# archivist: Managing Data Analysis Results
https://github.com/pbiecek/archivist

Marcin Kosiński[1,2],
Przemysław Biecek[2]

[1]IT Research and Development
Grupa Wirtualna Polska

[2]Faculty of Mathematics, Informatics
and Mechanics, University of Warsaw

14 October, 2015

# About Grupa Wirtualna Polska

## The Leader Of The Polish Internet



Figure : wp.pl

## IT Research and Development

- large-scale online learning
- personalized news article recommendation
- e-mail targeting
- text mining
- web user behavior identification

# Main features

- allows to **store** and **archive** objects in repositories
  (stored on a local disk or via GitHub/Dropbox)
- provides handy tools facilitating objects' **search** and **recovery**
- ideally performs as **cache**
- supports the philosophy of **reproducible research**

## Solves reproducible research problems

- sometimes raw data are large or with limited access
- computations take a lot of time or require specialized hardware
- reproducibility requires specific versions of packages

# archivist: Cache Use Case

```r
getMaxDistribution <- function(D = rnorm, N = 10,
                               R = 1000000) {
  res <- replicate(R, max(D(N)))
  summary(res)
}

# or get directly from GitHub
library(archivist)
loadFromGithubRepo(md5hash="c", user="MarcinKosinski",
                   repo="Museum")
```

# archivist: Cache Use Case

Using the archivist one can prepare a repository which stores calls and results of the cache() function to avoid their re-call in the future.

```
cacheRepo <- tempdir()
createEmptyRepo(cacheRepo)

system.time(cache(cacheRepo, getMaxDistribution,
                                rnorm, 10) )
   user   system  elapsed
  5.230    0.090    5.315
system.time(cache(cacheRepo, getMaxDistribution,
                                rnorm, 10) )
   user   system  elapsed
  0.009    0.000    0.008
```

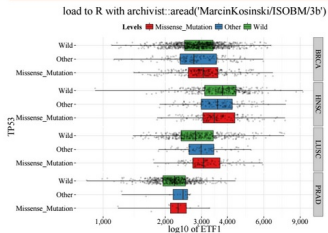# archivist: Retrieving an Object Use Case

## Fragment of a poster



Figure : ggplot object

## Using archivist to retrieve an object

```
aread('MarcinKosinski/ISOBM/3b')
    -> mutationsPlot
plot(mutationsPlot)
```
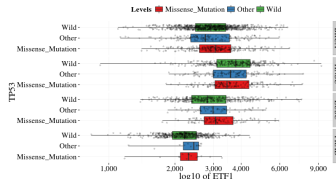


Figure : downloaded object

# archivist: Object's Pedigree Use Case

```
createEmptyRepo("FORUM_BI", default = TRUE)
invisible(aoptions("silent", TRUE))

data(iris)
iris %a%
  dplyr::filter(Sepal.Length < 16) %a%
  lm(Petal.Length~Species, data=.) %a%
  summary() -> obj

ahistory(obj)

     iris                                [ff575c261c949d073b2895b05d1097c3]
-> dplyr::filter(Sepal.Length < 16)      [9f7045b7322cdf3a9071377c6fe9c175]
-> lm(Petal.Length ~ Species, data = .)  [0a82efeb8250a47718cea9d7f64e5ae7]
-> summary()                             [671a0b89fccdf02087acb002374a0fcd]
```

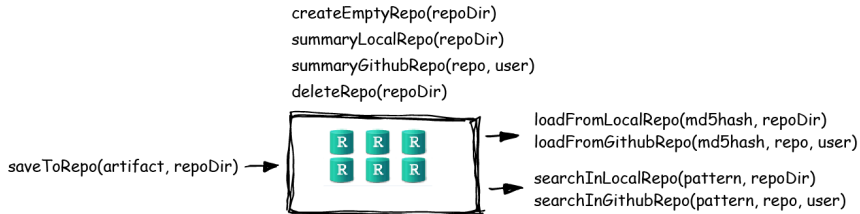# archivist: Objects' Exploration Within a Repository

```
models <- asearch("pbiecek/graphGallery",
 patterns = c("class:lm", "coefname:Sepal.Length"))
lapply(models, coef)
```

```
[[1]]
 (Intercept) Sepal.Length
   -7.101443     1.858433

[[2]]
     (Intercept)     Sepal.Length  Speciesversicolor  Speciesvirginica
      -1.7023422        0.6321099          2.2101378         3.0900021
```

# archivist: How does it work?

createEmptyRepo(repoDir)
summaryLocalRepo(repoDir)
summaryGithubRepo(repo, user)
deleteRepo(repoDir)

loadFromLocalRepo(md5hash, repoDir)
loadFromGithubRepo(md5hash, repo, user)

saveToRepo(artifact, repoDir) →

searchInLocalRepo(pattern, repoDir)
searchInGithubRepo(pattern, repo, user)

Each repository contains a database with objects metadata.
Objects are stored as binary files.
Each object has a unique key - md5 hash.
Metadata, like object class, name, creation date, relations with other objects
are useful when searching for an object in a repository.

```
library("archivist")
```

# archivist: Plans & Prototypes

Automated repository creation on github, commit, push and a return of a hook to an object.

```r
archive(iris, "MarcinKosinski",
 "archivist-Museum-RforeveR_last3",
        USER_EMAIL, USER_PASSWORD, app_key, app_secret,
           github_token) -> aread_input
# function returns a hook
archivist::aread("MarcinKosinski/archivist-Museum-
    RforeveR_last3/ff575c261c949d073b2895b05d1097c3")
archivist::aread(aread_input) -> x
digest::digest(x)

[1] "ff575c261c949d073b2895b05d1097c3"

identical(iris,x)

[1] TRUE
```

# archivist: Learn More

## http://pbiecek.github.io/archivist/