





Contents

1	roduction	4			
2	Cor	nnecting to SciDB and Running Queries	4		
	2.1	Connecting to SciDB	4		
	2.2	Listing and removing SciDB arrays	5		
	2.3	Running SciDB queries	5		
	2.4	Iterating over query results	6		
3	SciDB Arrays for R Users				
4	The	e scidbdf data frame-like class	8		
	4.1	Examples of scidbdf objects	9		
5	Data manipulation functions				
	5.1	Subset operations	9		
	5.2	Database joins	10		
	5.3	Aggregation	11		
	5.4	Binding new values to an array	11		
	5.5	Sort, unique and index_lookup	12		
6	The scidb array class				
	6.1	Subsetting and indexing scidb array objects	13		
		6.1.1 Arithmetic operations	15		
		6.1.2 Available arithmetic operations	17		
	6.2	Persistence of dynamically-allocated scidb arrays	18		
	6.3	Miscellaneous array functions	18		
7	Using SciDB arrays with existing R code				
8	Package installation 2				

The scidb Package

8.1	Installing the R package from CRAN				
8.2	8.2 Installing the simple network service for SciDB				
	8.2.1 Installation on RHEL/CentOS 6	21			
	8.2.2 Installation on Ubuntu 12.04	21			
8.3	Error handling	22			
8.4	Package options, miscellaneous notes, and software license	22			

1 Introduction

SciDB is an open-source database that organizes data in *n*-dimensional arrays. SciDB features include ACID transactions, parallel processing, distributed storage, efficient sparse array storage, and native linear algebra operations. The **scidb** package for R provides two ways to interact with SciDB from R:

- 1. By running SciDB queries from R, optionally transferring data through data frames or data frame iterators.
- 2. Through several SciDB array object classes for R. The arrays mimic standard R arrays and data frames, but operations on them are performed by the SciDB engine. Data are materialized to R only when requested.

In some cases, R scripts and packages may be used with little or no modification with scidb arrays, allowing SciDB to power large-scale parallel R computation. This vignette illustrates using SciDB from R by example. For more detailed information on the functions described in this vignette, see the manual pages in the package.

2 Connecting to SciDB and Running Queries

This section outlines the most basic interaction between R and SciDB: running queries and transferring one-dimensional SciDB arrays between R and SciDB through R data frames.

2.1 Connecting to SciDB

The scidbconnect function establishes a connection to a simple HTTP network service called shim running on a SciDB coordinator instance (see Section 8.2). The function may be safely called multiple times. Once a connection is established, connection information is maintained until a different connection is established or the R session ends.

The network interface optionally supports TLS encryption and user authentication. Users are defined by the operating system users on the SciDB coordinator instance.

Listing 1: Connecting to SciDB

```
library("scidb")
# Connect to localhost by default on unencrypted port 8080:
scidbconnect()

# Connect to SciDB on an encrypted port 8083 with example authentication:
scidbconnect(host="localhost", port=8083, username="scidbuser", password="test")
```

The shim service can be configured to support either open/unencrypted or encrypted/authenticated ports, or both. We recommend using only encrypted/authenticated sessions when communicating with SciDB over public networks.

2.2 Listing and removing SciDB arrays

The scidblist function lists SciDB objects (arrays, instances, etc.), optionally showing detailed schema information for arrays. Returned results may be filtered using regular expression-style syntax.

The scidbremove function removes a SciDB array, or optionally a set of SciDB arrays defined by regular expression. The function accepts a vector of array names, resulting in the removal of all the specified arrays. Combine this feature with the regular expression filtering output of scidblist to remove sets of arrays matching the filter.

2.3 Running SciDB queries

The iquery function executes SciDB queries using either the SciDB array functional language (AFL) or declarative array query language (AQL) syntax. When AFL is used, the iquery function optionally returns query results in an R data frame if the argument return=TRUE is specified. Returned output is similar to output obtained by the SciDB iquery command-line program with the -olcsv+ option. The iquery function does not return anything by default.

Query results returned by the iquery function are internally presented to R using a generic CSV format, providing very flexible support for many data types. (The *n*-dimensional array class described in the next section uses a binary data exchange method between R and SciDB.) Note that, although R and SciDB have a number of common data types, each system contains types not supported by the other. Thus, conversion errors may arise. The iquery function is designed to reasonably minimize such issues and simplify basic data transfer between the systems. Data types common to R and SciDB include double-precision numeric, character string, logical, and 32-bit integers. They iquery function supports standard R read.table parameter options to facilitate type conversion.

