# How To: Run the ENCODE long-RNA-seq analysis pipeline on DNAnexus

**Overview:** In this exercise, we will run the ENCODE Uniform Processing Long RNA-seq Pipeline on a small test dataset containing reads from chromosome 21 sampled from an ENCODE RNA-seq experiment on a stomach tissue sample.

The ENCODE Portal page for the experiment is here: (https://www.encodeproject.org/experiments/ENCSR000AFI/)

The pipeline was specified by the ENCODE RNA Working Group and implemented at the ENCODE Data Coordinating Center (DCC). Today we will run the pipeline on the DNAnexus cloud platform. Typically, full ENCODE RNA experiments run on this pipeline are whole genome 30x read depth and take around 10 hours. This demonstration dataset can be processed in about 46 minutes.
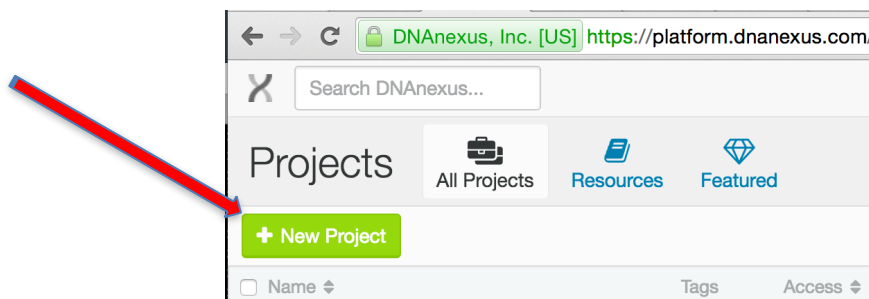
The ENCODE pipeline code is open-source and lives on github at: https://github.com/ENCODE-DCC/long-rna-seq-pipeline. The pipeline is modeled on the encode portal which provides links directly to the exact scripts that define each step: https://www.encodeproject.org/pipelines/ENCPL002LPE/.

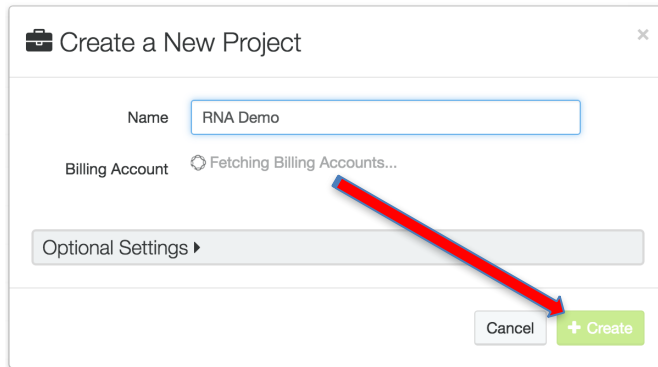**Summary of Steps:** Here is a high-level summary of what you will learn to do in this exercise.

- **Find** the ENCODE Uniform Processing Pipeline project on DNAnexus.
- **Copy** the pipeline software and files from that project to a new project in your account.
- **Complete** the specification of inputs to the workflow.
- **Run** the pipeline workflow on the cloud.
- **Monitor** the run's progress.
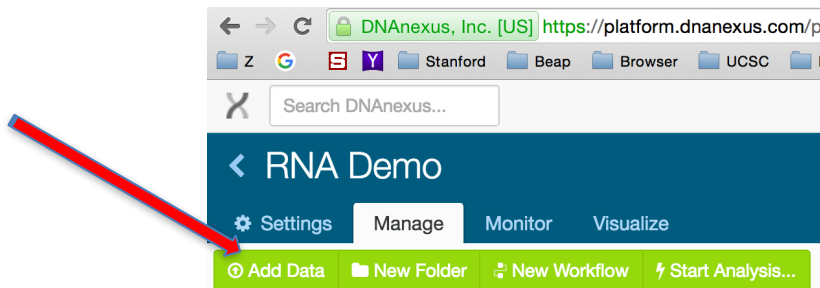- **Visualize** the output.

**Step-by-step:**

1) You will need to create an account on the DNAnexus website www.dnanexus.com. Log in to your DNAnexus account.

