



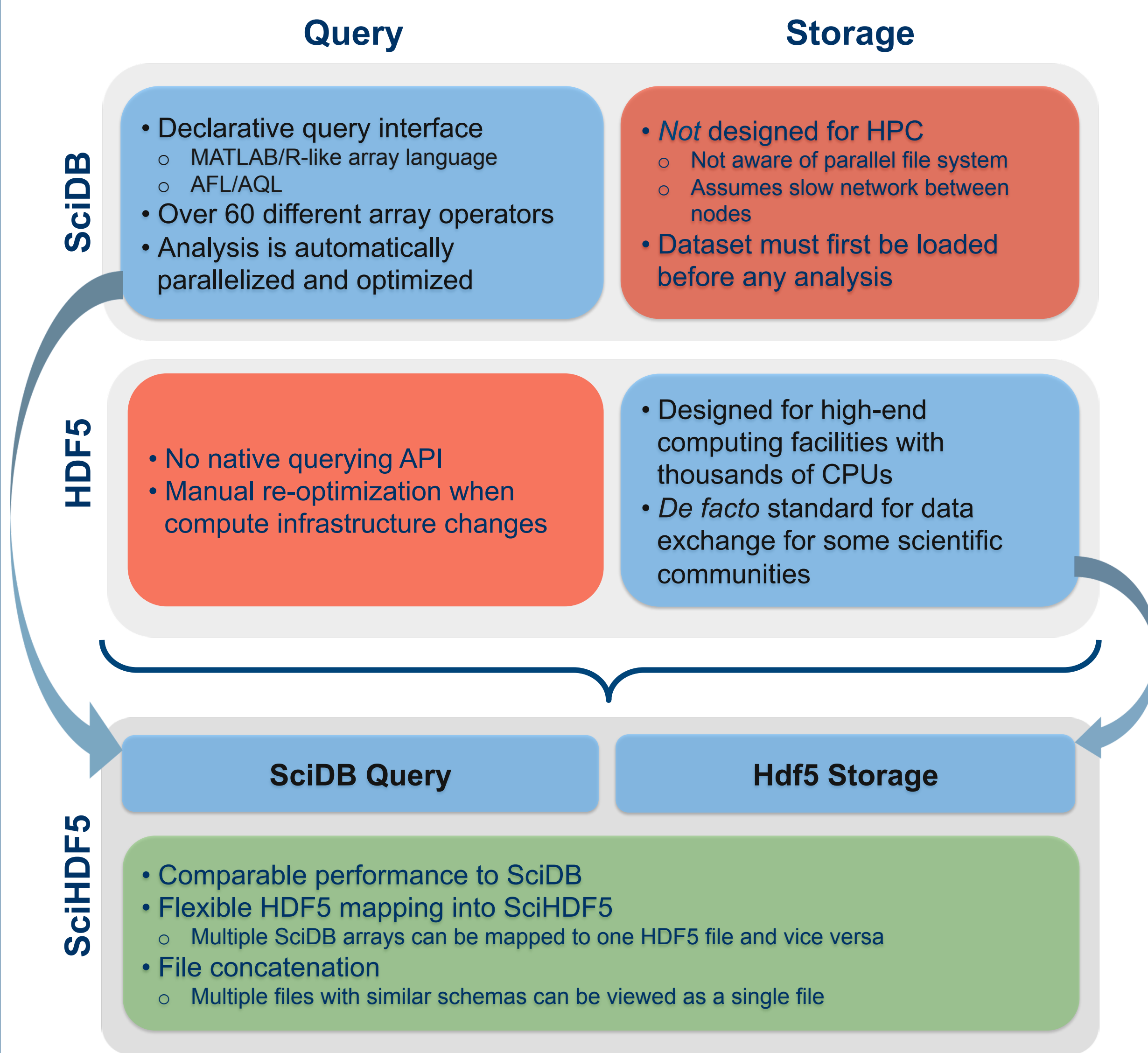
# SciHDF5: Integrating SciDB with the HDF5 File Format

S. Floratos<sup>1,2</sup>, S. Blanas<sup>2</sup>, S. Byna<sup>1</sup>, B. Dong<sup>1</sup>, K. Wu<sup>1</sup>, Prabhat<sup>1</sup>

<sup>1</sup>Lawrence Berkeley National Laboratory, <sup>2</sup>The Ohio State University

## Motivation

Large-scale experiments are often accompanied with massive datasets. In many cases, scientific experiments need to meet very fast access requirements in complex datasets. A well-known technology that is being used to meet these requirements is the hierarchical Data Format (HDF5). Another system for large scale analysis often adopted by the scientific community is SciDB, a parallel array database.



