

CRBHits: Conditional Reciprocal Best Hits in R

Kristian K Ullrich¹

1 Max Planck Institute for Evolutionary Biology, Scientific IT group, August Thienemann Str. 2, 24306 Plön

DOI: [10.21105/joss.02424](https://doi.org/10.21105/joss.02424)

Software

- [Review](#) ↗
- [Repository](#) ↗
- [Archive](#) ↗

Editor: [Kristina Riemer](#) ↗

Reviewers:

- [@clauswilke](#)
- [@a-r-j](#)

Submitted: 26 May 2020

Published: 02 July 2020

License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC BY 4.0](#)).

Summary

CRBHits is a reimplementation of the Conditional Reciprocal Best Hit (CRBH) algorithm **crb-blast** in **R** (R Core Team, 2019). The new R package targets ecology, population and evolutionary biologists working in the field of comparative genomics.

The Reciprocal Best Hit (RBH) approach is commonly used in bioinformatics to show that two sequences evolved from a common ancestral gene. In other words, RBH tries to find orthologous protein sequences within and between species. These orthologous sequences can be further analysed to evaluate protein family evolution, infer phylogenetic trees and to annotate protein function (Altenhoff, Glover, & Dessimoz, 2019). The initial sequence search step is classically performed with the Basic Local Alignment Search Tool (blast) (Altschul, Gish, Miller, Myers, & Lipman, 1990) and due to evolutionary constraints, in most cases protein coding sequences are compared between two species. Downstream analysis use the resulting RBH to cluster sequence pairs and build so-called orthologous groups like e.g. **OrthoFinder** (Emms & Kelly, 2015) and other tools.

The CRBH algorithm was introduced by Aubry, Kelly, Kümper, Smith-Unna, & Hibberd (2014) and builds upon the traditional RBH approach to find additional orthologous sequences between two sets of sequences. As described earlier (Aubry et al., 2014; Scott, 2017), CRBH uses the sequence search results to fit an expect-value (e-value) cutoff given each RBH to subsequently add sequence pairs to the list of bona-fide orthologs given their alignment length.

Unfortunately, as mentioned by Scott (2017), the original implementation of CRBH (**crb-blast**) lag improved blast-like search algorithm to speed up the analysis. As a consequence, Scott (2017) ported CRBH to python **shmlast**, while **shmlast** cannot deal with IUPAC nucleotide code so far.

CRBHits constitutes a new R package, which build upon previous implementations and ports CRBH into the **R** environment, which is popular among biologists. **CRBHits** improve CRBH by additional implemented filter steps (Rost, 1999) and the possibility to apply custom filters.

Downstream functionalities

Calculating synonymous (dS) and nonsynonymous substitutions (dN) per orthologous sequence pair is a common task for evolutionary biologists, since its ratio dN/dS can be used as an indicator of selective pressure acting on a protein (Kryazhimskiy & Plotkin, 2008). However, this task is computational more demanding and consist of at least two steps, namely codon sequence alignment creation and dN/dS calculation. Further, the codon sequence alignment step consist of three subtasks, namely coding nucleotide to protein sequence translation, pairwise protein sequence alignment calculation and converting the protein sequence alignment back into a codon based alignment.

Downstream of CRBH creation, [CRBHits](#) features all above mentioned steps and subtasks. [CRBHits](#) has the ability to directly create codon alignments within R with the help of the widely used R package [Biostrings](#) (Pagès, Aboyoun, Gentleman, & DebRoy, 2017) (more than 200k downloads per year since 2014). These codon alignments can be subsequently used to calculate synonymous and nonsynonymous substitutions per sequence pair and is implemented in a multithreaded fashion either via the R package [seqinr](#) (Charif & Lobry, 2007) or the use of an R external tool [KaKs_Calculator2.0](#) (Wang, Zhang, Zhang, Zhu, & Yu, 2010).

Implementation

Like [shmlast](#), [CRBHits](#) benefits from the blast-like sequence search software [LAST](#) (Kiełbasa, Wan, Sato, Horton, & Frith, 2011) and plots the fitted model of the CRBH e-value based algorithm. In addition, users can filter the hit pairs prior to CRBH fitting for other criteria like query coverage, protein identity and/or the twilight zone of protein sequence alignments according to Rost (1999). The implemented filter uses equation 2 (see Rost, 1999):

$$f(x_{\text{hit pair}}) = \begin{cases} 100, & \text{for } L_{\text{hit pair}} < 11 \\ 480 * L^{-0.32 * (1 + e^{-\frac{L}{1000}})}, & \text{for } L_{\text{hit pair}} \leq 450 \\ 19.5, & \text{for } L_{\text{hit pair}} > 450 \end{cases}$$

where $x_{\text{hit pair}}$ is the expected protein identity given the alignment length $L_{\text{hit pair}}$. If the actual protein identity of a hit pair exceeds the expected protein identity ($\text{pident}_{\text{hit pair}} \geq f(x_{\text{hit pair}})$), it is retained for subsequent CRBH calculation.

