset Element	Dataset Region
Dataset Element \vee Attribute	Combination
Dataset Element \vee Object	Combination
Dataset Element \vee Combination	Combination
Attribute \vee Attribute	Attribute
Attribute \vee Object	Combination
Attribute \vee Combination	Combination
Object \vee Object	Object
Object \vee Combination	Combination
Combination \vee Combination	Combination
Combination \wedge Dataset Element	None
Combination \wedge Attribute	None
Combination \wedge Object	None
Combination \wedge Combination	None

The indexing API can work in conjunction with the view creation. When an `H5Qapply` call is made on a given location, an index attached to any dataset queried for element value ranges will be used to speed up the query process and return a dataspace selection to the library for later use. There are different techniques for creating data element indices, and the most efficient method will vary depending on the type of the data that is to be indexed, its layout, etc. A new interface for the HDF5 library that uses a plugin mechanism is therefore defined.

3.1 Indexing Interface and Plugins

This interface is defined for adding third-party indexing plugins, such as FastBit [2], ALACRITY [1], etc. The interface provides indexing plugins with efficient access to the contents of the container for both the creation and the maintenance of indices. In addition, the interface allows third-party plugins to create private

data structures within the container for storing the contents of the index. The current API as well as the plugin interface are presented below:

```

1  herr_t  H5Xregister(const H5X_class_t *idx_class);
2  herr_t  H5Xunregister(unsigned plugin_id);
3  herr_t  H5Xcreate(hid_t scope_id, unsigned plugin_id, hid_t xcpl_id);
4  herr_t  H5Xremove(hid_t scope_id, unsigned n /* Index n to be removed */);
5  herr_t  H5Xget_count(hid_t scope_id, hsize_t *idx_count);
6  herr_t  H5Xget_info(hid_t scope_id, unsigned n, H5X_info_t *info);
7  hsize_t H5Xget_size(hid_t scope_id);
8
9  typedef struct {
10     unsigned version;          /* Version number of the index plugin class struct */
11                                /* (Should always be set to H5X_CLASS_VERSION, which
12                                 * may vary between releases of HDF5 library) */
13     unsigned id;               /* Index ID (assigned by The HDF Group, for now) */
14     const char *idx_name;      /* Index name (for debugging only, currently) */
15     H5X_type_t type;           /* Type of data indexed by this plugin */
16
17     /* Callbacks, described above */
18     void *(*create)(hid_t dataset_id, hid_t xcpl_id, hid_t xapl_id,
19                     size_t *metadata_size, void **metadata);
20     herr_t (*remove)(hid_t dataset_id, size_t metadata_size, void *metadata);
21     void *(*open)(hid_t dataset_id, hid_t xapl_id, size_t metadata_size,
22                   void *metadata);
23     herr_t (*close)(void *idx_handle);
24     herr_t (*copy)(hid_t src_dataset_id, hid_t dest_dataset_id, hid_t xcpl_id,
25                    hid_t xapl_id, size_t src_metadata_size, void *src_metadata,
26                    size_t *dest_metadata_size, void **dest_metadata);
27     herr_t (*pre_update)(void *idx_handle, hid_t dataspace_id, hid_t xxpl_id);
28     herr_t (*post_update)(void *idx_handle, const void *buf, hid_t dataspace_id,
29                            hid_t xxpl_id);
30     herr_t (*query)(void *idx_handle, hid_t query_id, hid_t xxpl_id,
31                     hid_t *dataspace_id);
32     herr_t (*refresh)(void *idx_handle, size_t *metadata_size, void **metadata);
33     herr_t (*get_size)(void *idx_handle, hsize_t *idx_size);
34 } H5X_class_t;

```

Index objects are stored in the HDF5 container that they apply to, but are not visible in the container's group hierarchy⁴. Instead, index objects are part of the metadata for the indexed dataset. New index objects are created by passing an HDF5 location (group or dataset) to be indexed and the index plugin ID to the `H5Xcreate` call. Alternatively an index may be created at the same time as a dataset gets created by passing a property to the dataset creation property list. Index information (such as plugin ID and index metadata) is stored at index creation time⁵, and when the user later calls `H5Dopen`, the plugin open callback will retrieve this stored information and make use of the corresponding index plugin for all subsequent operations⁶. Similarly, calling `H5Dclose` will call the plugin index close callback and close the objects used to store the index data.

⁴Plugin developers, note that the HDF5 library's existing anonymous dataset and group creation calls can be used to create objects in HDF5 files that are not visible in the container's group hierarchy.

⁵Adding index information introduces a file format change.

⁶Note that the dataset's index object is opened only when a query call is made on the dataset.

