

# Updates to rhdf5

---

Mike Smith

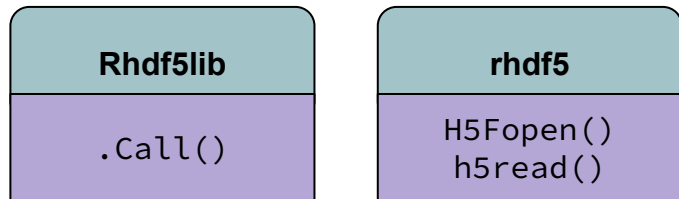
  @grimbough

# Introduction to **rhdf5**



One of several packages providing this functionality

# Introduction to **rhdf5**



C / C++ Library

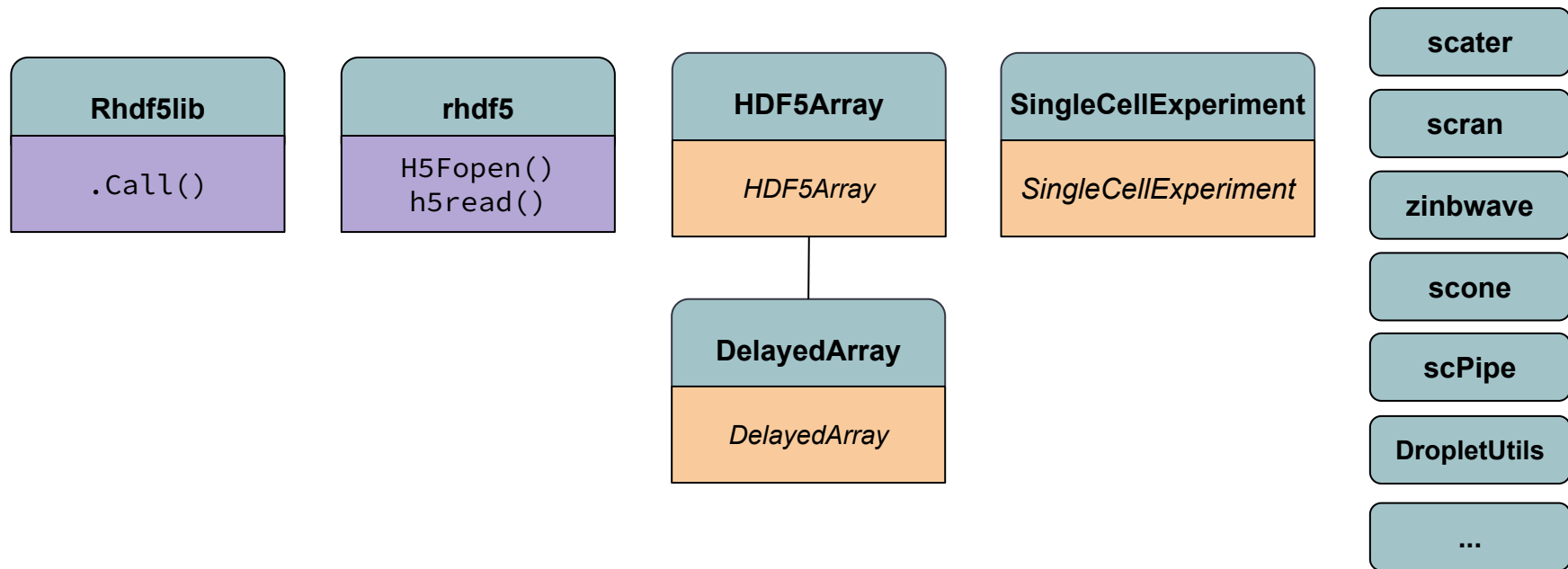
R Interface

Counts Matrix

Complete SC  
Dataset

Analysis Tools

# Introduction to rhdf5



C / C++ Library

R Interface

Counts Matrix

Complete SC  
Dataset

Analysis Tools

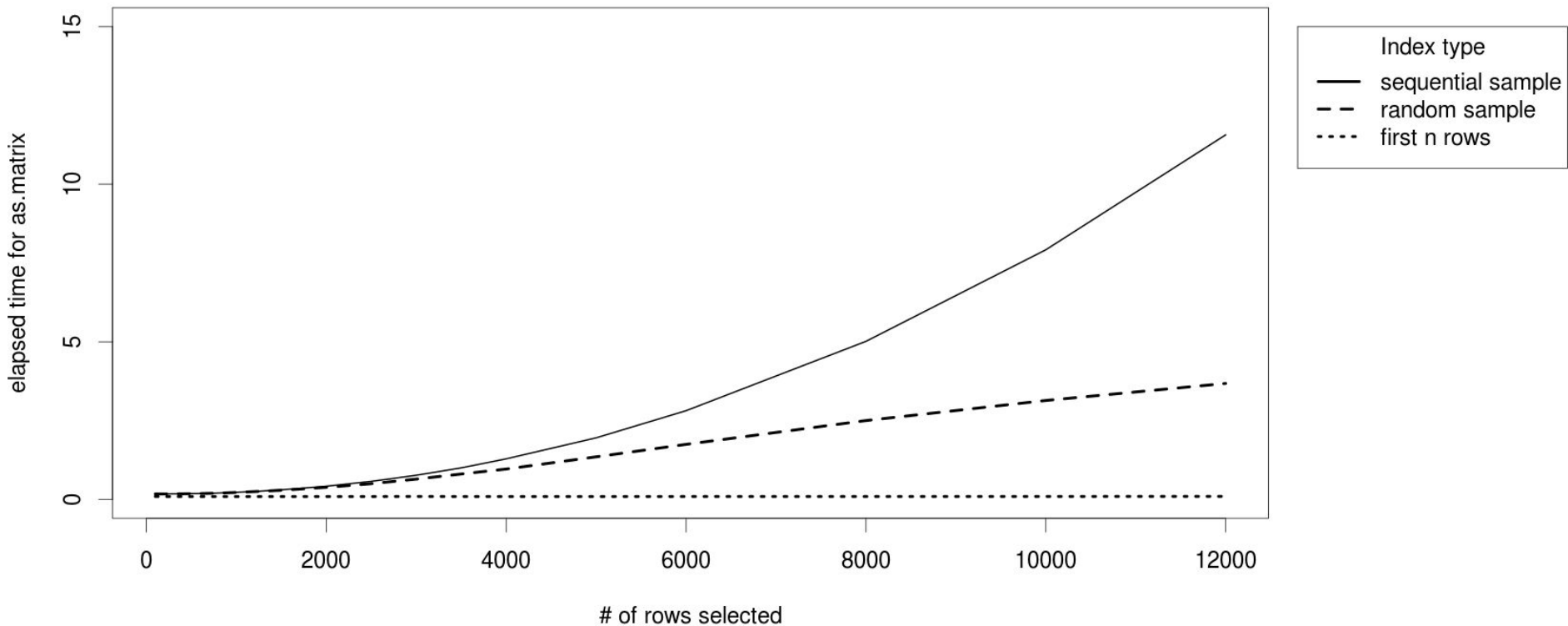
# Update to underlying HDF5 library

- Switch from HDF5 version 1.8 to 1.10