## How to Use SciHDF5

1

Declare SciDB array using the prefix “SciHDF5”

2

CREATE ARRAY now creates a SciDB object that contains “SciHDF5 attributes”

3

Every “SciHDF5 attribute” should be defined in the `scihdf5.config` file

4

Every “SciHDF5 attribute” corresponds to one HDF5 dataset

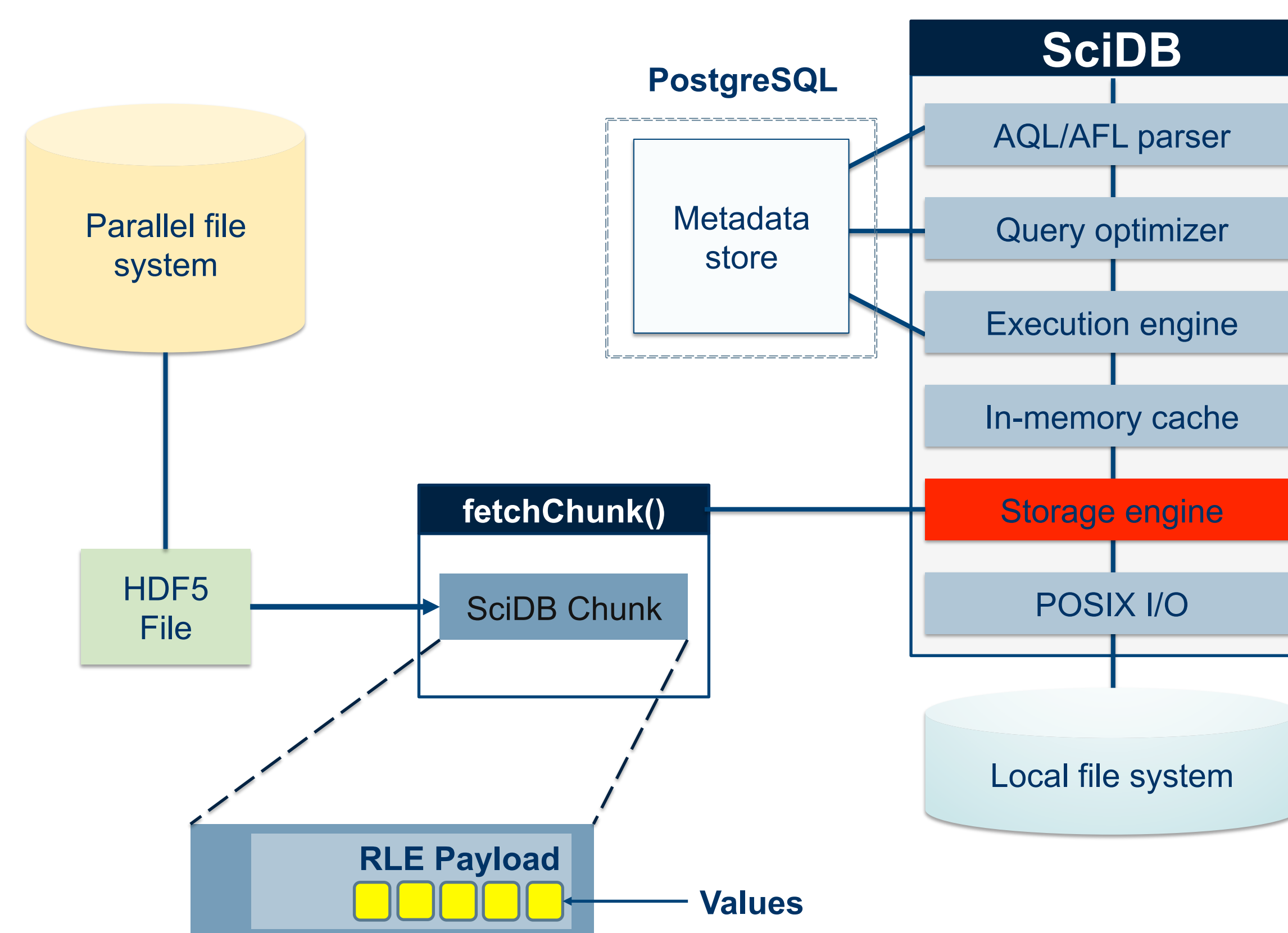
5

An “SciHDF5 attribute” can be mapped to multiple HDF5 files that contain the same HDF5 dataset