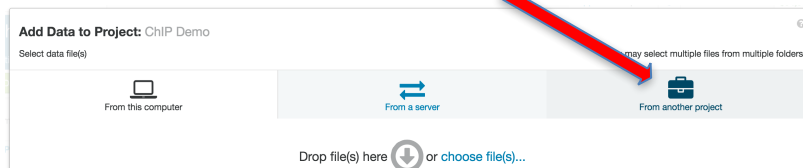2) Once logged into your DNAnexus account, create a new project. Select "All Projects" and then click "New Project":
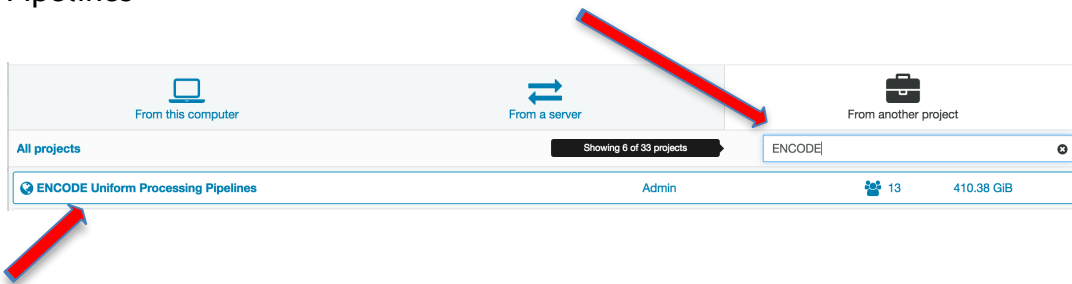
3) Give your project a new name and click "Create".
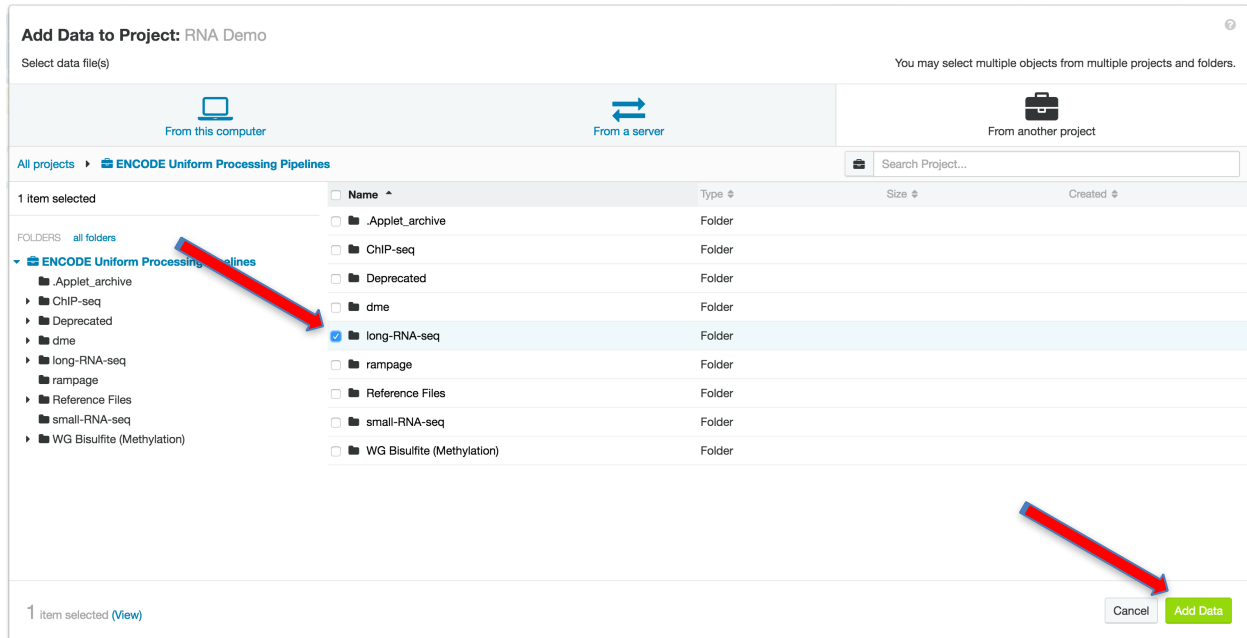


