

AnnotationHub: A client package for retrieving data from the AnnotationHub web service

Marc Carlson

April 12, 2014

1 AnnotationHub Objects

The *AnnotationHub* package provides a client interface to resources stored at the AnnotationHub web service.

```
> library(AnnotationHub)
```

The *AnnotationHub* package is straightforward to use. The 1st thing you need to do to make use of it is to create an `AnnotationHub` object like this:

```
> ah = AnnotationHub()
```

Now at this point you have already done everything you need in order to get annotations. If you know exactly what the resource you want is called (and where it can be found), you could get it right now by just tab completing to it using the `$` operator.

Lets suppose that you knowd the following is the path to your data:

```
ah$goldenpath.hg19.encodeDCC.wgEncodeUwTfbs.wgEncodeUwTfbsMcf7CtcfStdPkRep1.narrowPeak_0.0.1.RData
```

Simply tab completing to the above path (followed by hitting enter), as demonstrated below will actually retrieve an object and then assign its contents to a local variable called `res`.

```
> res <- ah$goldenpath.hg19.encodeDCC.wgEncodeUwTfbs.wgEncodeUwTfbsMcf7CtcfStdPkRep1.
```

As you can see it's pretty easy to get data out using `AnnotationHub` objects. The rest of this vignette is mostly about helping you to make sure you are accessing the version of `AnnotationHub` objects that you intend to use and also about making sure that you can filter down the huge number of objects to the few that you are really interested in.

2 Configuring AnnotationHub objects

When you create the AnnotationHub object, it will set up the object for you with some default settings. If you look at the object you will see some helpful information about it.

```
> ah

class: AnnotationHub
length: 10778
filters: none
hubUrl: http://annotationhub.bioconductor.org/ah
snapshotVersion: 2.14/1.4.0; snapshotDate: 2014-04-01
hubCache: /home/biocbuild/.AnnotationHub
```

By default, you can see that the AnnotationHub object is set to the latest snapshotDate and a snapshot version that matches the version of Bioconductor that you are using. You can also learn about these data with the appropriate methods.

```
> snapshotVersion(ah)
```

```
[1] "2.14/1.4.0"
```

```
> snapshotDate(ah)
```

```
[1] "2014-04-01 PDT"
```

If you are interested in using an older version of a snapshot, you can list previous versions with the possibleDates like this:

```
> pd <- possibleDates(ah)
```

```
> pd
```

```
[1] "2013-03-20 GMT" "2013-03-21 GMT" "2013-03-22 GMT" "2013-03-27 GMT"
[5] "2013-04-05 GMT" "2013-04-30 GMT" "2013-06-24 GMT" "2013-06-25 GMT"
[9] "2013-06-26 GMT" "2013-06-27 GMT" "2013-06-28 GMT" "2013-06-29 GMT"
[13] "2013-10-30 GMT" "2013-11-21 GMT" "2013-12-20 GMT" "2013-12-27 GMT"
[17] "2014-04-01 GMT"
```

And then you can set the dates like this:

```
> snapshotDate(ah) <- pd[1]
```

3 Exploring and setting filters for `AnnotationHub`

If you are interested in how many annotation resources are currently available for your `AnnotationHub` object, you can just take the length like this:

```
> length(ah)
```

```
[1] 10778
```

Similarly, there are also methods to show the resource names, or even the full set of resource URLs for available resources.

```
> names <- head(names(ah), n=1)
```

```
> names
```

```
ensembl.release.69.fasta.ailuropoda_melanoleuca.cdna.Ailuropoda_melanoleuca.ailMel1.69.cdna.all.fa.rz
```

```
> urls <- head(snapshotUrls(ah), n=1)
```

```
> urls
```

```
http://annotationhub.bioconductor.org/ah/ensembl/release-69/fasta/ailuropoda_melanoleuca/cdna/Ailuropoda_melanoleuca.ailMel1.69.cdna.all.fa.rz
```

For humans, the number of resources available is going to be overwhelming. How should we cut this data set down to size? For this task, we introduce filters. Every `AnnotationHub` object contains a list of filters that can be configured to control which resources it can return. By default this list is empty, which means you get everything.

```
> filters(ah)
```

```
list()
```

How can we learn which things are available for filtering? For this we have defined `columns` and `keytypes` methods, which will list all the kinds of data that can be filtered on.

```
> columns(ah)
```

```

[1] "BiocVersion"      "DataProvider"      "Title"
[4] "SourceFile"       "Species"           "SourceUrl"
[7] "SourceVersion"    "TaxonomyId"       "Genome"
[10] "Description"      "Tags"             "RDataClass"
[13] "RDataPath"        "Coordinate_1_based" "Maintainer"
[16] "RDataVersion"     "RDataDateAdded"   "Recipe"

```

```
> keytypes(ah)
```

```

[1] "BiocVersion"      "DataProvider"      "Title"
[4] "SourceFile"       "Species"           "SourceUrl"
[7] "SourceVersion"    "TaxonomyId"       "Genome"
[10] "Description"      "Tags"             "RDataClass"
[13] "RDataPath"        "Coordinate_1_based" "Maintainer"
[16] "RDataVersion"     "RDataDateAdded"   "Recipe"

