

rhd5client – HDF5 server access

Vincent J. Carey, stvjc at channing.harvard.edu, Shweta Gopaulakrishnan, reshg at channing.harvard.edu, Samuel Pollack, spollack at jimmy.harvard.edu

February 27, 2018

Contents

1	HDF5 server	2
2	Some details	2
2.1	Motivation	2
2.2	Executive summary	2
2.3	Hierarchy of server resources	3

1 HDF5 server

HDF Server “extends the HDF5 data model to efficiently store large data objects (e.g. up to multi-TB data arrays) and access them over the web using a RESTful API.” In this package, several data structures are introduced

- to model the server data architecture and
- to perform targeted extraction of numerical data from HDF5 arrays stored on the server.

We maintain, thanks to a grant from the National Cancer Institute, the server <http://h5s.channingremotedata.org:5000/>. Visit this URL to get a flavor of the server structure: datasets, groups, and datatypes are high-level elements to be manipulated to work with data values from the server.

A key application of the rhdf5client package is support for the [restfulSE](#) package that defines an interface between the SummarizedExperiment class and the HDF5 Server. The server provides content for `assay()` requests to `RESTfulSummarizedExperiment` instances.

2 Some details

2.1 Motivation

Extensive human and computational effort is expended on downloading and managing large genomic data at site of analysis. Interoperable formats that are accessible via generic operations like those in RESTful APIs may help to improve cost-effectiveness of genome-scale analyses.

In this report we examine the use of HDF5 server as a back end for assay data.

A modest server configured to deliver HDF5 content via a RESTful API has been prepared and is used in this vignette.

2.2 Executive summary

We want to provide rapid read-only access to array-like data. To do this, the hierarchy and additional formalities of the HDF5 server data architecture are exposed through R functions and related classes. Full details on the HDF5 server are available at [the HDFgroup site](#).

The `dsmeta` function returns top-level groups and datasets.

```
library(rhdf5client)
bigec2 = H5S_source("http://h5s.channingremotedata.org:5000")
## analyzing groups for their links...
## done
bigec2
## HDF5 server domain: http://h5s.channingremotedata.org:5000
## There are 10 groups.
## Use groups(), links(), ..., to probe and access metadata.
## Use dsmeta() to get information on datasets within groups.
## Use [[ [dsname] ]] to get a reference suitable for [i, j] subsetting.
dsmeta(bigec2)[1:2,]      # two groups
```

r hdf5client – HDF5 server access

```
## DataFrame with 2 rows and 3 columns
##   groupnum
##   <integer>
## 1      1
## 2      2
##
## 1                               tenx_400k_sorted,mike,neurons400k, ...
## 2 newlinkedtype.datatypetest.test,compound.datatypetest.test,newdtype.datatypetest.test, ...
##           grp.uuid
##           <character>
## 1 ea87fd6e-ffba-11e7-985f-0aa4f45de806
## 2 ea8901c8-ffba-11e7-985f-0aa4f45de806
dsmeta(bigec2)[1,2][[1]] # all dataset candidates in group 1
## [1] "tenx_400k_sorted" "mike"          "neurons400k"
## [4] "tenx_full"        "tissues"       "assays"
## [7] "patelGBMSC"       "neurons100k"    "tenx_100k_sorted"
```

2.3 Hierarchy of server resources

2.3.1 Server

Given the URL of a server running HDF5 server, we create an instance of `H5S_source`:

```
mys = H5S_source(serverURL="http://h5s.channingremotedata.org:5000")
## analyzing groups for their links...
## done
mys
## HDF5 server domain: http://h5s.channingremotedata.org:5000
## There are 10 groups.
## Use groups(), links(), ..., to probe and access metadata.
## Use dsmeta() to get information on datasets within groups.
## Use [[ [dsname] ]] to get a reference suitable for [i, j] subsetting.
```

2.3.2 Groups

The server identifies a collection of ‘groups’. For the server we are working with, only one group, at the root, is of interest.

```
groups(mys)
## DataFrame with 10 rows and 2 columns
##   groups      nlinks
##   <character> <integer>
## 1 ea87fd6e-ffba-11e7-985f-0aa4f45de806      11
## 2 ea8901c8-ffba-11e7-985f-0aa4f45de806       4
## 3 ea88d144-ffba-11e7-985f-0aa4f45de806      28
## 4 ea88a142-ffba-11e7-985f-0aa4f45de806       3
## 5 ea885872-ffba-11e7-985f-0aa4f45de806      87
```

r hdf5client – HDF5 server access

```
## 6 ea88b8b2-ffba-11e7-985f-0aa4f45de806      1
## 7 ea887096-ffba-11e7-985f-0aa4f45de806      1
## 8 ea88ea1c-ffba-11e7-985f-0aa4f45de806      1
## 9 ea882fdc-ffba-11e7-985f-0aa4f45de806      1
## 10 ea888914-ffba-11e7-985f-0aa4f45de806     1
```

2.3.3 Links for a group

There is a class to hold the link set for any group:

```
lks = links(mys,1)
lks
## HDF5 server link set for group ea87fd6e-ffba-11e7-985f-0aa4f45de806
## There are 11 links.
## Use targets([linkset]) to extract target URLs.
```

2.3.4 Dataset access

We use double-bracket subscripting to grab a reference to a dataset from an H5S source. A dataset must be two-dimensional, i.e., accessible with two subscripts.

```
dta = bigec2[["tenx_100k_sorted"]]
dta
## H5S_dataset instance:
##           dsname intl.dim1 intl.dim2          created      type.base
## 1 tenx_100k_sorted    1e+05    27998 2017-05-11T15:30:23Z H5T_STD_I32LE
```

2.3.5 Data

Data are accessed by subscripting. The subscripts are arrays of increasing, positive integers.

```
x = dta[ 15:20, 1905:1906 ]
x
##      [,1] [,2]
## [1,]    40    1
## [2,]    35    0
## [3,]    13    1
## [4,]   118    2
## [5,]    25    1
## [6,]    26    1
```

(Obsolete) The subscripts are colon-delimited character with initial, final and optional stride.

```
x = dta["15:20", "1904:1906"]
x
##      [,1] [,2]
## [1,]    35    0
## [2,]    13    1
```

r hdf5client – HDF5 server access

```
## [3,] 118    2
## [4,] 25     1
## [5,] 26     1
```