- Motivated by users unable to open files created by other software
  - <https://support.bioconductor.org/p/109845>
  - <https://github.com/pachterlab/sleuth/issues/175>
- Not a simple drop-in replacement!
  - *“The hid\_t type was changed from 32-bit to a 64-bit value.”*
  - Change should be transparent to users
- Potential interesting new features e.g. SWMR

# Performance improvements

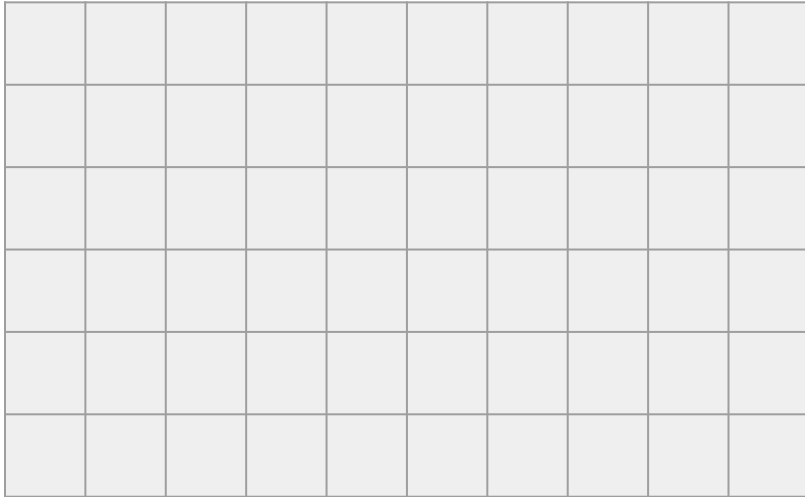
- Low hanging fruit - classic R inefficiencies
  - Unnecessary reordering
  - Copying rather than preallocating