4) Select "Add Data" ...
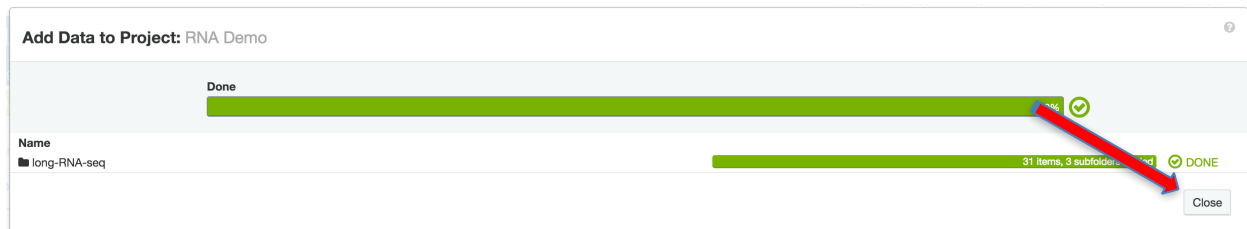


5) ... select "From another project" ...



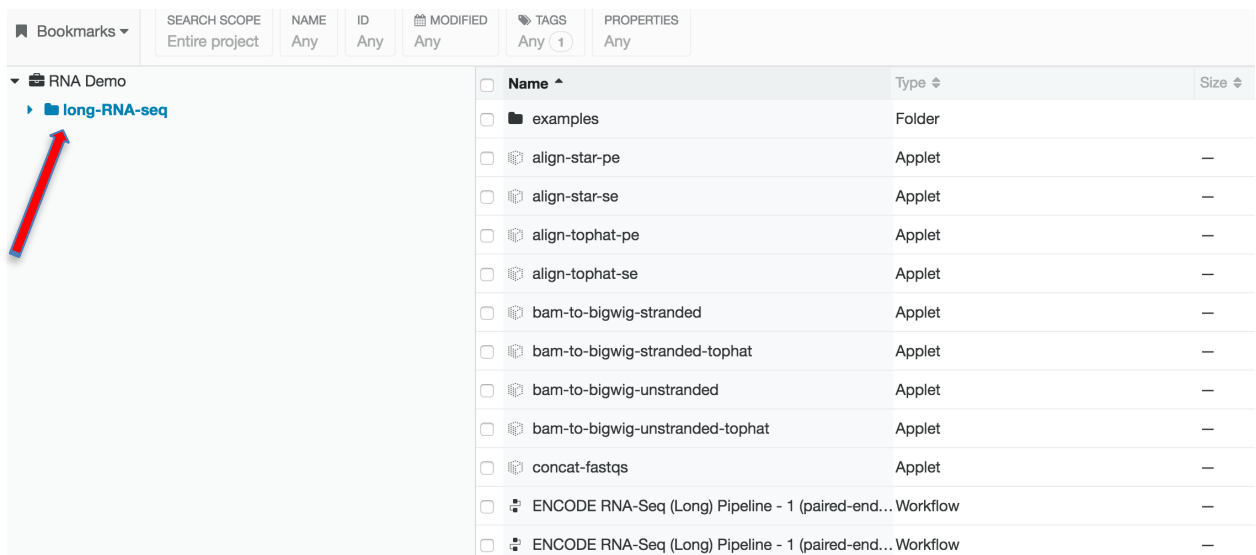6) Type "ENCODE" in the search box and then select "ENCODE Uniform Processing Pipelines"