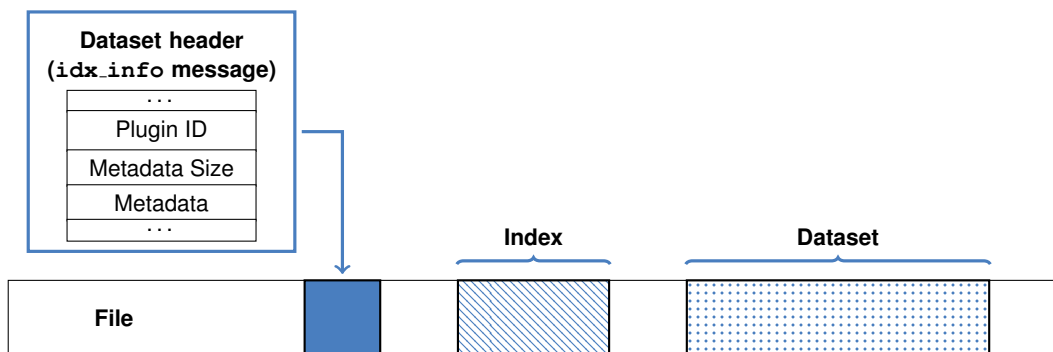


Figure 4 – Index information (plugin ID and metadata) is stored along the object header.

When a call to `H5Dwrite` is made, the index plugin `pre_update` and `post_update` callbacks will be triggered, allowing efficient index update by first telling the index plugin the region that is going to be updated with new data, and then realizing the actual index update, after the dataset write has completed. This allows various optimization to be made, depending on the data selection passed and the index plugin used. For example, a plugin could store the region and defer the actual index update until the dataset is closed, hence saving repeated index computation/update calls.

When a call to `H5Qapply` is made, the index plugin query callback will be invoked to create a selection of elements in the dataset that match the query parameters. Applications can also directly use the internally called `H5Dquery` routine defined below to directly execute a query on a particular dataset (accelerated by any index defined on the dataset), and retrieve the selection that matches the query.

```
1 | herr_t H5Dquery(hid_t dset_id, hid_t space_id, hid_t query_id, hid_t xapl_id,
2 |               hid_t *space_id);
```

Because the amount of space taken by the index cannot be directly retrieved by the user (since the datasets storing the indices are known only by the plugin itself), the `get_size` callback can query the amount of space that the index takes in the file and users may correspondingly query that information using the `H5Xget_size` routine.

3.1.1 Current and future plugins

Implementations for FastBit and ALACRITY index packages are already present, as well as a dummy brute force indexing plugin. Early performance results are presented in Figure 5.

In the future more plugins will be added, with or without external dependency (e.g., PyTables indexing, bitmap indexing). To satisfy that need, dynamic plugin loading and registration will be supported, allowing external libraries to plug to the current interface.

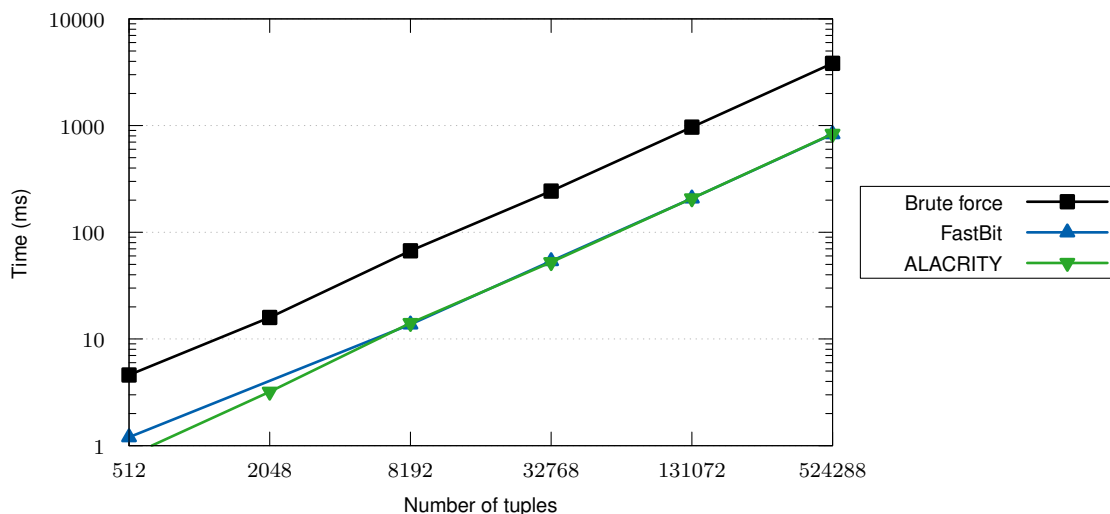


Figure 5 – Indexing performance.

