# Instructions for controlled access to TOPMed sequence data on the cloud

Tom Blackwell - University of Michigan - October 21, 2019

There are six separate components which need to be lined up in advance.

## Permissions:

(1)  An NIH eRA Commons login.

(2)  An approved dbGaP Data Access Request ("DAR") for the TOPMed study or studies you wish to use.

(3)  If someone else made the dbGaP Data Access Request, ask them to assign you as a "Designated Downloader" for the Data Access Request.

(4)  An ".ngc" key file from the dbGaP "My Research Projects" page.  To obtain this, follow the separate instructions for sample identifiers in item (6) below.

(5)  Access to either a Google or Amazon billing account.  The dbGaP access mechanism is designed to run identically on either platform, and TOPMed sequence data are stored in duplicate on both platforms.

## Sample Identifiers:

(6)  dbGaP access to whole genome sequence data identifies sequenced DNA samples by their NCBI Sequence Read Archive "SRR" run accession numbers.  The correspondence to familiar TOPMed "NWD" sample identifiers is shown in the NCBI Sequence Read Archive "Run Selector".  This is a java web interface to NCBI's internal database and it runs in a browser on the user's desktop.  Separate, detailed instructions are provided in an appendix below,  "Listing TOPMed sequence data available in dbGaP".  These instructions give the only way I know to obtain an ".ngc" key file.  It would be difficult to script access to the Run Selector, or to access it directly from a cloud instance.  The dbGaP access mechanism also allows combining TOPMed data with sequences from other studies.  After locating the samples of interest, create a plain text file showing all of their SRR numbers, one SRR number per line.

## Cloud Access:

(7)  The fusera software is available from  github.com/mitre/fusera.  Detailed usage and troubleshooting instructions are in a wiki:  github.com/mitre/fusera/wiki.

(8)  Alternately, we provide a Docker machine image with samtools and the fusera software already installed:  https://hub.docker.com/r/statgen/cram-access-tools,  currently running fusera v-1.0.0.  This also contains the build 38 human genome reference sequence, which is needed in order to read each .cram file.

(9)  Start a computing instance on either Google or Amazon.  From my desktop, this is:

```
gcloud init
gcloud compute instances create  fusera-trial-00  --zone  us-central1-c      \
    --image-project  ubuntu-os-cloud  --image-family  ubuntu-1604-lts        \
    --create-disk  size=50GB
```

If you are using Google's web interface instead of the command line, do not check "Container", do check "Allow default access" and do not check either box under "Firewall".

(10)  ssh into the running compute instance, create a working directory and a subdirectory, say '/working' and '/working/genomes', and copy both the '*.ngc' key file and the list of SRR numbers from steps 4 and 6 into the working directory on the compute instance from your desktop.

(11)  Install Docker-CE ("Community Edition") on the compute instance, following instructions from      https://docs.docker.com/install/linux/docker-ce/ubuntu/#set-up-the-rep      down   through 'sudo docker run hello-world'.

(12)  Copy the statgen docker image from docker hub to the compute instance and start it with special permissions:

```
sudo docker pull  statgen/cram-access-tools
sudo docker run --rm -it -v /working:/working --privileged --cap-add SYS_ADMIN   \
    --cap-add MKNOD --device /dev/fuse  statgen/cram-access-tools
```

(13)  Run fusera inside the docker container as:

```
nohup fusera mount --verbose  --ngc  <ngc.key.file>  --accession  <srr.list>     \
    /working/genomes  >  /working/fusera.log  2>&1  &
```

<srr.list> may be either a comma-separated list on the command line with one or more literal SRR numbers or a path to a text file of SRR numbers.  The log file collects any error messages.

(14)  This creates a 'fuse' virtual file system in the Docker image, mounted at /working/genomes, which contains the sequence data in a directory structure like:

```
/working/genomes/
     SRR12345678/
         NWD113355.b38.irc.v1.cram
         NWD113355.b38.irc.v1.cram.crai
     SRR12345679/
         NWD224466.b38.irc.v1.cram
         NWD224466.b38.irc.v1.cram.crai
     SRR12345680/
         NWD335577.b38.irc.v1.cram
         NWD335577.b38.irc.v1.cram.crai
```

The fusera process must stay running in the background while other tools access the sequence data.

# Appendix: Listing TOPMed sequence data available in dbGaP

This page provides expanded instructions to list the NCBI "SRR" accession numbers associated with TOPMed samples and to download a "*.ngc" key file which gives access to these data on the cloud. This is step (6) in the preceding "Instructions for controlled access to TOPMed sequence data on the cloud".
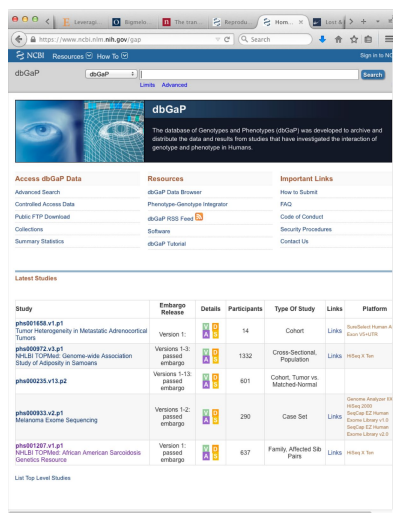
Start at the dbGaP main page: www.ncbi.nlm.nih.gov/gap            (panel A)
Select "Controlled Access Data" (second item in the left hand column).
This resolves to: dbgap.ncbi.nlm.nih.gov/aa/wga.cgi?page=login   (panel B)