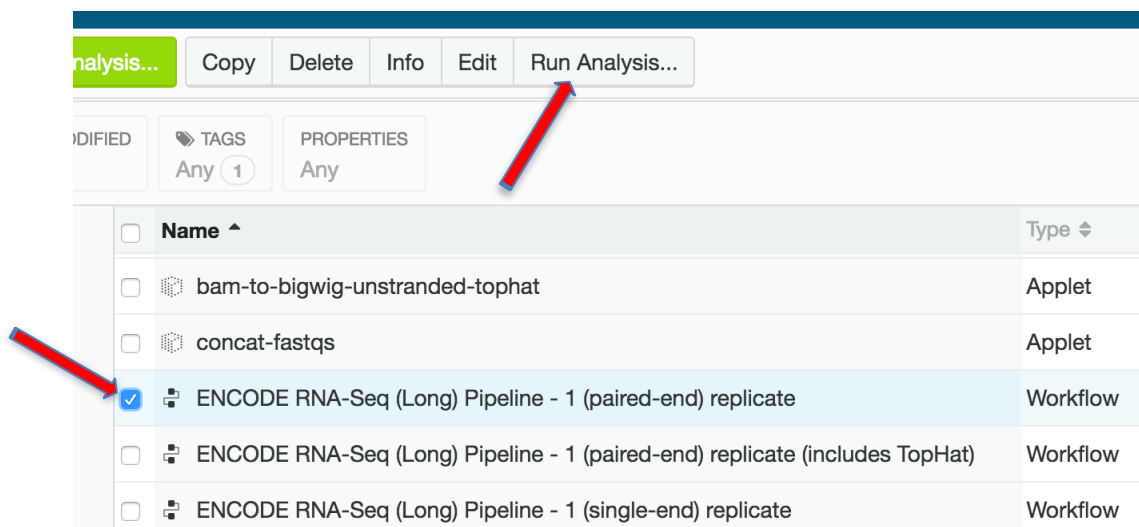7) Click the box next to "long-RNA-seq" and select "Add Data".



8) When finished, the following pop-up window should appear. Click "Close".
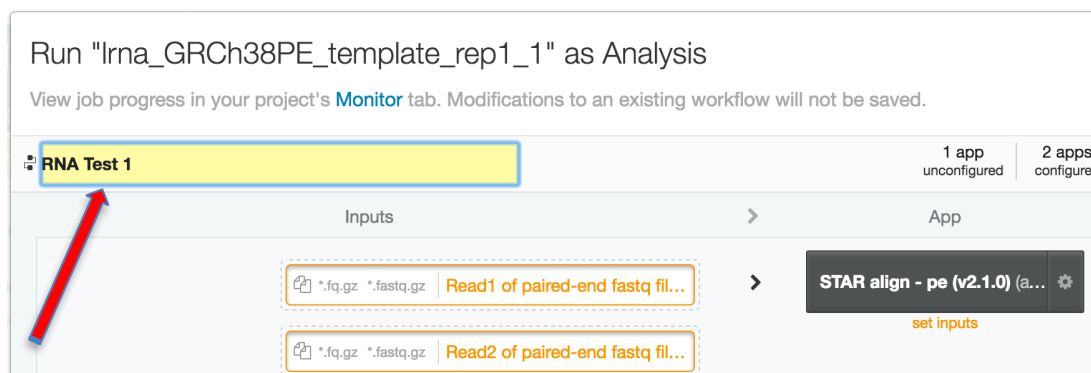


9) Click the long-RNA-seq text to open the folder. You should see the elements of the pipeline copied to your project.
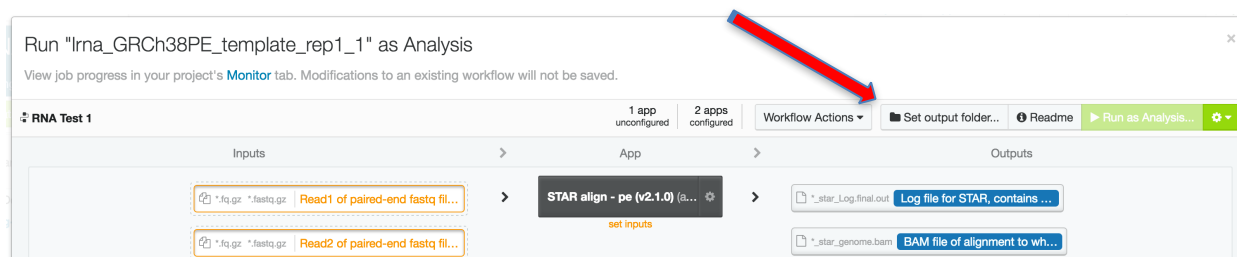
10) Select the workflow named "ENCODE RNA-seq (Long) Pipeline – 1 replicate (paired-end)". (You may need to resize the "Name" column and scroll to distinguish among the several versions of the pipeline.) Upon selecting the workflow, press "Run Analysis…".
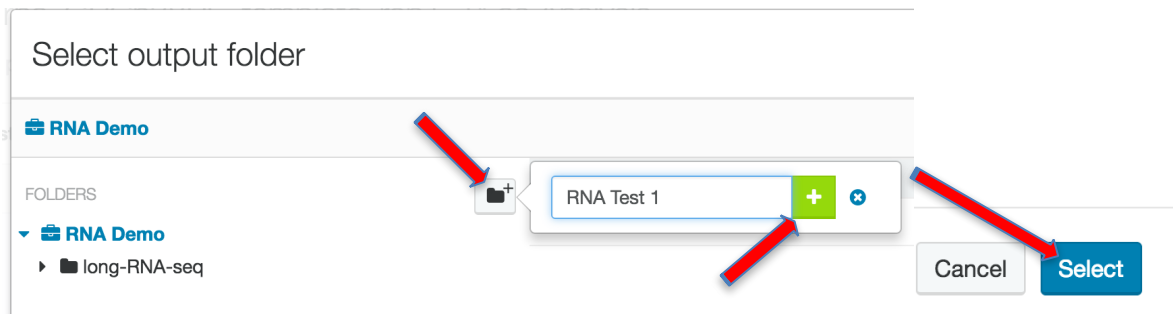