Listing 2 illustrates basic use of iquery. (Program code is shown in black text and example results are shown in blue in all following example listings.)

Listing 2: Data frame example

```
library("scidb")
scidbconnect()
                             # Connect to SciDB on localhost
scidblist()
                             # List SciDB arrays (nothing there yet)
[1] NULL
# Build a 1-D SciDB array named "P:"
iquery("store(build(<x:double>[i=0:99,100,0],asin(1)*i/25),P)")
# Return to R the result of an apply operator:
S = iquery("apply(P,y,sin(x))",return=TRUE)
head(S)
1 0 0.0000000 0.0000000
2 1 0.0628319 0.0627905
3 2 0.1256640 0.1253330
4 3 0.1884960 0.1873810
5 4 0.2513270 0.2486900
```

2.4 Iterating over query results

The iquery function returns query results into a single R data frame by default. Large results expected to contain lots of rows may be iterated over by setting the iterative=TRUE argument. When iquery is used with the iterative=TRUE setting, it returns an iterator that iterates over chunks of rows of the result data frame. Iterators are defined by the iterators package. Their data may be directly accessed with the nextElem method, or indirectly with foreach. See the iterators and foreach packages for many examples and further documentation of their use.

Listing 3: Iterating over an iquery result

```
# Build a small 1-D SciDB test array:
iquery("store(build(< x:double > [i=1:10,10,0],i/10.0),A)")
# Return the result of a SciDB apply operator in an R iterator with a
# chunk size of at most 7 rows at a time:
it = iquery("apply(A,y,sqrt(x))", return=TRUE, iterative=TRUE, n=7)
nextElem(it)
  i
      X
1 1 0.1 0.316228
 2 0.2 0.447214
 3 0.3 0.547723
 4 0.4 0.632456
 5 0.5 0.707107
 6 0.6 0.774597
nextElem(it)
      x
  7 0.7 0.836660
  8 0.8 0.894427
    0.9 0.948683
 10 1.0 1.000000
nextElem(it)
Error: StopIteration
```

3 SciDB Arrays for R Users

Data are organized by SciDB in n-dimensional sparse arrays. "Sparse" in SciDB arrays means that array elements may be left undefined, and such array elements are omitted from computations. Note that this interpretation of sparse differs in a subtle way from that used by sparse matrices defined by R's Matrix package (whose sparse matrix elements are implicitly zero).

The elements of a SciDB array, called *cells*, contain one or more *attributes* (similar to R variables). The number and data types of attributes are uniform across all cells in an array. Thus, a one-dimensional SciDB array is conceptually similar to a data frame in R: the SciDB dimension index corresponds to data frame row index, and SciDB attributes to data frame columns. Higher-dimensional arrays in SciDB don't correspond directly to objects in R; the scidb *n*-dimensional array class described below is limited to working with one attribute at a time.

The integer coordinate systems used to index SciDB arrays are similar to R, except that SciDB integer indices may be zero or negative, and are represented by 62-bit signed integers (R indices are unsigned positive 31-bit integer or 52-bit integer-valued double values).

SciDB attribute values within a cell may be explicitly marked missing, indicated by a special SciDB missing code also referred to as a NULL code. SciDB internally supports a large number of possible missing codes. All SciDB missing code values are mapped to NA values in R.

In addition to available SciDB missing codes, SciDB double-precision floating point values also provide a value indicating missingness identically to R, and use the identical NA representation that R uses. Unlike R however, other SciDB data types do not define NA (use SciDB NULL instead).

The scidb package defines two array classes for R with data backed by SciDB arrays.

4 The scidbdf data frame-like class

The scidbdf class defines a data frame-like class with data backed by one-dimensional SciDB arrays. Like data frames, the columns represent variables of distinct types and the rows represent observations. Each attribute in the backing SciDB array represents a column in the scidbdf object. The scidbdf object elements are read-only (the backing SciDB array may be manually updated, for example using the iquery function). Non-integer row indices are not supported.

Use either the df2scidb or as.scidb functions to create new SciDB arrays and corresponding scidbdf R objects by copying R data frames into SciDB. The types and nullable options may be used to explicitly specify the SciDB type and nullability values of each data frame column. See the R help page for df2scidb for more information.

The scidb function returns an R scidbdf or scidb object representation of an existing SciDB array.