3.2 Limitations

There are some existing limitations in the use of indices in the current implementation: FastBit and ALACRITY do not support incremental updates, an index is a shared resource for a dataset. Taken together, these conspire to put limits on application updates to datasets with indices. Additionally, because FastBit and ALACRITY do not allow incremental updates to an index, each modification to an existing index forces the index to be entirely rebuilt. The limitation in FastBit and ALACRITY will need to be addressed in the base packages' implementation, so that incremental updates to their index information can be made.

Questions

Some questions are still open regarding the handling of indices:

1. How to handle index updates when the specified index plugin is not available? *In traditional databases, stored procedures are saved with the data and therefore available at any time, but that is not the case here. We could mark the index as out of date and rebuild the index when the plugin is available again.*
2. How to handle index queries when the specified index plugin is not available? *We could fallback to another plugin and do a brute force query on the data.*

3.3 Support for HDF5 Compound Types

In a simple scenario, the HDF5 datatype used for creating the dataset can be defined as a native and simple type. Therefore, building an index on this dataset implies building that index from the entire dataset. However, in more complex scenarios, the dataset may have been created by using a compound datatype, hence defining multiple fields composed of native and simple datatypes within that same dataset. Consequently, creating an index from that dataset requires the user to select a particular field to be indexed, which may lead to having multiple indices per dataset depending on the number of fields that it contains. This can be done by

passing the `datatype_id` of the field to be indexed to the `xcpl_id`, the index creation property list, of the `H5Xcreate` call, which passes it down to the `plugin create` callback. As multiple fields can be defined, the field `datatype_id` must be stored along with the existing metadata, within the `idx_info` message (see Figure 4) so that the index associated to the field can be retrieved at the time of the query. When doing a query, the compound type is passed to the `H5Qcreate` call. The corresponding index is then used and the query is passed to the `query` callback of the plugin.

Consequently, when removing an index, one may choose to remove the index that corresponds to a particular field. This can be achieved by calling `H5Xget_info`, compare the `datatype` returned within the `info` structure, and pass the index number that needs to be removed.

3.4 Support for HDF5 Chunking

To support indexing of HDF5 chunks, we make each chunk a local *sub-dataset* of the original dataset. In that sense, handling every chunk can be seen as handling a dataset from the indexing plugin point of view. If the dataset is chunked, at the time of the index creation, we create a B-tree⁷ (physically stored on disk) that maps the coordinates of the chunks to the index plugin metadata. When the `create` callback is called (by representing the chunk as a local dataset, i.e., making the dataset layout point to the address of the chunk), metadata information is returned and stored. In the case of contiguous datasets, the index metadata as well as the index plugin ID is stored within the dataset header of the index info message (see Figure 4). In the case of chunked datasets, multiple metadata that correspond to each index created from each chunk may be accessed. Therefore, only the address of the B-tree that contains the metadata pieces for each chunk is stored in that header message and the index metadata itself is stored in the B-tree.

When the dataset is opened and the index reopened, we can lookup the index information in the B-tree that corresponds to each chunk and call the `open` callback using the associated metadata. Similarly, when a query is issued and needs to be answered, the chunks that correspond to the selection passed to the `H5Dquery` call are selected and their index is used to answer that query. The selection returned is then added to a global selection, which is then in turn returned to the user. Finally, when `H5Xremove` is called, the `delete` plugin callback is invoked for each chunk, by using the index metadata information stored in the B-tree.

3.5 Support for Parallel Indexing

An important design choice to support parallel indexing is to give as much freedom as possible to the indexing plugin developer so that in the case when the indexing library supports parallel indexing, it is still possible to take advantage of it. Three options are available to support parallel indexing from the previous interface:

- Let the index creation be collective. This however implies having synchronization points, which is the main constraint.
- Let the index creation be independent. However, creating datasets to store the index data must be done collectively⁸.
- Make the index creation in two phases. One that consists of building the index in parallel, independently, and gathering the information (index size, etc) at the end. The other that consists of letting the dataset

⁷The B-tree could also be replaced by a map object.

⁸An option could be to use the metadata server VOL plugin but this option is not easily doable yet.

creation be done by a master process, which can then let the other processes write the index data independently. There could however be a memory constraint in this case if the index data has to be kept in memory between these two phases (though building the index twice is not a good solution either).

Changes in the plugin interface include passing the parallel context to the callback (MPI communicator), which can be done by using the index access property list. In the case of chunked datasets, one may also want to operate on several chunks at the same time, in parallel. This can be done by passing a list of IDs that corresponds to each chunk to the callbacks. However, this also makes the plugin interface more heavy.

4 HDF5 and tools

Existing HDF5 tools must be compatible and take into account the existence of indices in the file if there are any. For reference, the following behavior for the tools is given:

- h5copy: copy
- h5dump: report index information
- h5ls: report index information
- h5diff: ignore index information?
- h5repack: copy index information / or generate index
- h5edit: ignore index information?
- h5toh4: ignore index information?
- h5import: ignore index information

Additional tools for indexing data and answering queries will also be added in the future.