11) This window represents an "Analysis", which is an instance of the long-RNA-seq workflow. Give the analysis an informative name, like "Total RNA for chr21 of human fetal stomach tissue", or "RNA Test 1".
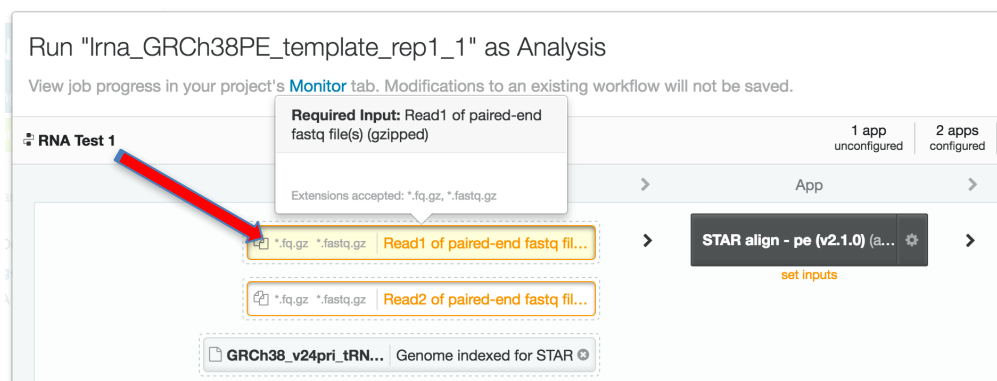


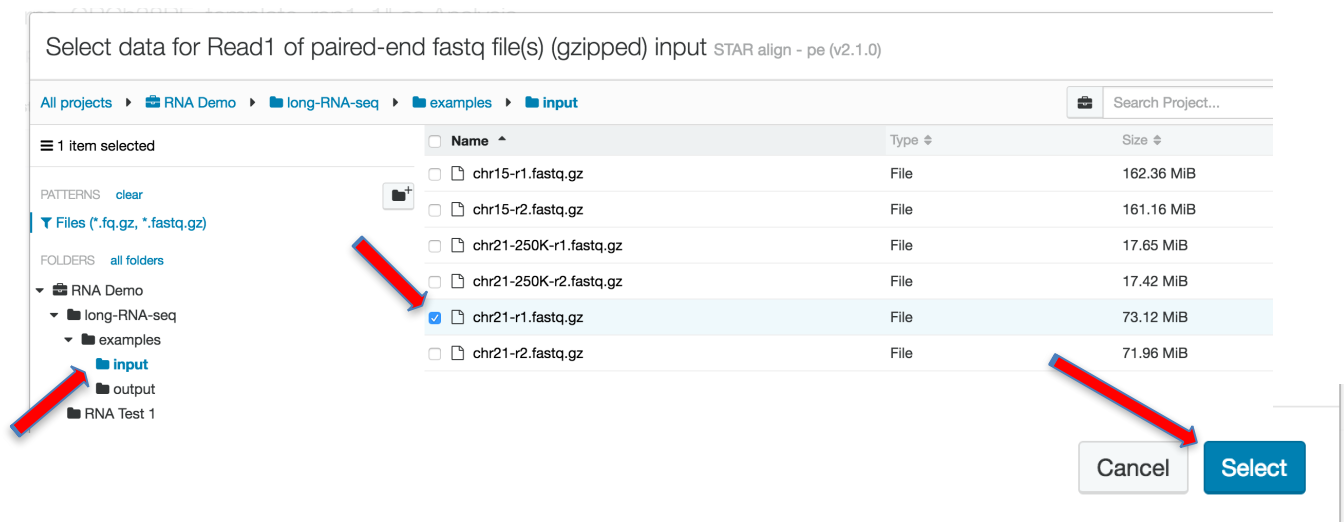12) Next, select a folder where your result files should be placed. Click on "Set output folder …"