Objects of class scidbdf obey a subset of R indexing operations. Columns may be selected by numeric index or attribute name, but the short-hand R \$-style variable selection notation is not supported. Rows may only be selected by integer.

Subsets of scidbdf objects are returned as new scidbdf objects of the appropriate size (dimension, number of attributes/columns). The package uses the special empty-bracket notation, [], to indicate that data should be materialized to R as an R data frame. Illustrations are provided in the examples.

4.1 Examples of scidbdf objects

Listing 4: SciDB data frame-like objects

```
library("scidb")
scidbconnect()
# Copy the Michelson-Morley experiment data to SciDB, returning a scidbdf object
X = as.scidb(morley, name="morely")
str(X)
SciDB array name: morley
Attributes:
  attribute
            type nullable
       Expt int32
                      FALSE
        Run int32
                      FALSE
      Speed int32
                     FALSE
Row dimension:
  No name start length chunk_interval chunk_overlap low high type
                   100
# Materialize the first five rows of X to R (using [] to return results to R):
X[1:5, c("Run", "Speed")][]
   Run Speed
     1
         850
1
     2
         740
     3
         900
3
       1070
         930
```

5 Data manipulation functions

The package defines a number of common SciDB data manipulation functions for scidbdf objects.

5.1 Subset operations

The package supports standard array subsetting operations along dimension indices. Use the subset function to filter array contents by a boolean expression similarly to the usual R subset function. Under the hood, this function uses the SciDB filter operator—the function name subset more closely matches standard R syntax.

The subset function requires two arguments, a SciDB array reference and a valid SciDB logical expression represented as a string. Here is an example:

Listing 5: A 1-d SciDB array object using one of several available attributes

```
# Upload the iris data to SciDB and create a data frame-like object
# that refers to it:
data(iris)
df = as.scidb(iris,name="iris",gc=FALSE)
subset(df,"Petal_Width>2.4")[]
    Sepal_Length Sepal_Width Petal_Length Petal_Width
100
             6.3
                          3.3
                                        6.0
                                                     2.5 virginica
109
             7.2
                          3.6
                                        6.1
                                                     2.5 virginica
144
             6.7
                          3.3
                                        5.7
                                                     2.5 virginica
```

5.2 Database joins

The package overloads the R merge function, enabling a number of database join-like operations on SciDB array objects. Use the R help("merge",package="scidb") function for detailed help. Here is a simple example that performs an inner join on array attributes:

Listing 6: Merge example

```
authors = data.frame(
             surname = c("Tukey", "Venables", "Tierney", "Ripley", "McNeil"),
             nationality = c("US", "Australia", "US", "UK", "Australia"),
             deceased = c("yes", rep("no", 4)),
           stringsAsFactors=FALSE)
books = data.frame(
           name = c("Tukey", "Venables", "Tierney",
                     "Ripley", "Ripley", "McNeil", "R Core"),
           title = c("Exploratory Data Analysis",
                      "Modern Applied Statistics ...",
                      "LISP-STAT", "Spatial Statistics", "Stochastic Simulation",
                      "Interactive Data Analysis", "An Introduction to R"),
           other.author = c(NA, "Ripley", NA, NA, NA, NA, "Venables & Smith"),
           stringsAsFactors=FALSE)
a = as.scidb(authors)
b = as.scidb(books)
merge(a,b,by=list("surname", "name"))[,c(4:6,8:9)][]
   surname nationality deceased
                                                          title other_author
    McNeil
             Australia
                                     Interactive Data Analysis
                                                                         < NA >
1
   Ripley
                     UK
                              nο
                                             Spatial Statistics
                                                                         < NA >
2
   Ripley
                     UK
                                          Stochastic Simulation
                                                                         <NA>
                              no
   Tierney
                     US
                                                      LISP-STAT
                                                                         < NA >
                              no
                     US
                                     Exploratory Data Analysis
                                                                         <NA>
     Tukey
                             yes
5 Venables
                              no Modern Applied Statistics ...
             Australia
                                                                       Ripley
```

The merge implementation has some limitations outlined in its man page. Joins on attributes are presently the most limited case, and are restricted to inner joins at the moment. Use merge on array dimensions for the most flexibility.

5.3 Aggregation

SciDB aggregation functions are available generally for SciDB array objects with the aggregate function. Aggregation may be defined over array dimensions and/or array attributes very similarly to standard R aggregation syntax. Aggregation functions must be valid SciDB aggregate expressions, represented as strings.