5 Usage Example

In the following example, we show how one can make use of the query and indexing interface to retrieve a dataspace selection within a dataset. For simplicity's sake, we first create a dataset within the file, then open it to create and attach a new index, in order to finally query data from it. Note that for convenience, calls to directly read data that corresponds to the result of the index query may be moved to the high-level API in the future.

```
1 | #define NTUPLES 256
2 |
3 | int
4 | main(int argc, char *argv[])
5 | {
6 |     float data[NTUPLES];
7 |     hsize_t dims[1] = {NTUPLES};
8 |     hid_t t file_id, dataspace_id, dataset_id;
9 |     hid_t query_id, result_space_id;
10 |     size_t result_npoints;
11 |     float *result;
12 |     int i;
13 | }
```

```

14  /* Initialize data. */
15  for(i = 0; i < NTUPLES; i++) data[i] = (float) i;
16
17  /* Create file. */
18  file_id = H5Fcreate(file_name, H5F_ACC_TRUNC, H5P_DEFAULT, H5P_DEFAULT);
19
20  /* Create the data space for the dataset. */
21  dataspace_id = H5Screate_simple(rank, dims, NULL);
22
23  /* Create dataset. */
24  dataset_id = H5Dcreate(file_id, "Pressure", H5T_NATIVE_FLOAT, dataspace_id,
25      H5P_DEFAULT, H5P_DEFAULT, H5P_DEFAULT);
26
27  /* Write dataset. */
28  H5Dwrite(dataset_id, H5T_NATIVE_FLOAT, H5S_ALL, H5S_ALL, H5P_DEFAULT, data);
29
30  /* Close the dataset. */
31  H5Dclose(dataset_id);
32
33  /* Close dataspace. */
34  H5Sclose(dataspace_id);
35
36  /* Open dataset. */
37  dataset_id = H5Dopen(file_id, "Pressure", H5P_DEFAULT);
38
39  /* Create index using FastBit. */
40  H5Xcreate(file_id, H5X_PLUGIN_FASTBIT, dataset_id, H5P_DEFAULT);
41
42  /* Close the dataset. */
43  H5Dclose(dataset_id);
44
45  /* Create a simple query */
46  query_id = H5Qcreate(H5Q_TYPE_DATA_ELEM, H5Q_MATCH_EQUAL, H5T_NATIVE_FLOAT,
47      &query_value);
48
49  /* Open dataset. */
50  dataset_id = H5Dopen(file_id, "Pressure", H5P_DEFAULT);
51
52  /* Use query to select elements in the dataset. */
53  result_space_id = H5Dquery(dataset_id, query_id);
54
55  /* Allocate space to read data. */
56  result_npoints = (size_t) H5Sget_select_npoints(result_space_id);
57  result = malloc(result_npoints * sizeof(float));
58
59  /* Read data using result_space_id. */
60  H5Dread(dataset_id, H5T_NATIVE_FLOAT, H5S_ALL, result_space_id,
61      H5P_DEFAULT, result);
62
63  /* Use result. */
64
65  /* Free result. */
66  free(result);

```

```
67 |  
68 |     /* Close the dataset. */  
69 |     H5Dclose(dataset_id);  
70 |  
71 |     /* Close dataspace. */  
72 |     H5Sclose(result_space_id);  
73 |  
74 |     /* Close query. */  
75 |     H5Qclose(query_id);  
76 |  
77 |     /* Close the file. */  
78 |     H5Fclose(file_id);  
79 | }
```

Revision History

- Jul. 17, 2014:* Version 1 circulated for comment within The HDF Group.
- Aug. 20, 2014:* Version 2 circulated for comment within The HDF Group.
- Nov. 14, 2014:* Version 3 circulated for comment within The HDF Group.
- Feb. 25, 2015:* Version 4 edits and typo fixing.
- Jul. 24, 2015:* Version 5 edits to view object.
- Oct. 14, 2015:* Version 6 rework of query/view sections.

References

- [1] J. Jenkins, I. Arkatkar, S. Lakshminarasimhan, I. Boyuka, David A., E. Schendel, N. Shah, S. Ethier, C.-S. Chang, J. Chen, H. Kolla, S. Klasky, R. Ross, and N. Samatova. ALACRITY: Analytics-Driven Lossless Data Compression for Rapid In-Situ Indexing, Storing, and Querying. In A. Hameurlain, J. Kng, R. Wagner, S. Liddle, K.-D. Schewe, and X. Zhou, editors, *Transactions on Large-Scale Data- and Knowledge-Centered Systems X*, volume 8220 of *Lecture Notes in Computer Science*, pages 95–114. Springer Berlin Heidelberg, 2013.
- [2] K. Wu. FastBit: an efficient indexing technology for accelerating data-intensive science. *Journal of Physics: Conference Series*, 16(1):556, 2005.