13) Click on the new folder button to create a new folder and name it something catchy like "Chromosome 21 of ENCODE experiment ENCSR000AFI" or "RNA Test 1". Be sure to press the "+" sign and then "Select" at lower right.



14) Now it is time to add the sample data files to your analysis. Select the input box with orange text that says "Read1 of paired-end fast file(s) (gzipped)".



15) A new window opens where you will navigate to the input files. Expand the "long-RNA-seq" folder and within that "examples" and then "input". From here select the file named "chr21-r1.fastq.gz", then press "Select".

16) Now repeat for "Read2", choosing "chr21-r2.fastq.gz".



17) Notice that there is no longer any orange text declaring required inputs. Nevertheless, there are *optional* inputs that one may declare. Select the box labeled "STAR align - pe (v2.1.0). Here you may enter an "Identifier for biosample library", which will wind up being embedded in the result bam file's header. There are other options for running this alignment step on DNAnexus. It is not recommended that you alter these.



18) This analysis is already set to run on GRCh38. It may also be run on the hg19 assembly. To do this would require replacing three reference input files with the corresponding hg19 files, which are located in the same "ENCODE Uniform Processing Pipelines" project, from which you copied this pipeline.

**19)** With all the required inputs provided, it is time to run your analysis, so click "Run as Analysis".



**20)** Starting the analysis will bring up the "Monitor" tab which will display the details of the pipeline steps as they run. Click on the "+" box to see the analysis subjobs. If necessary, the "Terminate" button can be used to cancel the analysis.

21) Click on the analysis name to watch the progress of each step.



22) Within the output folder you specified above, all result files will accumulate as the steps of the pipeline complete. Many of the files will have additional information. For example, select the file whose name ends with "_star_genome.bam", then choose "Info".



23) This alignment file reports 6,591,972 reads but there is even more detail to be found by clicking on the "{…}".

24) By expanding on the "samtools_flagstats {…}", we learn that 6,591,972 is actually read pairs, and by expanding "STAR_final_log {…}", we can learn that 61.59% were uniquely mapped.

```
Details  {
    STAR_log_final : {...},
    samtools_flagstats : {
        diff_chroms : 0,
        diff_chroms_qc_failed : 0,
        duplicates : 0,
        duplicates_qc_failed : 0,
        mapped : 6587886,
        mapped_pct : "99.94%",
        mapped_qc_failed : 0,
        paired : 6591972,
        paired_properly : 6587886,
        paired_properly_pct : "99.94%"
        paired_properly_qc_failed : 0
        paired_qc_failed : 0,
        read1 : 3295986,
        read1_qc_failed : 0,
        read2 : 3295986,
        read2_qc_failed : 0,
        singletons : 0,
        singletons_pct : "0.00%",
        singletons_qc_failed : 0,
        total : 6591972,
        total_qc_failed : 0,
        with_itself : 6587886,
        with_itself_qc_failed : 0
    }
}
```

```
Details  {
    STAR_log_final : {
        % of reads mapped to multiple loci : "38.21%",
        % of reads mapped to too many loci : "0.17%",
        % of reads unmapped: other : "0.00%",
        % of reads unmapped: too many mismatches : "0.00%",
        % of reads unmapped: too short : "0.02%",
        Average input read length : 202,
        Average mapped length : 200.43,
        Deletion average length : 1.83,
        Deletion rate per base : "0.02%",
        Finished on : "Jun 06 20:35:55",
        Insertion average length : 1.38,
        Insertion rate per base : "0.01%",
        Mapping speed, Million of reads per hour : 53.66,
        Mismatch rate per base, % : "0.29%",
        Number of input reads : 1028564,
        Number of reads mapped to multiple loci : 393059,
        Number of reads mapped to too many loci : 1771,
        Number of splices: AT/AC : 145,
        Number of splices: Annotated (sjdb) : 158051,
        Number of splices: GC/AG : 1186,
        Number of splices: GT/AG : 158969,
        Number of splices: Non-canonical : 346,
        Number of splices: Total : 160646,
        Started job on : "Jun 06 20:30:44",
        Started mapping on : "Jun 06 20:34:46"
        Uniquely mapped reads % : "61.59%",
        Uniquely mapped reads number : 633462
    },
    samtools_flagstats : {...}
}
```