Listing 7: Aggregation examples

```
Upload the iris data to SciDB and create a data frame-like object
# that refers to it:
data(iris)
df = as.scidb(iris,name="iris",gc=FALSE)
aggregate(df, by="Species", FUN="avg(Petal_Length) as avg_pl, stdev(Petal_Width)
   ")[]
  Species_index avg_pl Petal_Width_stdev
                                              Species
0
              0
                 1.462
                                0.1053856
                                               setosa
1
              1
                 4.260
                                0.1977527 versicolor
2
                 5.552
                                0.2746501
                                            virginica
```

Note that variable names may be specified in the SciDB aggregate expression.

5.4 Binding new values to an array

The SciDB package defines the bind function to add variables to arrays similar to the R cbind function for data frames. In the 1-d SciDB array scidbdf class the bind function is equivalent to the cbind function for data frames. However, bind can also operate on higher-dimensional arrays.

Listing 8: Binding new variables to an array

```
df = as.scidb(iris,name="iris",gc=FALSE)
   = bind(df, "prod", "Petal_Length * Petal_Width")
head(y, n=3)
  Sepal_Length Sepal_Width Petal_Length Petal_Width Species prod
           5.1
                        3.5
                                      1.4
                                                   0.2
                                                        setosa 0.28
           4.9
                        3.0
2
                                      1.4
                                                   0.2
                                                        setosa 0.28
3
           4.7
                        3.2
                                      1.3
                                                   0.2
                                                        setosa 0.26
```

5.5 Sort, unique and index_lookup

Use the **sort** function to sort on a subset of dimensions and/or attribute of a SciDB array object, creating a new sorted array.

Use the unique function to return a SciDB array that removes duplicate elements of a single-attribute SciDB input array.

Use the index_lookup function along with unique to bind a new variable that enumerates unique values of a variable similarly to the R factor function.

Examples follow. Note that we use the SciDB project function in one example. project presents an alternative syntax to the functionally equivalent column subset selection of variables using brackets. In some cases using project can be more efficient.

Listing 9: Sorting variables and enumerating unique values

```
x = as.scidb(iris,name="iris")
 Sort x by Petal_Width and Species
a = sort(x, attributes=c("Petal_Width", "Species"))
head(a, n=3)
  Sepal_Length Sepal_Width Petal_Length Petal_Width Species
           4.9
                        3.1
                                      1.5
                                                   0.1
2
           4.9
                        3.6
                                      1.4
                                                   0.1
                                                        setosa
3
           4.8
                        3.0
                                      1.4
                                                   0.1
# Find unique values of Species:
unique(a[, "Species"])[]
     Species
      setosa
1 versicolor
   virginica
# Add a new variable that enumerates the unique values of Species:
head(index_lookup(a, unique(project(a, "Species")), "Species"))
  Sepal_Length Sepal_Width Petal_Length Petal_Width Species Species_index
           4.9
                        3.1
                                      1.5
                                                   0.1
                                                                            0
                                                        setosa
2
                                                                            0
           4.9
                        3.6
                                      1.4
                                                   0.1
                                                        setosa
3
                                                                            0
           4.8
                        3.0
                                      1.4
                                                   0.1
                                                        setosa
```

HOMER

6 The scidb array class

Similarly to the data frame-like class, the package defines the scidb array class for R that represents vectors, matrices and general n-dimensional arrays. Array objects defined by the scidb class behave in some ways like standard R arrays. But their data reside in SciDB and most operations on them are computed by SciDB.

The scidb array class supports working with a single array attribute at a time to conform to R arrays (which generally support a single value per cell). Consider the iris data presented in Listing 10, represented within SciDB as a 1-D array with five attributes. Listing 10 illustrates creating a scidbdf data frame-like object in R referring to the iris data and also a 1-D scidb array object in R referring to the same data using the Sepal_Width attribute. Unlike R, SciDB numeric array indices may be zero or negative (indices begin at zero in the example in Listing 10). Data from scidb array objects are not materialized to R until extracted with the empty indexing function, [].

