

Run EQTL with batchtools SSH cluster functions

Credentials

By default, boto3 loads credentials from `~/.aws/credentials`. To use different credentials, set the environment variable:

```
1 export AWS_SHARED_CREDENTIALS_FILE=/Users/spollack/.aws/credentials.new-group-acct
```

AMI

Public AMI `ami-e3cobf99` is Ubuntu 16.04, with R-3.4.3, Bioconductor 3.6, batchtools and the EQTL dependencies installed. It was made with these commands:

```
1 sudo apt-get update
2 sudo apt-get install -y r-base-core                # brings in the R dependencies
3 sudo apt-get install -y libcurl4-openssl-dev       # devtools dependency
4 sudo apt-get install -y libssl-dev                 # devtools dependency
5 sudo apt-get install -y nfs-kernel-server          # NFS server (also EFS dependency)
6 sudo apt-get install -y libxml2-dev                # XML dependency
7 sudo apt-get install -y default-jdk                # R-3.4.3 dependency
8 sudo apt-get install -y libmariadb-client-lgpl-dev # RMySQL dependency
9
10 # Install R-3.4.3 from source
11 wget https://cran.r-project.org/src/base/R-3/R-3.4.3.tar.gz
12 gunzip R-3.4.3.tar.gz
13 tar xvf R-3.4.3.tar
14 cd R-3.4.3
15 sudo ./configure --prefix=/usr/local/lib/R-3.4.3 --with-x=no
16 sudo make all
17 sudo make install
18
19 # Edit /etc/environment to put R-3.4.3 in the path:
20 PATH="/usr/local/lib/R-
3.4.3/bin:/usr/local/sbin:/usr/local/bin:/usr/sbin:/usr/bin:/sbin:/bin:/usr/games:/usr/local
/games"
21 # Comment out env_reset and secure_path (sudo visudo) so sudo also finds R-3.4.3
22
23 sudo /usr/local/lib/R-3.4.3/bin/R                # install to system library
24 > install.packages("devtools", repos="http://lib.stat.cmu.edu/R/CRAN")
25 > install.packages("XML", repos="http://lib.stat.cmu.edu/R/CRAN")      # rtracklayer
dependency
```

```

26
27 > source("https://bioconductor.org/biocLite.R")
28 > biocLite();
29 > biocLite("gQTLstats"); # ldblock and VariantAnnotation installed with this
30 > biocLite("geuvPack")
31
32 # use sampoll/batchtools to avoid the SSH linux-helper path quoting problem.
33 > library(devtools)
34 > install_github("sampoll/batchtools")

```

Note: The quoting in Worker.R at line 49 does not work on this system. (It is an open question whether it works on any other system.) To make this work, I use the path to the `linux-helper` script in the local build, on the assumption that all of the R source trees are the same, which is true as long as we are using the same AMI.

This means we need to install and use the batchtools from my GitHub repository.

Configuration

Bioconductor is too large for a t2.micro, so we use a t2.small. A sample config file looks like this. The config file name is the first command line argument.

```

1 [clusterdef]
2 name = cluster-eqtl
3 vpc = vpc-576d222f
4 nnode = 2
5 key = sam-new-account
6 type = t2.small
7 ami = ami-ec30bf99

```

The cluster has to be set up in two parts, because the nodes are not ready to accept connections on the ssh port until some time after boto3 classifies them as "running." The second command line argument is "1" for the first part of the setup (launch nodes) and "2" for the rest (once the nodes are ready for ssh.)

Setting nnode to 2 will result in one head node (which does not do work) and one compute node.

(A possible solution: <https://stackoverflow.com/questions/7405598/process-for-telling-when-a-new-ec2-host-can-be-connected-to>)

This cluster will use NFS. (Configuration file with EFS is below.) When the cluster is set up, `run.py` prints out the IP Addresses. For example:

```

1 $ python run.py eqtl.nfs.config 1
2 # wait until console shows node status checks are done
3 $ python run.py eqtl.nfs.config 2
4 # program prints output, ending in
5 Head node: public IP Address: 54.159.134.146 Private IP Address: 172.31.92.105
6 Compute node: public IP Address: 54.91.232.33 Private IP Address: 172.31.80.128

```

Run EQTL

ssh to the head node and cd to /scratch. The batchtools.conf.R file should already be there. sftp up eqtl.R:

```
1 # eqtl.R
2 library(gQTLstats)
3 library(ldblock)
4 library(VariantAnnotation)
5 library(geuvPack)
6 library(batchtools)
7
8 data(geuFPKM)
9 ss <- stack1kg()
10 v17 <- ss@files[[17]]
11 someGenes <- c("ORMDL3", "GSDMB", "IKZF3", "MED24", "CSF3", "ERBB2",
12               "GRB7", "MIEN1", "GSDMA", "THRA", "MSL1")
13 se17 <- geuFPKM[ which(rowData(geuFPKM)$gene_name %in% someGenes),]
14
15 n <- 10
16 vr0 <- 39.5e6
17 vr1 <- 40.5e6
18 v <- seq(vr0, vr1, length=n+1)
19 vl <- zipup(v[1:n], v[2:(n+1)]-1)
20
21 run.job <- function(v, se, vcf) {
22   library(gQTLstats)
23   library(ldblock)
24   library(VariantAnnotation)
25   library(geuvPack)
26
27   vr <- GRanges("17", IRanges(start=v[1], end=v[2], names=c("range")))
28   results <- AllAssoc(se, vcf, vr)
29 }
30
31 concat.job <- function(gr1, gr2) {
32   gr <- c(gr1, gr2)
33 }
```

Run from R command prompt:

```
1 R
2 > source("eqtl.R")
3 > system("rm -Rf ./registry") # in case there is an old registry
4 > reg <- makeRegistry();
5 > batchMap(fun = run.job, as.list(vl), more.args = list(se=se17, vcf=v17))
6 > submitJobs()
7 > waitForJobs()
8 > all.results <- reduceResults(fun = concat.job)
9 > save(all.results, file="all.results.RData")
```

Run with EFS

The EFS is `fs-e9a39da0`, in `VPC-5763222f`. DNS name: `fs-e9a39da0.efs.us-east-1.amazonaws.com`

There are two security groups for the EFS; one for the mount targets and one for the nodes that will mount it. The node SG is called `bioc-efs-node-sg` and the mount SG is called `bioc-efs-mount-sg`. on port 22 from anywhere. `bioc-efs-mount-sg` allows incoming NFS traffic on port 2049 from `bioc-cfn-sg`.

In the EFS there is a top-level directory called `/btrun` which has group and user owner `ubuntu:ubuntu` and read-write-execute permissions set to `775`. This will allow any node based on Ubuntu to read from and write to a batchtools registry.

(`bioc-efs-node-sg` allows incoming traffic on any port from `bioc-efs-mount-sg` and)

Config File

```
1 [clusterdef]
2 name = cluster-eqtl
3 vpc = vpc-576d222f
4 nnode = 3
5 key = sam-new-account
6 type = t2.small
7 ami = ami-e3c0bf99
8 sgn = bioc-efs-node-sg
9 efs = fs-e9a39da0.efs.us-east-1.amazonaws.com
```

Run EQTL

ssh to the head node and cd to `/efs/btrun.eqtl`. R is already there. Execute the same commands at the R command prompt as in the NFS cluster.

Note:

When submitting jobs, SSH asks for confirmation to continue connecting once per host. This is annoying. Something should go into the `known_hosts` file to prevent this.