25) To visualize the signal results as custom tracks at the UCSC Genome Browser, select the the 2 bigwig files ending in "_minusUniq.bw" and "_plusUniq.bw". These two files are the signal produced from only uniquely mapped reads for the minus and plus DNA strands respectively. Select "Download."

26) A new window will pop up. Select "Get bulk URLs" and copy the two URLs. These URL's will link to your output files and will remain active for 24 hours.



27) In a new web browser window or tab, go to http://genome.ucsc.edu/ and select "My Data" from the top options bar, then select "Custom Tracks".



28) Paste the URLs you copied above into the first text window. Be sure the reference genome is correct for your results (human GRCh38/hg38 for this demo). *Tip: The UCSC Genome Browser is sensitive to white-space at the end of URL's. If there are spaces after the URL's you've pasted, delete them and make sure each URL is on its own line.* Now press "Submit".



29) You have one more chance to verify that you have selected GRCh38, before pressing "go".

30) Your two custom tracks will be displayed at the top of the browser image. Because the raw data were subsampled to only chromosome 21, there should be no significant result anywhere else. Set the browser's position to chr21:41,167,801-41,276,597. This is the location of the BACE2 gene. To see more clearly your results, change both of the custom tracks to "full" (right-click on the track in the image).



31) The UCSC Browser image should clearly show signal spiking at the exons of BACE2 for only one of your two custom tracks. Notice that these signal tracks are autoscaling, so while the plus signal peaks at 447 RPM (reads per million mapped reads), the minus signal peaks at less than 5 RPM in this location. Try other genes located on chromosome 21. For example, HLCS, SOD1, ETS2, or AIRE.

Congratulations! You have replicated an ENCODE analysis starting with primary data. You can repeat this process on your own data, and be assured that your results will be directly comparable to all the experiments the ENCODE DCC has analyzed.

## Other DNAnexus Tools:

*To load data once you are in your own project*

1) Start a "New Project" or find your own project in the DNAnexus homepage.



2) If new, name project in the upper left corner.



3) Select "Add Data" to select the files you want to use for analysis to your project.



4) When the "Add Data to Project" window pops up, select "From another DNAnexus project."



5) Scroll down and select "ENCODE Universal Processing Pipeline" project to access the data.

| Broad Inst Viral NGS | Viewer | 1 | 0.11 GB |
|---|---|---|---|
| ENCODE Uniform Processing Pipelines | Viewer | 13 | 349.28 GB |

6) Choose "Add Data" to select these files.

2 items selected (View)   Cancel   **Add Data**

7) When these files are uploaded, the following window will pop up.

**Add Data to Project:** ENCODE_Demo

Done

100%

| Name | | |
|---|---|---|
| long-RNA-seq | 25 items, 4 subfolders copied | ✓ DONE |
| Reference Files | 54 items, 6 subfolders copied | ✓ DONE |

Close

8) These files and associated applets will now appear in the Manage tab of your browser.

## ENCODE_Demo

Manage | Monitor | Visualize

⊕ Add Data | 📁 New Folder | New Workflow | ⚡ Start Analysis...