Listing 10: A 1-d SciDB array object using one of several available attributes

```
# Upload the iris data to SciDB and create a data frame-like object
# that refers to it:
data(iris)
df = as.scidb(iris,name="iris",gc=FALSE)
dim(df)
[1] 150
# Now create a vector-like object using the same data in SciDB, but
# referring to only one of the available variables (Sepal_Width):
x = scidb("iris", attribute="Sepal_Width", data.frame=FALSE)
length(x)
[1] 150
x[99:103][]
[1] 2.5 2.8 3.3 2.7 3.0
attributes(x)
                    # List all available attributes in the SciDB array
[1] "Sepal_Length" "Sepal_Width"
                                  "Petal_Length" "Petal_Width"
```

6.1 Subsetting and indexing scidb array objects

SciDB arrays act in many ways like regular arrays in R. Rectilinear subarrays may be defined by ranges of integer indices. Subarrays of scidb array objects are returned as new scidb array objects of the appropriate size.

Despite the similarities, there are differences between regular R and scidb array object indexing. In particular:

- The empty indexing function, [] applied to a scidb object materializes its array data as an R array. If the data exceed a return size threshold, an iterator over the array indices and data will be returned instead. The package option options("scidb.max.array.elements") controls the threshold.
- Index ranges follow SciDB convention. Arrays may have non-positive integer indices. In particular, note that the starting SciDB integer index is arbitrary, but often zero. (By contrast, the upper left corner of R arrays is always indexed by [1,1,...].)
- Array length may exceed 2^{31} elements.
- scidb array objects are limited to double-precision and 32-bit signed integer numeric, logical, and single-byte character (char) element data types.

Listing 11 illustrates basic subsetting operations on a sparse matrix. Note that the default SciDB index is zero-based.

Listing 11: Basic scidb subarray indexing

```
scidbremove("A", error=invisible)
# Create a small, sparse 2-d array in SciDB:
iquery("store(build_sparse(<val:double>[i=0:9,10,0,j=0:9,5,0],i*j,i<=j),A)")
A = scidb("A")
# Indexing operations return new SciDB arrays:
A[0:5,2:4]
A reference to a 6x3 SciDB array
# But their data can be materialized into an R array with []:
A[0:5,2:4][]
6 x 3 sparse Matrix of class "dgCMatrix"
[1,] 0 0
[2,] 2 3
[3,] 4 6 8
[4,] . 9 12
[5,] . . 16
[6,] . .
```

Listing 12 illustrates basic integer indexing operations on a sparse 3-D SciDB array. Note that since R has no direct way to represent sparse arrays with dimension greater than 2, the above result is returned as a list of values and their coordinates.

Listing 12: Basic scidb subarray indexing

```
scidbremove("A", error=invisible)
# Create a small, sparse 3-d array:
iquery("store(build_sparse(<val:double>[i=0:9,10,0,j=0:9,5,0,k=0:9,2,0],k,k<5
   and (j=1 \text{ or } j=3 \text{ or } j=5 \text{ or } j=7)),A)")
A = scidb("A")
dim(A)
[1] 10 10 10
# Indexing operations return new SciDB arrays:
A[0:3,2:3,3:4]
A reference to a 4x2x2 SciDB array
# But their data can be materialized into an R array with []:
A[0:3,2:3,3:4][]
$values
[1] 3 4 3 4 3 4 3 4
$coordinates
     [,1] [,2]
                [,3]
        0
             - 1
[1,]
[2,]
        0
              1
                    1
[3,]
        1
              1
                    0
[4,]
        1
              1
                    1
[5,]
        2
                    0
              1
[6,]
        2
              1
                    1
[7,]
        3
                    0
              1
[8,]
```

6.1.1 Arithmetic operations

The scidb array class supports a few standard linear algebra operations for dense and sparse matrices and vectors.

Listing 13 shows a more interesting 2-d dense array example that compares matrix arithmetic in R and SciDB. Similarly to data frames, the as.scidb function can be used to export R matrices and vectors to SciDB arrays.

Note that we generally use as.scidb for convenience—note that it's not the most efficient way to import data into SciDB. For very large data, use the SciDB bulk data loader as outlined in the SciDB documentation instead.