# Performance improvements



# HDF5 Hyperlabs

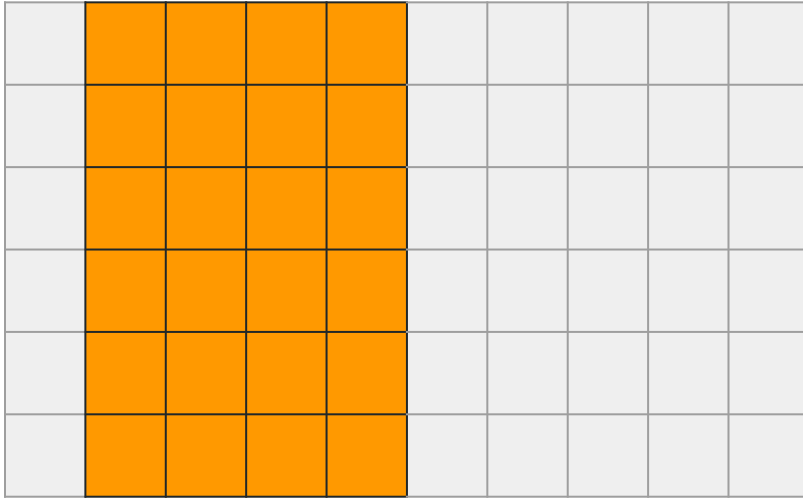
- Regularly spaced selections of elements





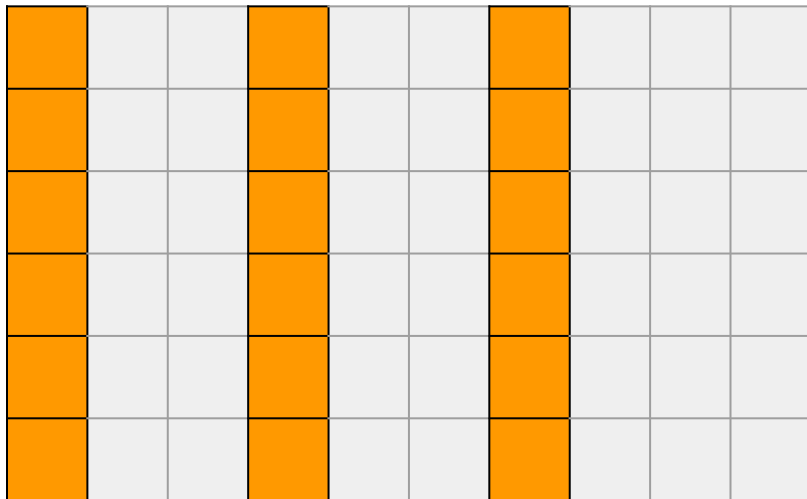
# HDF5 Hyperlabs

- Regularly spaced selections of elements



# HDF5 Hyperlabs

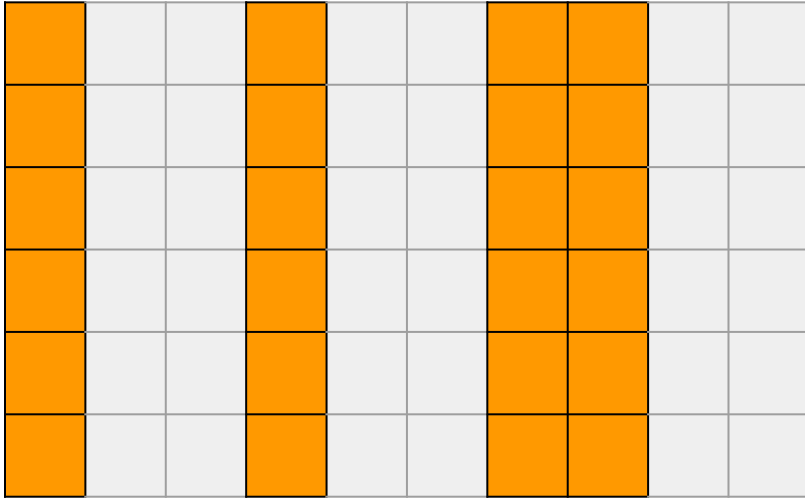
- Regularly spaced selections of elements