▼ 📦 ENCODE_Demo
  ▶ 📁 long-RNA-seq
  ▶ 📁 Reference Files

| | Name ▲ | Type ⇕ | Size ⇕ |
|---|---|---|---|
| ☐ | 📁 long-RNA-seq | Folder | |
| ☐ | 📁 Reference Files | Folder | |
| ☐ | 📦 align-star-se (Fri Dec 12 01:41:16 2014) | Applet | 1.16 MB |
| ☐ | 📦 align-tophat-pe (Fri Jan 9 01:28:56 2015) | Applet | 28.86 MB |
| ☐ | 📦 align-tophat-se (Fri Dec 12 01:41:04 2014) | Applet | 27.45 MB |

*To import a fastq file directly from the ENCODE portal to DNAnexus*

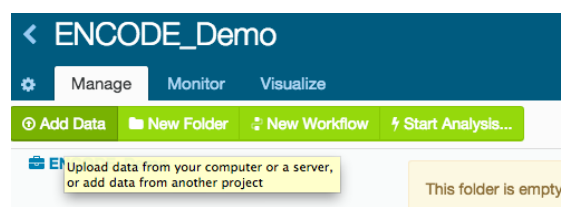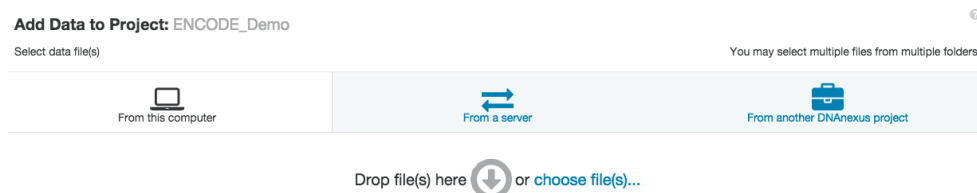1) Go to the ENCODE portal (encodeproject.org) and find the fastq file you are interested in using. Right click on this file and select "Copy Link Address."

**Files linked to ENCSR000AFI**

Raw data

| Accession ⇕ | File type ⇕ | Biological replicate ⇕ | Technical replicate ⇕ | Read length ⇕ | Run type ⇕ | Paired end ⇕ | Mapping assembly ⇕ | Lab ⇕ | Date added ⇕ | Validation status |
|---|---|---|---|---|---|---|---|---|---|---|
| ENCFF001RNE ⬇ Download 4.78 GB | fastq | 2 | 1 | 101 nt | paired-ended | 2 | | Thomas Gingeras, CSHL | 2013-07-17 | pending |
| ENCFF001F ⬇ Download 4.8 GB | | | | 101 nt | paired-ended | 1 | | Thomas Gingeras, CSHL | 2013-07-17 | pending |
| ENCFF001F ⬇ Download 5.15 GB | | | | 101 nt | paired-ended | 2 | | Thomas Gingeras, CSHL | 2013-07-18 | pending |
| ENCFF001F ⬇ Download | | | | 101 nt | paired-ended | 1 | | Thomas Gingeras, CSHL | 2013-07-18 | pending |

Open Link in New Tab
Open Link in New Window
Open Link in Incognito Window
Save Link As...
**Copy Link Address**
Copy
Search Google for 'Download'
Print...

2) In the manage tab, under "Add Data" select the "From a Server" option and paste the URL into the box. Select "Add Data" and the file will upload.

**Add Data to Project:** ENCODE_Demo

Select data file(s)                                                    You may add multiple URLs.

| From this computer | From a server | From another DNAnexus project |
|---|---|---|

https://www.encodeproject.org/files/ENCFF001RNE/@@download/ENCFF001RNE.fastq.gz    ✕

Enter a URL...

**Add Data to Project:** ENCODE DEMO_June24

**Done**

**Name**
https://www.encodeproject.org/files/ENCFF001RNE/@@download/ENCFF001RNE.fastq.gz    ✕

*To share project with another user*

1) In order to share your project, select the blue "Share" button at the upper right corner of the browser page.

**Admin** your access      🔒 **Private** access policy      **Share** 2 Members

2) This will bring up a pop-up window where you can add user names and select permissions to allow collaborators access to view, edit, or contribute to your projects.

Share project                                                    ✕

| Name | Access | Charges Allowed | |
|---|---|---|---|
| Benjamin Hitz (hitz) | Viewer | | Remove |
| Eurie Hong (euriehong) | Admin | $ | |

Add member...

**Examples:**
jsmith
user-jsmith

Close