Listing 13: Matrix arithmetic in R and SciDB

```
set.seed(1)
v = as.scidb(matrix(rnorm(25),5))
str(v)
SciDB array name:
                          attribute in use:
                    V
All attributes:
Array dimensions:
  No name start length chunk_interval chunk_overlap low high
                      5
   0
               0
                                      5
                                                     0
                                                          0
                                                               4 int64
               0
                      5
                                      5
                                                     0
  1
        j
                                                          0
                                                               4 int64
crossprod( v[] )
                                # Compute t(v) %*% v using R
            [,1]
                       [,2]
                                  [,3]
                                              [,4]
      3.7779406
                  0.8044385
[1,]
                            -3.518737
                                         0.7422526
                                                   -3.463709
[2,]
      0.8044385
                  1.8806651
                             -3.127749
                                         1.0173059
[3,] -3.5187373
                -3.1277487
                              8.993765 -1.8112513
                                                    6.750959
[4,]
      0.7422526
                  1.0173059
                             -1.811251
                                         1.9202310
[5,] -3.4637086 -1.6524219
                             6.750959 -1.2491613
                                                    5.803521
# Now compute using SciDB, and materialize the result to R:
crossprod(v)[]
                                  [,3]
                                              [,4]
            [,1]
                        [,2]
                                                         [,5]
      3.7779406
[1,]
                  0.8044385 - 3.518737
                                         0.7422526
                                                   -3.463709
[2,]
      0.8044385
                  1.8806651 -3.127749
                                         1.0173059
[3,] -3.5187373
                 -3.1277487
                              8.993765
                                        -1.8112513
                                                    6.750959
                  1.0173059
[4,]
     0.7422526
                             -1.811251
                                         1.9202310
[5,] -3.4637086 -1.6524219
                              6.750959
                                        -1.2491613
                                                    5.803521
```

Linear algebra operations like the cross product in Listing 13 store their results in new dynamically-named SciDB arrays. The underlying SciDB array name for a scidb array object is available from the array object's @name slot.

Basic matrix/vector arithmetic operations on SciDB arrays (addition, subtraction, matrix and matrix vector products, scalar products, crossprod and tcrossprod) use standard R syntax. You can mix R and SciDB matrices and vectors and the scidb package will try to do the right thing by assigning R data to temporary SciDB arrays conforming to required database schema. Listing 14 shows an example of computations that mix scidb array objects with R vectors.

Listing 14: Mixed R and SciDB array arithmetic

```
# Build an array inside SciDB named 'U':
iquery("store(build(<x:double>[i=0:4,5,0,j=0:4,5,0],(i+1)/double(j+1)),U)")
u = scidb("U")
                       # u is an R representation of a SciDB array
set.seed (1)
x = cbind(rnorm(5))
                        # An R column vector
y = u %*% x
                        # Computed by SciDB, returning a SciDB array object
y[,drop=FALSE]
                        # Return the computed result to R
           [,1]
[1,] -0.3484533
[2,] -0.6969065
[3,] -1.0453598
[4,] -1.3938131
[5,] -1.7422663
```

Although the examples may seem trivial, the simple linear algebra capability shown in Listings 13 and 14 enable quite a lot of interesting computation. Later sections illustrate using this idea to overload more substantial functions in existing R packages.

6.1.2 Available arithmetic operations

The scidb class supports the operations shown in Table 1.

Expression	Operation	Operands	Output
A %*% B	Matrix multiplication	A, B Conformable SciDB arrays or R ma-	SciDB array
		trices/vectors	
$A \pm B$	Matrix summation/difference	A, B SciDB arrays or R matrices/vectors	SciDB array
$\mathtt{crossprod}(A,\!B)$	Cross product t(A) %*% B	A, B SciDB arrays or R matrices/vectors	SciDB array
${ t tcrossprod}(A,\!B)$	Cross product A %*% t(B)	A, B SciDB arrays or R matrices/vectors	SciDB array
A */B	Elementwise product/quotient	A, B Conformable SciDB arrays or R ma-	SciDB array
		trices/vectors	
$\alpha */A$	Scalar multiplication/division	SciDB array A , scalar α	SciDB array
t(A)	Transpose	SciDB array A	SciDB array
$A[\text{range}, \text{range}, \ldots]$	Subarray	SciDB array A	SciDB Array
A[]	Materialize	SciDB array A	R array
$\mathtt{svd}(A)$	Singular value decomposition	Dense SciDB array A	SciDB arrays

Table 1: SciDB Array Class Operations

The subarray/materialize operations [] support the standard drop argument.

6.2 Persistence of dynamically-allocated scidb arrays

Previous examples illustrate that new scidb arrays may be created after some R operations. For example, the subarray of a scidb array is a new scidb array. SciDB arrays created as the result of R operations do not persist by default—they are removed from SciDB when their corresponding R objects are deleted in R. Consider the example in Listing 15, shown with debugging turned on.