In contrast to previous implementations, [CRBHits](#) only take coding nucleotide sequences (CDS) as the query and target inputs. This is due to the downstream functionality of [CRBHits](#) to directly calculate codon alignments within R, which rely on CDS. The inputs are translated into protein sequences, aligned globally (Smith, Waterman, & others, 1981) and converted into codon alignments.

Functions are completely coded in R and only the external prerequisites ([LAST](#) and [KaKs_Calculator2.0](#)) need to be compiled. Further, users can create their own filters before CRBH calculation.

Functions and Examples

The following example shows how to obtain CRBH between the coding sequences of *Schizosaccharomyces pombe* (Wood et al., 2012) and *Nematostella vectensis* (Apweiler, Bairoch, & Wu, 2004) by using two URLs as input strings and multiple threads for calculation.

```
library(CRBHits)
cds1 <- paste0("ftp://ftp.pombase.org/pombe/genome_sequence_and_features/",
               "feature_sequences/cds.fa.gz")
cds2 <- paste0("ftp://ftp.ebi.ac.uk/pub/databases/reference_proteomes/Qf0/",
               "Eukaryota/UP000001593_45351_DNA.fasta.gz")
#get help ?cdsfile2rbh
cds1.cds2.crbh <- cdsfile2rbh(cds1, cds2, plotCurve = TRUE, threads = 4)
```

Accept / Reject secondary hits as homologs

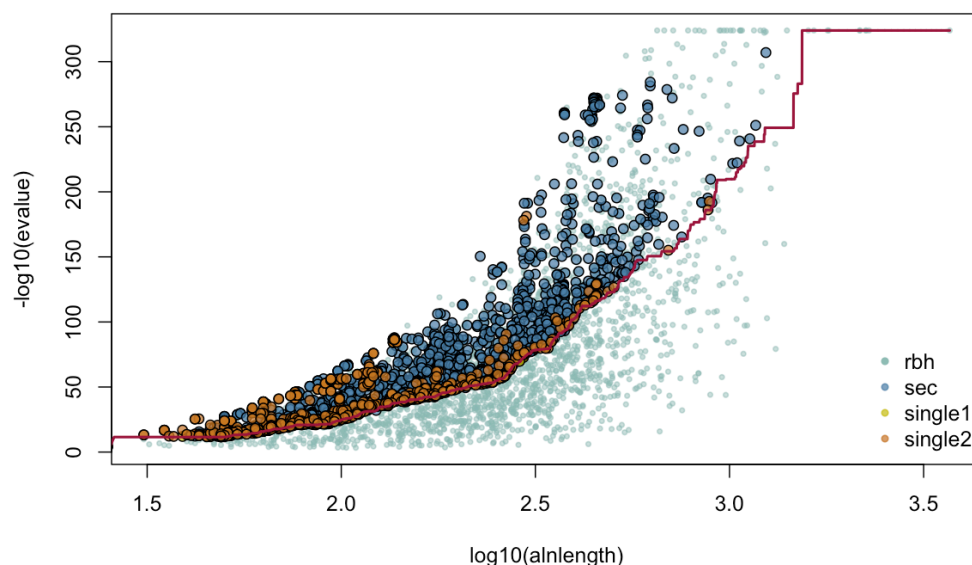


Figure 1: Accepted secondary reciprocal best hits based on CRBH fitting.

The obtained CRBH can also be used to calculate synonymous (dS/ks) and nonsynonymous (dN/ka) substitutions per hit pair using either the model from Li (1993) or from Yang & Nielsen (2000).

```
cds1 <- Biostrings::readDNAStringSet(cds1)
cds2 <- Biostrings::readDNAStringSet(cds2)
#get help ?rbh2kaks
cds1.cds2.kaks.Li <- rbh2kaks(cds1.cds2.crbh$crbh.pairs, cds1, cds2,
                             model = "Li", threads = 4)
cds1.cds2.kaks.YN <- rbh2kaks(cds1.cds2.crbh$crbh.pairs, cds1, cds2,
                             model = "YN", threads = 4)
```

Table 1: Performance comparison for CRBH and dN/dS calculations (Intel Xeon CPU E5-2620 v3 @ 2.40GHz; 3575 hit pairs).

Number of Threads	1	2	4	8
Runtime of CRBH(shmlast) in sec	38 (s)	30 (s)	28 (s)	28 (s)
Runtime of CRBH(CRBHits) in sec	32 (s)	26 (s)	24 (s)	22 (s)
Runtime of kaks.Li in sec	357 (s)	167 (s)	87 (s)	49 (s)
Runtime of kaks.YN in sec	474 (s)	230 (s)	121 (s)	63 (s)

Conclusions

CRBHits implements CRBH in R (see Figure 1) and also can be used to calculate codon alignment based nucleotide diversities in a multithreaded fashion (see Table 1).

Availability

CRBHits is an open source software made available under the MIT license. It can be installed from its gitlab repository using the [devtools](#) package.

```
devtools::install_gitlab("mpievolbio-it/crbhits",  
  host = "https://gitlab.gwdg.de", build_vignettes = TRUE)
```

The R package website, which contain a detailed HOWTO to install the prerequisites (mentioned above) and package vignettes are available at <https://mpievolbio-it.pages.gwdg.de/crbhits>.