```

Once we know which things can be used to filter on, we can extract values that these things can be required to match. For this task, we have defined a `key` method.

```
> head(keys(ah, keytype="Species"))
```

```

[1] "9606"              "Acromyrmex echinator" "Acyrtosiphon pisum"
[4] "Aedes aegypti"     "Agaricus bisporus"   "Ailuropoda melanoleuca"

```

Now we are able to construct and assign a filter to our `AnnotationHub` object. Lets set it up to only find resources from humans.

```
> filters(ah) <- list(Species="Homo sapiens")
```

And now if we look we will see that our `AnnotationHub` object is only exposing resources from `Homo sapiens`.

```
> length(ah)
```

```
[1] 5339
```

```
> names <- head(names(ah),n=1)
```

```
> names
```

```

ensembl.release.69.fasta.homo_sapiens.cdna.Homo_sapiens.GRCh37.
69.cdna.all.fa.rz

```

```
> urls <- head(snapshotUrls(ah), n=1)
```

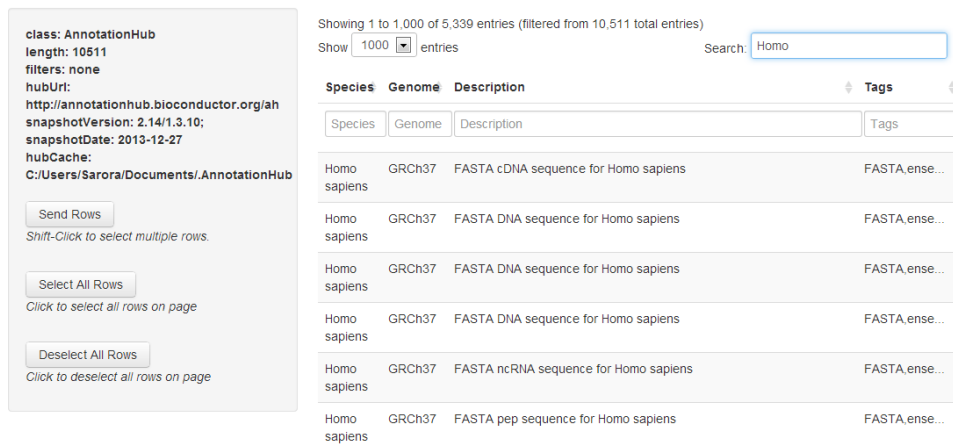
```
> urls
```

```
http://annotationhub.bioconductor.org/ah/ensembl/release-69/
fasta/homo_sapiens/cdna/Homo_sapiens.GRCh37.69.cdna.all.fa.rz
```

We can also look at the `AnnotationHub` object in a browser using the `display` function. We can then filter the `AnnotationHub` object for ‘Homo sapiens’ by either using the Global search field on the top right corner of the page or the in-column search field for ‘Species’.

```
> d <- display(ah)
```

Data Tables binding



The screenshot displays the AnnotationHub web interface. On the left, a sidebar shows the object's class (`AnnotationHub`), length (10611), filters (none), and hub URL (`http://annotationhub.bioconductor.org/ah`). It also includes snapshot version (2.14/1.3.10), snapshot date (2013-12-27), and a local cache path. Action buttons for 'Send Rows', 'Select All Rows', and 'Deselect All Rows' are present. The main content area shows a search bar with 'Homo' entered, displaying 1,000 of 5,339 filtered entries. A table with columns 'Species', 'Genome', 'Description', and 'Tags' is shown, listing FASTA sequences for Homo sapiens across different genome builds and transcript types.

Species	Genome	Description	Tags
Homo sapiens	GRCh37	FASTA cDNA sequence for Homo sapiens	FASTA,ense...
Homo sapiens	GRCh37	FASTA DNA sequence for Homo sapiens	FASTA,ense...
Homo sapiens	GRCh37	FASTA DNA sequence for Homo sapiens	FASTA,ense...
Homo sapiens	GRCh37	FASTA DNA sequence for Homo sapiens	FASTA,ense...
Homo sapiens	GRCh37	FASTA ncRNA sequence for Homo sapiens	FASTA,ense...
Homo sapiens	GRCh37	FASTA pep sequence for Homo sapiens	FASTA,ense...

Figure 1: Displaying and filtering the Annotation Hub object in a browser

By default 1000 entries are displayed per page, we can change this using the filter on the top of the page or navigate through different pages using the page scrolling feature at the bottom of the page.

We can also select the rows of interest to us and send them back to the R session using ‘Send Rows’ button ; this sets a filter internally which filters the `AnnotationHub` object.

4 Using `AnnotationHub` to retrieve data

So now that we have our `AnnotationHub` object configured to expose only the data for humans how would we go about getting that data downloaded?