- Defined by *offset*, *count*, *stride* and *block*
- Available in **rhdf5**

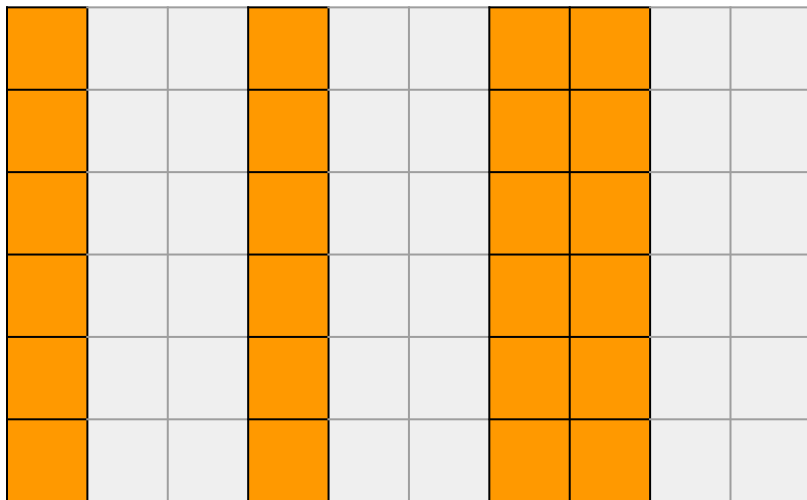
# HDF5 Hyperslab Unions

- More complex selections require unions of hyperslabs

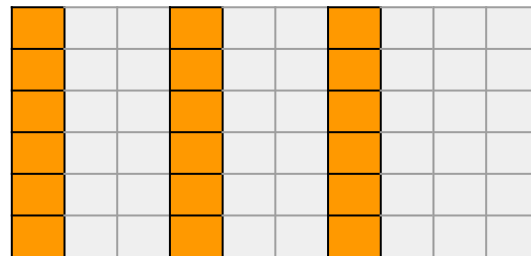


# HDF5 Hyperslab Unions

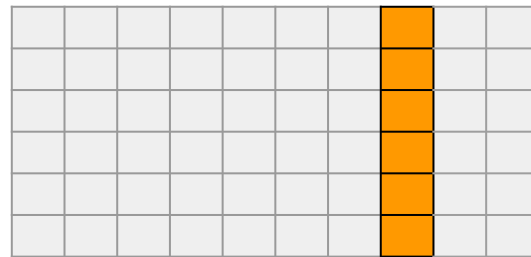
- More complex selections require unions of hyperslabs



=



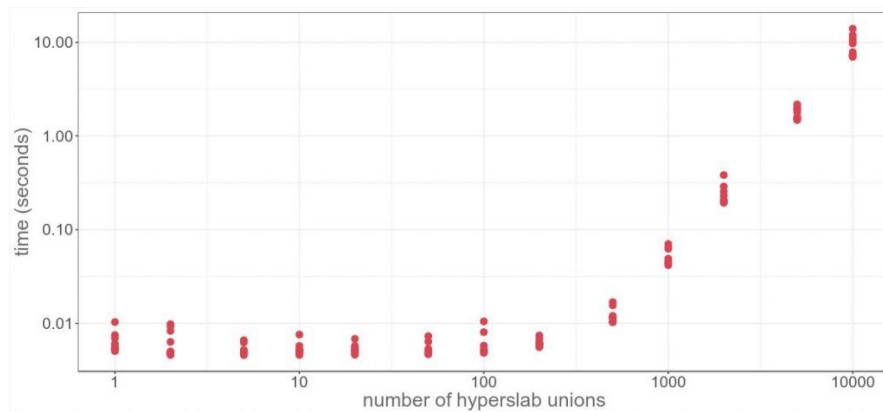
+



# HDF5 Hyperslab Unions

- Performing many unions gets very slow
- Acknowledged by HDF5 group - but no solution suggested e.g.