References

- Altenhoff, A. M., Glover, N. M., & Dessimoz, C. (2019). Inferring orthology and paralogy. In *Evolutionary genomics* (pp. 149–175). Springer. doi:[10.1007/978-1-4939-9074-0_5](https://doi.org/10.1007/978-1-4939-9074-0_5)
- Altschul, S. F., Gish, W., Miller, W., Myers, E. W., & Lipman, D. J. (1990). Basic local alignment search tool. *Journal of molecular biology*, 215(3), 403–410. doi:[10.1016/S0022-2836\(05\)80360-2](https://doi.org/10.1016/S0022-2836(05)80360-2)
- Apweiler, R., Bairoch, A., & Wu, C. H. (2004). Protein sequence databases. *Current opinion in chemical biology*, 8(1), 76–80. doi:[10.1016/j.cbpa.2003.12.004](https://doi.org/10.1016/j.cbpa.2003.12.004)
- Aubry, S., Kelly, S., Kümpers, B. M., Smith-Unna, R. D., & Hibberd, J. M. (2014). Deep evolutionary comparison of gene expression identifies parallel recruitment of trans-factors in two independent origins of c4 photosynthesis. *PLoS genetics*, 10(6). doi:[10.1371/journal.pgen.1004365](https://doi.org/10.1371/journal.pgen.1004365)
- Charif, D., & Lobry, J. R. (2007). SeqinR 1.0-2: A contributed package to the r project for statistical computing devoted to biological sequences retrieval and analysis. In *Structural approaches to sequence evolution* (pp. 207–232). Springer. doi:[10.1007/978-3-540-35306-5_10](https://doi.org/10.1007/978-3-540-35306-5_10)
- Emms, D. M., & Kelly, S. (2015). OrthoFinder: Solving fundamental biases in whole genome comparisons dramatically improves orthogroup inference accuracy. *Genome biology*, 16(1), 157. doi:[10.1186/s13059-015-0721-2](https://doi.org/10.1186/s13059-015-0721-2)
- Kiełbasa, S. M., Wan, R., Sato, K., Horton, P., & Frith, M. C. (2011). Adaptive seeds tame genomic sequence comparison. *Genome research*, 21(3), 487–493. doi:[10.1101/gr.113985.110](https://doi.org/10.1101/gr.113985.110)
- Kryazhimskiy, S., & Plotkin, J. B. (2008). The population genetics of dN/dS. *PLoS genetics*, 4(12). doi:[10.1371/journal.pgen.1000304](https://doi.org/10.1371/journal.pgen.1000304)
- Li, W.-H. (1993). Unbiased estimation of the rates of synonymous and nonsynonymous substitution. *Journal of molecular evolution*, 36(1), 96–99. doi:[10.1007/BF02407308](https://doi.org/10.1007/BF02407308)
- Pagès, H., Aboyoun, P., Gentleman, R., & DebRoy, S. (2017). Biostrings: Efficient manipulation of biological strings. *R package version*, 2(0). doi:[10.18129/B9.bioc.Biostrings](https://doi.org/10.18129/B9.bioc.Biostrings)
- R Core Team. (2019). R: A language and environment for statistical computing. Retrieved from <https://www.r-project.org/>
- Rost, B. (1999). Twilight zone of protein sequence alignments. *Protein engineering*, 12(2), 85–94. doi:[10.1093/protein/12.2.85](https://doi.org/10.1093/protein/12.2.85)

- Scott, C. (2017). Shmblast: An improved implementation of conditional reciprocal best hits with last and python. *Journal of Open Source Software*, 2(9), 142. doi:[doi:10.21105/joss.00142](https://doi.org/10.21105/joss.00142)
- Smith, T. F., Waterman, M. S., & others. (1981). Identification of common molecular subsequences. *Journal of molecular biology*, 147(1), 195–197. doi:[10.1016/0022-2836\(81\)90087-5](https://doi.org/10.1016/0022-2836(81)90087-5)
- Wang, D., Zhang, Y., Zhang, Z., Zhu, J., & Yu, J. (2010). KaKs_Calculator 2.0: A toolkit incorporating gamma-series methods and sliding window strategies. *Genomics, proteomics & bioinformatics*, 8(1), 77–80. doi:[10.1016/S1672-0229\(10\)60008-3](https://doi.org/10.1016/S1672-0229(10)60008-3)
- Wood, V., Harris, M. A., McDowall, M. D., Rutherford, K., Vaughan, B. W., Staines, D. M., Aslett, M., et al. (2012). PomBase: A comprehensive online resource for fission yeast. *Nucleic acids research*, 40(D1), D695–D699. doi:[10.1093/nar/gkr853](https://doi.org/10.1093/nar/gkr853)
- Yang, Z., & Nielsen, R. (2000). Estimating synonymous and nonsynonymous substitution rates under realistic evolutionary models. *Molecular biology and evolution*, 17(1), 32–43. doi:[10.1093/oxfordjournals.molbev.a026236](https://doi.org/10.1093/oxfordjournals.molbev.a026236)