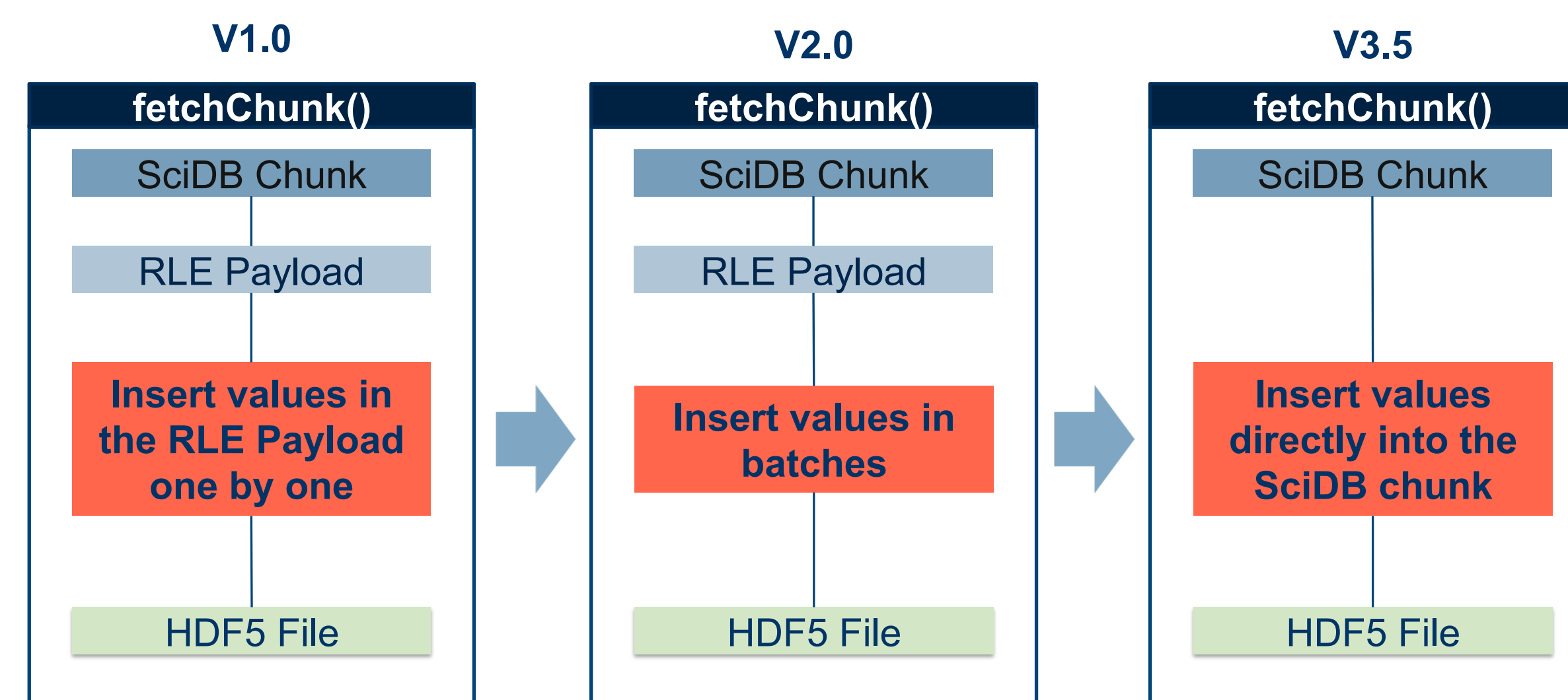
6

Queries touching this SciDB object read/write data directly from the HDF5 file

## SciHDF5 Architecture

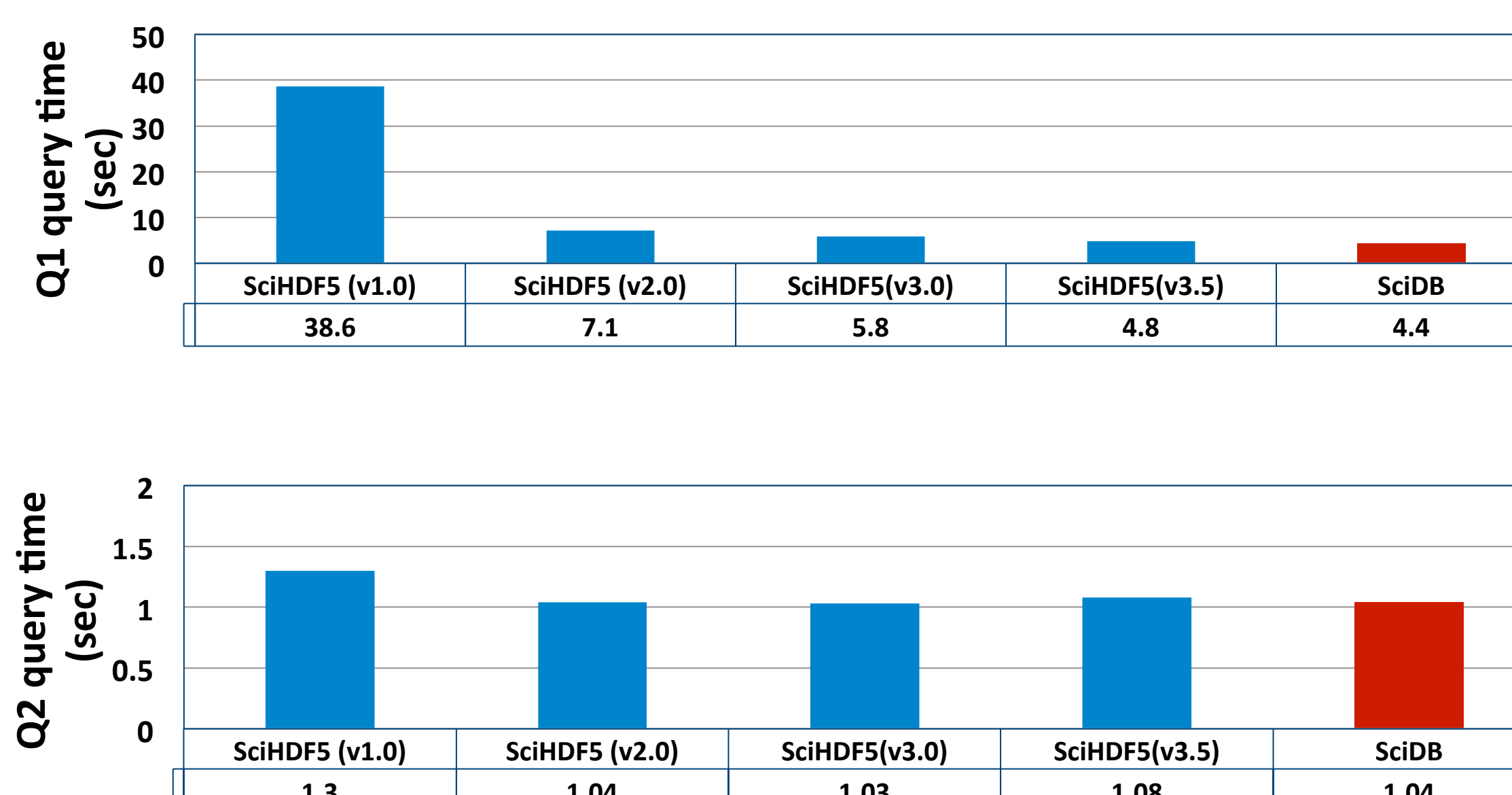


## SciHDF5 Optimizations



## SciHDF5 versus SciDB

We compare the above SciHDF5 optimizations to the original SciDB using a synthetic dataset of 10GB and 8 SciDB/SciHDF5 instances. The first query (Q1) is an aggregation over the entire dataset and the second one (Q2) is a random selection of 1 million elements.



## Additional Functionality

- **Users can define an HDF5 object in different ways**
  - One SciDB array with hdf5 attributes that are mapped to different HDF5 files.
  - Many SciDB arrays with hdf5 attributes that are mapped to the same HDF5 file.