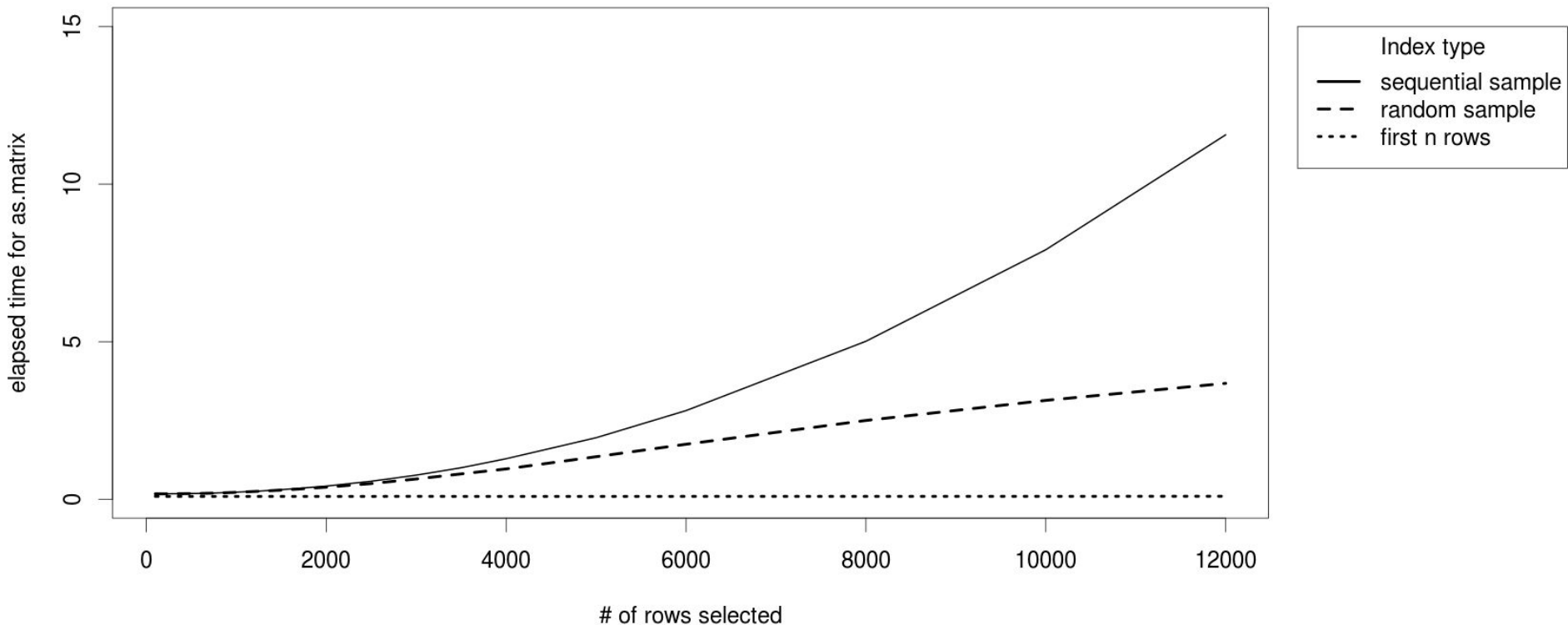
Union of non-consecutive hyperslabs is very slow