As mentioned above, we can use the \$ operator and tab completion to pull down a data source of interest like this.

```
ah$goldenpath.hg19.encodeDCC.wgEncodeUwTfbs.wgEncodeUwTfbsMcf7CtcfStdPkRep1.narrowPeak_0.0.1.RData
```

Just by using tab completion like this:

```
> res <- ah$goldenpath.hg19.encodeDCC.wgEncodeUwTfbs.wgEncodeUwTfbsMcf7CtcfStdPkRep1.
```

And once you have done this, you can look at the object stored in res and use it etc.. Any dependencies that you need to use this kind of object should automatically try to load at this time.

```
> res
```

GRanges with 82163 ranges and 6 metadata columns:

	seqnames	ranges	strand	name	score
	<Rle>	<IRanges>	<Rle>	<character>	<integer>
[1]	chr1	[237640, 237790]	*	.	0
[2]	chr1	[544660, 544810]	*	.	0
[3]	chr1	[567480, 567630]	*	.	0
[4]	chr1	[569820, 569970]	*	.	0
[5]	chr1	[714200, 714350]	*	.	0
...
[82159]	chrX	[154764540, 154764690]	*	.	0
[82160]	chrX	[154807400, 154807550]	*	.	0
[82161]	chrX	[154881060, 154881210]	*	.	0
[82162]	chrX	[154892100, 154892250]	*	.	0
[82163]	chrX	[154916040, 154916190]	*	.	0
	signalValue	pValue	qValue	peak	
	<numeric>	<numeric>	<numeric>	<integer>	
[1]	30	26.89200	-1	-1	
[2]	6	8.16393	-1	-1	
[3]	100	56.71760	-1	-1	
[4]	85	49.65350	-1	-1	
[5]	17	13.18360	-1	-1	
...	
[82159]	26	25.2917	-1	-1	
[82160]	22	27.6521	-1	-1	
[82161]	17	16.4194	-1	-1	
[82162]	72	101.6090	-1	-1	
[82163]	32	32.5209	-1	-1	

seqlengths:

```
      chr1      chr10      chr11      chr12 ...      chr8      chr9      chrX
249250621 135534747 135006516 133851895 ... 146364022 141213431 155270560
```

Also, since you have previously downloaded this object at the start of this vignette, the 2nd time it should pull this object from a local cache that *AnnotationHub* will have created for you. This is a feature of *AnnotationHub* that is meant to provide better performance by removing the need to pull a large amount of data from a distant server every time. However, this does not mean that once you have used *AnnotationHub* to retrieve data that you no longer need to have internet access. This is because whenever you create a *AnnotationHub* object, it needs to talk to the metadata server to learn about things like the latest available version etc. So if you intend to access your objects on the plane you will need to either save them to a convenient location or else take note of where your local cache is located so that you can load them up manually later.

5 Session Information

R version 3.1.0 (2014-04-10)

Platform: x86_64-unknown-linux-gnu (64-bit)

locale:

```
[1] LC_CTYPE=en_US.UTF-8      LC_NUMERIC=C
[3] LC_TIME=en_US.UTF-8      LC_COLLATE=C
[5] LC_MONETARY=en_US.UTF-8  LC_MESSAGES=en_US.UTF-8
[7] LC_PAPER=en_US.UTF-8     LC_NAME=C
[9] LC_ADDRESS=C             LC_TELEPHONE=C
[11] LC_MEASUREMENT=en_US.UTF-8 LC_IDENTIFICATION=C
```

attached base packages:

```
[1] parallel  stats      graphics  grDevices  utils      datasets  methods
[8] base
```

other attached packages:

```
[1] GenomicRanges_1.16.0 GenomeInfoDb_1.0.0  AnnotationHub_1.4.0
[4] IRanges_1.21.45      BiocGenerics_0.10.0
```

loaded via a namespace (and not attached):

[1] AnnotationDbi_1.26.0	Biobase_2.24.0	BiocInstaller_1.14.0
[4] Category_2.30.0	DBI_0.2-7	GSEABase_1.26.0
[7] MASS_7.3-31	Matrix_1.1-3	RBGL_1.40.0
[10] RColorBrewer_1.0-5	RCurl_1.95-4.1	RJSONIO_1.0-3
[13] RSQLite_0.11.4	Rcpp_0.11.1	XML_3.98-1.1
[16] XVector_0.4.0	annotate_1.42.0	bitops_1.0-6
[19] caTools_1.16	colorspace_1.2-4	dichromat_2.0-0
[22] digest_0.6.4	genefilter_1.46.0	ggplot2_0.9.3.1
[25] graph_1.42.0	grid_3.1.0	gridSVG_1.4-0
[28] gtable_0.1.2	httpuv_1.3.0	httr_0.3
[31] interactiveDisplay_1.2.0	labeling_0.2	lattice_0.20-29
[34] munsell_0.4.2	plyr_1.8.1	proto_0.3-10
[37] reshape2_1.2.2	rjson_0.2.13	scales_0.2.3
[40] shiny_0.9.1	splines_3.1.0	stats4_3.1.0
[43] stringr_0.6.2	survival_2.37-7	tools_3.1.0
[46] xtable_1.7-3		