### File concatenation

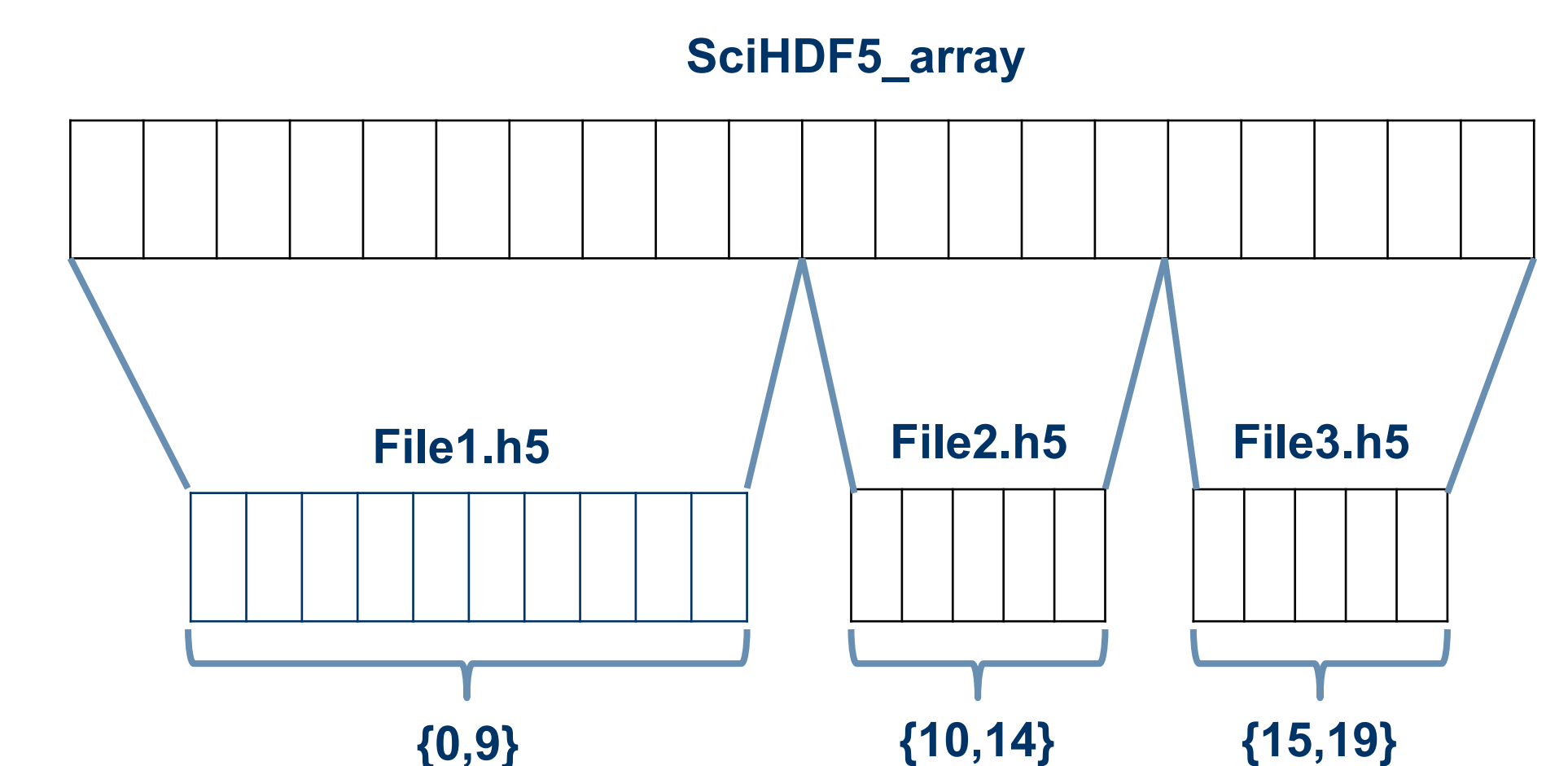
- If specified, SciHDF5 can assume that multiple HDF5 files are concatenated.
- First file in the “scihdf5.config” followed by the second file etc.

### Declaration of the SciDB array

```
“CREATE ARRAY SciHDF5_array <attribute_1, attribute_2> [i=0:20,10,0]”
```

### scihdf5.config

```
line 1: #HDF5 attributes
line 2: [attribute_1]
line 3: file = /path/to/file1.hdf5;/path/to/file2.hdf5;/path/to/file3.hdf5
line 4: dataset = /h5_group1/h5_dataset1
line 5:
line 6: [attribute_2]
line 7: file = /path/to/file1.hdf5;/path/to/file2.hdf5;/path/to/file3.hdf5
line 8: dataset = /h5_group2/h5_dataset2
```



## LUX/LZ Use Case Results

Our results are based on calibration data of the LUX/LZ detector that was built to analyze interactions from galactic dark matter. The original data sample contains 233 files with a total size of 1GB. For the purpose of this experiment, the dataset was duplicated 16 times. A cluster of 16 SciHDF5 instances were used to perform a query that analyzes the entire dataset such as calculating min, max and average values.