Listing 15: Non-persistence of intermediate arrays

```
# Build an array inside SciDB named 'U':
iquery("store(build(<x:double>[i=0:4,5,0,j=0:4,5,0],(i+1)/double(j+1)),U)")

u = scidb("U")
V = u[1:3,1:5]  # The subarray V is a new SciDB array
options(scidb.debug=TRUE)
rm(V)
gc()  # Force R to run garbage collection
remove(R_array817ceef8481100900000138)
```

The debugging message in Listing 15 illustrates that the temporary SciDB array that V represented was removed from SciDB when R garbage collection was run. In order to set any SciDB array associated with a scidb array object as persistent, set the array <code>@gc\$remove</code> setting to <code>FALSE</code>—for example, <code>V@gc\$remove=FALSE</code> in the above example.

Dynamically allocated arrays use the following naming convention: They all begin with "R_array" and end with a unique numeric identifier determined by the current SciDB session (1100900000138 in the above example).

6.3 Miscellaneous array functions

The count function applied to a scidb array object returns the count of non-empty cells in in the backing SciDB array.

crossprod and tcrossprod are defined for scidb array objects and mixtures of scidb and matrices.

The image function displays a heatmap of a regrid of a 2-D scidb array object, and returns the regridded array to R. The grid=c(m,n) function parameter specifies the regrid window sizes in each array dimension, and defaults to the array chunk sizes. The regrid aggregation function may be specified using the op function argument, and by default averages the array values over the regrid windows.

7 Using SciDB arrays with existing R code

This section illustrates using SciDB together with R and standard R packages from CRAN to compute solutions to large-scale problems. R functions that rely on linear algebra and aggregation operations may be adapted to use SciDB arrays in place of native R vectors and matrices in order to benefit from the large-scale parallel computing capabilities of SciDB.

The truncated singular value distribution (TSVD) is an important, widely used analysis method. Truncated SVD lies at the heart of principle components and other analysis methods.

We use the irlba package from CRAN, http://cran.r-project.org/web/packages/irlba/index.html to efficiently compute a truncated SVD. The IRLB algorithm used by the package relies on mostly matrix-vector products, and is well-suited for use with SciDB. In fact, we can use SciDB matrices with the irlba package without modifying the package at all.

The irlba package includes an option for user-defined matrix-vector products between a matrix A and a vector x, that is the R computation A %*% x, and for computation of t(A) %*% x. Because matrix vector and transpose operations are defined for the scidb array class, we don't technically need to use the user-defined option in the irlba package. However, by using the option, we can greatly improve efficiency by avoiding explicitly forming the matrix transpose by computing t(t(x) %*% A) instead of t(A) %*% x. Listing 16 illustrates this.

Listing 16: Efficient matrix vector product for IRLBA

```
# Let A be a scidb matrix
# Let x be a numeric vector
# Compute A %*% x if transpose=FALSE
# Compute t(A) %*% x if transpose=TRUE
# Return an R numeric vector.
matmul = function(A, x, transpose=FALSE)
{
   if(transpose)
   {
      return(t(crossprod(x,A))[,drop=FALSE]))
   }
   (A %*% x)[,drop=FALSE]
}
```

After defining the custom matrix or transpose matrix product in Listing 16, we can load and use the irlba package with SciDB arrays. Listing 17 illustrates computation of a few largest singular values and associated singular vectors of a $50,000 \times 50,000$ matrix with random entries (consuming about 18 GB). That problem large enough that it can't be computed easily in R-the matrix is too large to even represent in R version 2 (although that changes soon in R version 3).

Listing 17: Example large truncated SVD computation

After a while, the algorithm returns the truncated SVD in the variable S. The result obtained is comparable to what the svd(A, nu=3, nv=3) command would have produced, if it could handle the large matrix. Further optimizations are possible, but this simple example shows that it can be easy to get large-scale computation working without rewriting R code.

8 Package installation

Installation proceeds in two steps: installing the R package on any computer that has a network connection to a SciDB database, and installing a simple network service on the SciDB database coordinator computer.

8.1 Installing the R package from CRAN

The scidb package is available on CRAN. Start an R session and run:

Listing 18: Installing the R package from CRAN

```
install.packages("scidb")
```

8.2 Installing the simple network service for SciDB

The SciDB R package requires installation of a simple open-source HTTP network service called shim on the computer that SciDB is installed on. The service needs to be installed only on the

SciDB coordinator computer, not on client computers that connect to SciDB from R. It's available in packaged binary form for supported SciDB operating systems, and as source code which can be compiled and deployed on any SciDB installation.