The tab "Authorized Access" is already open. Select "Log In to dbGaP". (This is on the left, below the yellow banner and above "dbGaP Data Download".) It takes you to the login page with a very long address beginning:
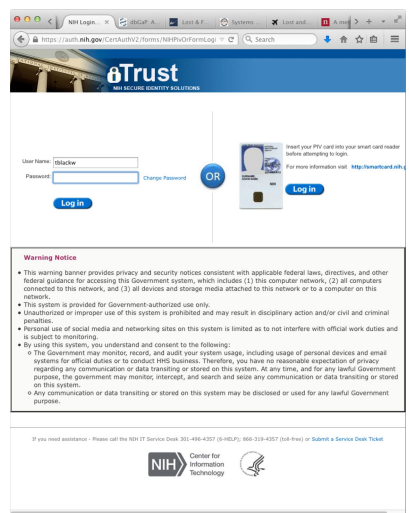auth.nih.gov/CertAuthV2/forms/NIHPivOrFormLogin.aspx?...            (panel C)



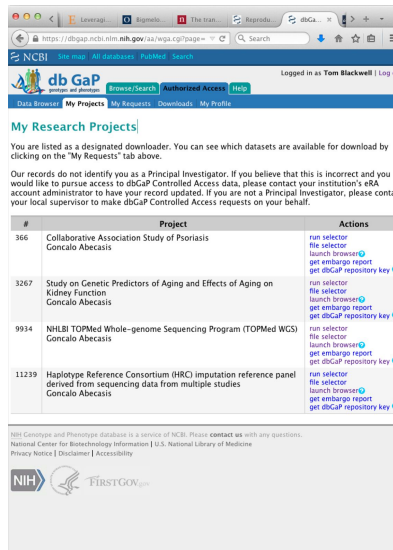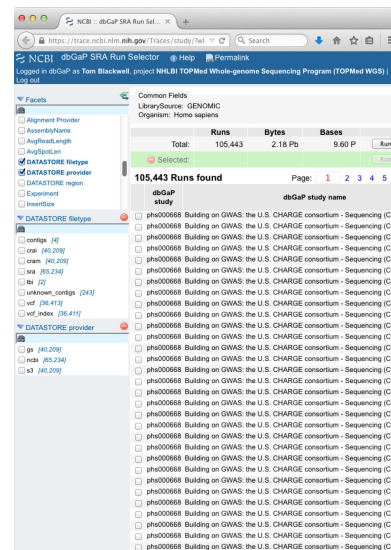A                                    B                                    C

Provide your eRA commons user name and password. This takes you to:
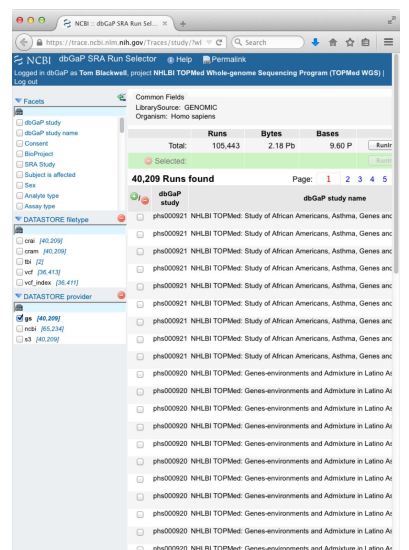dbgap.ncbi.nlm.nih.gov/aa/wga.cgi/page=list_wishlists with the "Authorized Access" and "My Projects" tabs already open (panel D).



D                                    E                                    F

From here on, the screenshots will inevitably show what I see and not exactly what you will see. TOPMed is the third among the four projects shown in panel D. Before going further, click "get dbGaP repository key" for the TOPMed project (last item in the right hand column for the TOPMed project). This generates and downloads to your desktop the "*.ngc" key file needed for cloud access to TOPMed sequence data. Set local permissions on this file so that only you can read it. Giving this file to anyone else will give them access to the sequence data.

Next, select "run selector", first item in the right hand column for the TOPMed project. This takes you to the "dbGaP SRA Run Selector" trace.ncbi.nlm.nih.gov/Traces/study/?wlid=... .

This page (panel E) may take 60 seconds or more to load. Always switch to the "old" Run Selector if the new one fails to find any data. For me, as of December 2018, it returns a table of 105,440 rows and 44 columns. In order to negotiate this table, use faceted search from the left hand column of the web page. The top left panel, "Facets", with a small vertical scrollbar, shows all 44 column names. In panel E, both "DATASTORE filetype" and "DATASTORE provider" are checked. Each one opens a panel below showing all possible data values in that column and the number of records with each value.

TOPMed sequence data available on the cloud are exclusively .cram files with associated .crai index files. Files with .sra filetype are TOPMed phase 1 data mapped to build 37 plus older ESP exome sequence data from some TOPMed parent studies. These are hosted by NCBI in .sra format. .vcf and .vcf_index filetypes refer to TOPMed freeze 5b genotype data in single-sample .vcf format provided in March 2018 as an experiment for the NHLBI Data STAGE developers.

The "DATASTORE provider" tab shows that the same 40,209 records are available from both "gs" and "s3", while 65,234 (.sra) records are available from "ncbi".

Checking the box for "gs" (shown as panel F) restricts the display to just the 40,209 records with that provider. Then pressing "+" in the pair of green and red buttons left of the column headings in the main table, just below "40,209 Runs found", selects all 40,209 records, turns them green and enables you to select "RunInfo Table" above (panel G). This downloads a tab-delimited file of 40,209 rows x 44 columns to your desktop, from which you can select the studies and samples of interest. Their NWD identifiers are shown in column 27 "Sample name" and the corresponding "SRR" run accession numbers are shown in column 24 "Run" (panel H). These SRR identifiers are what you will need for cloud access to the sequence data.



G                                    H                                    I