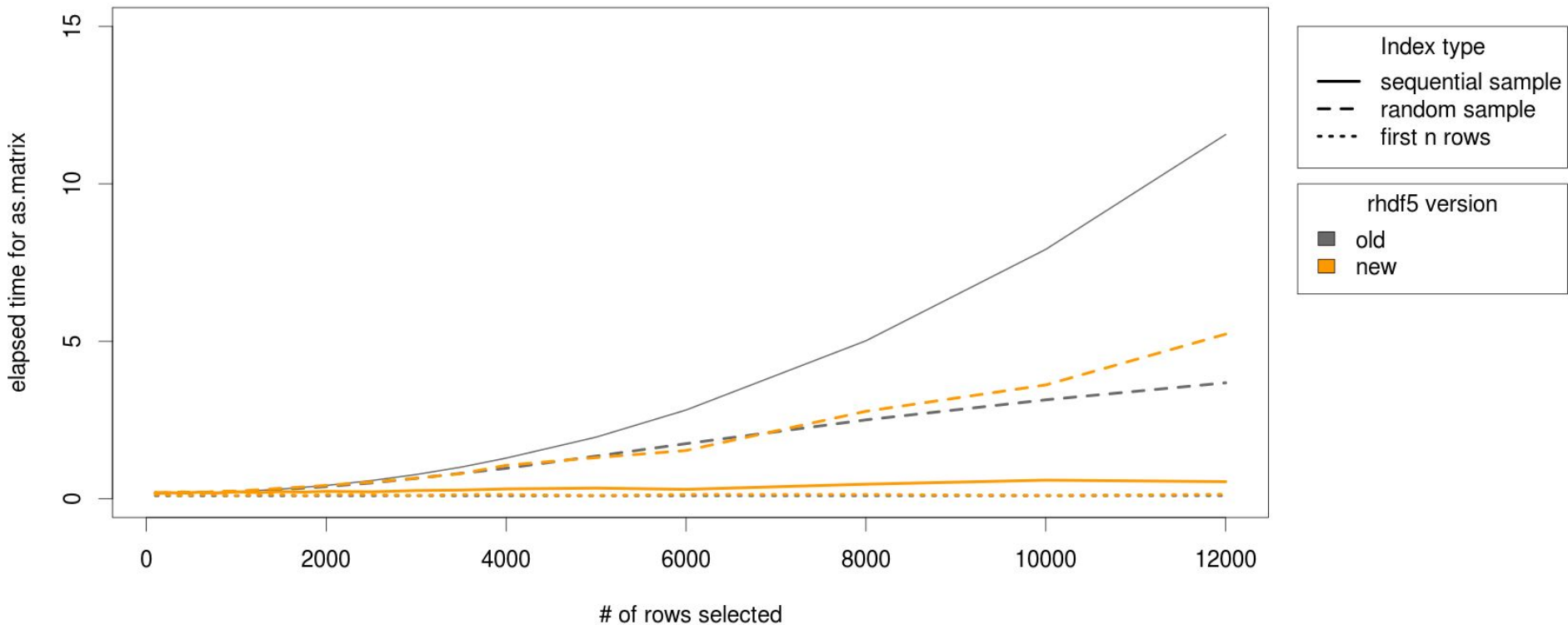


# Improvements to **rhdf5** indexing





# Improvements to **rhdf5** indexing



# Other additions

- Reading 'long' vectors e.g. Original 10x 1M Neuron h5 files
  - <https://github.com/grimbough/rhdf5/issues/8>
- Writing 'large' datasets
  - <https://github.com/grimbough/rhdf5/issues/30> & [#32](#)
- Tests & settings for file locking issues on Lustre & Solaris
  - `h5testFileLocking()`
  - `h5enableFileLocking()`
  - `h5disableFileLocking()`

# Expanding documentation

- Important to share knowledge / offer advice to users
  - Practical Tips [vignette](#)
  - Blog [posts](#)
- Other suggestions?

# Acknowledgements

Wolfgang Huber

Martin Morgan

Vince Carey

Daniel van Twisk

Hervé Pagès

Levi Waldron

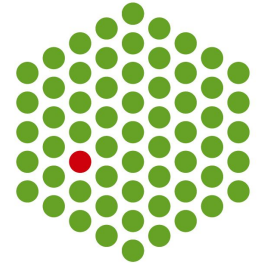
Aaron Lun

Mike Jiang

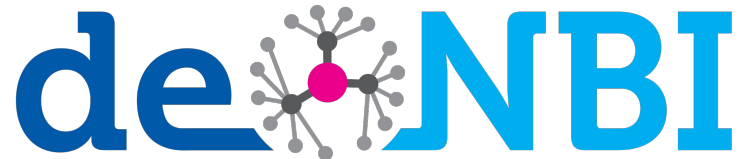


CHAN  
ZUCKERBERG  
INITIATIVE

EMBL



Bioconductor  
OPEN SOURCE SOFTWARE FOR BIOINFORMATICS



GERMAN NETWORK FOR BIOINFORMATICS INFRASTRUCTURE