Both installation approaches install the **shim** network service on the SciDB coordinator computer. Installing as a service requires root permission. The compiled source code version requires no special permissions to run.

Installation from binary software packages for SciDB-supported operating systems is easiest. Detailed up-to-date information can be found on Paradigm4's laboratory on Paradigm4's Github repository at https://github.com/Paradigm4/shim/wiki/Installing-shim. We outline installation for each supported operating system below. See our github page for source code. The open source package author, Bryan Lewis, maintains binary packages for SciDB-supported operating systems. They are tied to specific versions of SciDB. The present version is 13.9 (September, 2013).

8.2.1 Installation on RHEL/CentOS 6

Listing 19: Installing the simple HTTP service on RHEL

```
# Install with:
wget http://illposed.net/shim-13.2-1.x86_64.rpm
rpm -i shim-13.9-1.x86_64.rpm

# (Uninstall, if desired, with:)
yum remove shim
```

8.2.2 Installation on Ubuntu 12.04

Listing 20: Installing the simple HTTP service on Ubuntu

```
# Install with:
wget http://illposed.net/shim_13.2_amd64.deb
sudo gdebi shim_13.9_amd64.deb
# (Uninstall, if desired, with:)
apt-get remove shim
```

See the Wiki and web pages at https://github.com/Paradigm4/shim/ for up to date package information and source code.

The installed shim network service exposes SciDB as a very simple HTTP API. It includes a simple browser-based status and query tool. After installing shim, point your browser to the I.P. address of the SciDB coordinator machine and port 8080, for example: http://localhost:8080 on the

coordinator machine itself. Note that this API is not official and may change in the future. Help drive those changes by contributing ideas, code and bugfixes to the project on github, or feel free to discuss the service on the SciDB.org/forum.

8.3 Error handling

SciDB errors are trapped and converted to R errors that can be handled by standard R mechanisms. Some operations might try to return too much data to R, exceeding R's indexing limitations, system memory, or both. The package tries to avoid this kind of error using package options that limit returned data size shown in the next section.

8.4 Package options, miscellaneous notes, and software license

The scidb package defines several global package options. Package options may be set and retrieved with the R options function, and are listed in Table 2.

Option	Default value	Description
scidb.debug	NULL	Set to TRUE to display all queries issued to the SciDB engine and other
		debugging information.
scidb.index.sequence.limit	100 000 000	Maximum allowed scidb array object sequential indexing limit (for larger
		ranges, use between)
scidb.max.array.elements	100 000 000	Maximum allowed non-empty elements to return in a subsetting opera-
		tion of a scidb array object.

Table 2: Package options

Miscellaneous notes follow:

- R does not support 64-bit integer types. 64-bit signed and unsigned integers smaller than 2⁵³ in magnitude will be represented as double-precision floating point numbers. 64-bit integers outside that range appear as +/-Inf. All other integers (int8, uint8, int16, uint16, etc.) are represented in R by 32-bit signed integers. The uint32 type is not directly supported.
- R doesn't support single-precision floating point numbers. iquery results convert single-precision numbers within SciDB to double-precision floating-point numbers in R. Single-precision SciDB numbers are not supported by the scidb array class.
- SciDB does not natively support complex numbers. Loading complex numbers directly into SciDB from R is not defined
- The iquery function provides the most flexible mechanism for type conversion between the systems, fully under user control using read.table options.
- Allowed array naming conventions vary between R and SciDB. For example, SciDB does not allow decimal points in attribute names. The package may alter names with character substitution to reconcile names when it is reasonable to do so. A warning is emitted whenever an object is automatically renamed in this way.

Copyright (C) 2008-2013 SciDB, Inc.

The SciDB package for R is free software: you can redistribute it and/or modify it under the terms of the AFFERO GNU General Public License as published by the Free Software Foundation.

The SciDB package for R is distributed "AS-IS" AND WITHOUT ANY WARRANTY OF ANY KIND, INCLUDING ANY IMPLIED WARRANTY OF MERCHANTABILITY, NON-INFRINGEMENT, OR FITNESS FOR A PARTICULAR PURPOSE. See the AFFERO GNU General Public License for the complete license terms.

You should have received a copy of the AFFERO GNU General Public License along with the package. If not, see http://www.gnu.org/licenses/agpl-3.0.html