Become the Maestro of your Genomics Workflow with Bioconductor and Microsoft R Server

Colby Ford

Microsoft R Server is an enterprise-class tool for hosting and managing parallel and distributed workloads of R processes on servers. Organizations that need to process large amounts of data or perform complex processing on the data benefit the most from a parallel architecture like Microsoft R Server. It uses the Revoscaler package, which makes parallelization easy.

In genomics research, we often interact with large amounts of data from complex pipelines in a diverse array of formats. Luckily, <u>Bioconductor</u> helps make this process simpler by packaging up common sets of processes in ready-to-use R code.

Harnessing the power of both Bioconductor and Microsoft R Server together can help streamline the processing of your genomics data.

All About Bioconductor



Bioconductor is an open-source, open-development software project to provide tools for the analysis and comprehension of high-throughput genomic data. It is based primarily on the R programming language. In other words, it's an extension of R that is specialized for bioinformatics and genomics analyses.

Installation

Once R (either base R or Microsoft R Server) is installed on your local machine, installing Bioconductor is simple. Open R and use the following commands to grab the latest version of Bioconductor.

```
source("https://bioconductor.org/biocLite.R")
biocLite()
```

Once the biocLite script has loaded, you can now call any desired packages.

```
library("BiocInstaller")
biocLite("RforProteomics", dependencies = TRUE)
```

Now you are ready to use over 1,000 Bioconductor packages!

Workflows

The field of genomics is very broad, but Bioconductor will often have a solution for every

area. The Bioconductor site is rich with workflow examples to help connect the dots for your research. Check out <u>Bioconductor's help section</u> for a list of the available workflows.

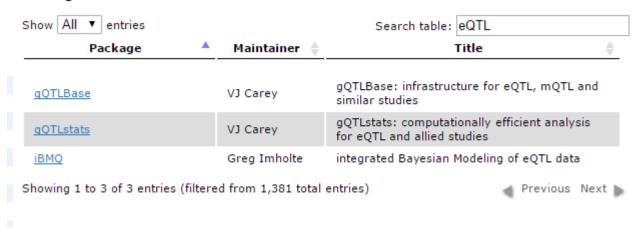
Here are a few of my favorites:

- <u>Sequence Analysis</u> Import fasta, fastq, BAM, gff, bed, wig, and other sequence formats. Trim, transform, align, and manipulate sequences. Perform quality assessment, ChIP-seq, differential expression, RNA-seq, and other workflows.
- <u>Variant Annotation</u> Read and write VCF files. Identify structural location of variants and compute amino acid coding changes for non-synonymous variants and predict consequence of amino acid coding changes.
- <u>High Throughput Assays</u> Import, transform, edit, analyze and visualize flow cytometric, mass spec, HTqPCR, cell-based, and other assays.
- <u>Transcription Factor Binding</u> Find candidate binding sites for known transcription factors via sequence matching.
- <u>Cancer Genomics</u> Download, process, and prepare <u>TCGA</u>, <u>ENCODE</u>, and <u>Roadmap</u> data to interrogate the epigenome of cultured cancer cell lines as well as normal and tumor tissues with high genomic resolution.

Package Database

Didn't see a workflow that exactly fit your needs? No problem! Bioconductor has a nice list of available packages that will help you find the right one for the problem at hand. Using the search box shown below, I searched for the term "*eQTL*" (Expression Quantitative Trait Loci) to find packages related to that topic.

Packages found under Software:



For more information, you can check out the package list <u>here</u>.

Whether you are using the pre-made workflows or ended up creating your own, you can likely speed up processing time by running your Bioconductor/R scripts in parallel. Microsoft R Server and RevoscaleR make this easy.

Let's take the **Annotating Genomic Variants** workflow, for example.

.vcf files are often very large and sometimes difficult to process or summarize due to their size. Using the <code>VariantAnnotation::locateVariants</code> function from Bioconductor makes this process more automated. We can use this function to identify where a variant falls with respect to gene structure, e.g., exon, utr, splice site, etc. We use the gene model from the <code>TxDb.Hsapiens.UCSC.hg19.knownGene</code> package loaded earlier.

```
## of interest. See ?locateVariants for details.
cds <- locateVariants(vcf, txdb, CodingVariants())
five <- locateVariants(vcf, txdb, FiveUTRVariants())
splice <- locateVariants(vcf, txdb, SpliceSiteVariants())
intron <- locateVariants(vcf, txdb, IntronVariants())
all <- locateVariants(vcf, txdb, AllVariants())</pre>
```

If we want to start summarizing the variants, we could use sapply to repetitively perform some operation over the entire data object. Take a look at the highlighted lines below.

This code easily summarizes the .vcf file in a few seconds. However, the NAO6985_17.vcf.gz file is only a small (35MB) sample from human Chromosome 17. What if you were to use multiple files to assess a sample population such as the ones available from 1000 Genomes? sapply might take a while...

We can use the Revoscaler package to parallelize the summarization function in the code. By using <code>rxexec</code> to distribute the processing over multiple cores of a processor or even multiple nodes on a Hadoop cluster, we can speed up the processing time tremendously. In the sample code below, we use the <code>rxexec</code> function to split up the processing by <code>GeneID</code>.

```
## Summarize variant location by gene using rxExec from Microsoft R Server:
vcflocationsummary <- function(nm) {
         d <- all[mcols(all)$GENEID %in% nm, c("QUERYID", "LOCATION")]
         table(mcols(d)$LOCATION[duplicated(d) == FALSE])
    }
rxExec(vcflocationsummary, rxElemArg(GENEID))</pre>
```

Note: this is only sample code. To fully use this workflow, visit <u>Bioconductor's workflow</u> variants.

Try it Out!

Now that we have explored how easy it is to speed up your genomics workflows using Microsoft R Server, you can try it out for yourself. Pick a workflow that fits your needs and then use it. Once you start seeing where the processing bottlenecks are, think about using Revoscaler's parallelization functions to speed things up. Look for loops and apply functions as prime candidates for distributed processing.

First time using Microsoft R Server and the RevoscaleR? <u>Microsoft's documentation</u> is a great place to start. To compare the RevoScaleR functions, read <u>Explore R and RevoScaleR</u> <u>in 25 functions</u>. If you still have questions, please <u>reach out to us</u> and we will be happy to help!