



EDGE bioinformatics

# EDGE Documentation

*Release Notes 1.1*

**EDGE Development Team**

**Feb 26, 2019**

---

## Contents

---

<b>1</b>	<b>EDGE ABCs</b>	<b>1</b>
1.1	About EDGE Bioinformatics . . . . .	1
1.2	Bioinformatics overview . . . . .	1
1.3	Computational Environment . . . . .	3
<b>2</b>	<b>Introduction</b>	<b>4</b>
2.1	What is EDGE? . . . . .	4
2.2	Why create EDGE? . . . . .	4
<b>3</b>	<b>System requirements</b>	<b>6</b>
3.1	Ubuntu 14.04 . . . . .	6
3.2	CentOS 6.7 . . . . .	7
3.3	CentOS 7 . . . . .	8
<b>4</b>	<b>Installation</b>	<b>10</b>
4.1	EDGE Installation . . . . .	10
4.2	EDGE Docker image . . . . .	18
4.3	EDGE VMware/OVF Image . . . . .	18
<b>5</b>	<b>Graphic User Interface (GUI)</b>	<b>20</b>
5.1	User Login . . . . .	20
5.2	Upload Files . . . . .	21
5.3	Initiating an analysis job . . . . .	22
5.4	Choosing processes/analyses . . . . .	24
5.5	Submission of a job . . . . .	31
5.6	Checking the status of an analysis job . . . . .	31
5.7	Monitoring the Resource Usage . . . . .	33
5.8	Management of Jobs . . . . .	33
5.9	Other Methods of Accessing EDGE . . . . .	34
<b>6</b>	<b>Command Line Interface (CLI)</b>	<b>37</b>
6.1	Configuration File . . . . .	38
6.2	Test Run . . . . .	40
6.3	Descriptions of each module . . . . .	42
6.4	Other command-line utility scripts . . . . .	49
<b>7</b>	<b>Output</b>	<b>50</b>

7.1	Example Output	51
<b>8</b>	<b>Databases</b>	<b>52</b>
8.1	EDGE provided databases	52
8.2	Building bwa index	54
8.3	SNP database genomes	54
8.4	Ebola Reference Genomes	61
<b>9</b>	<b>Third Party Tools</b>	<b>62</b>
9.1	Assembly	62
9.2	Annotation	62
9.3	Alignment	64
9.4	Taxonomy Classification	65
9.5	Phylogeny	66
9.6	Visualization and Graphic User Interface	66
9.7	Utility	67
<b>10</b>	<b>FAQs and Troubleshooting</b>	<b>69</b>
10.1	FAQs	69
10.2	Troubleshooting	70
10.3	Discussions / Bugs Reporting	70
<b>11</b>	<b>Copyright</b>	<b>71</b>
<b>12</b>	<b>Contact Us</b>	<b>72</b>
<b>13</b>	<b>Citation</b>	<b>73</b>

# CHAPTER 1

---

## EDGE ABCs

---

A quick About EDGE, overview of the Bioinformatic workflows, and the Computational environment

### 1.1 About EDGE Bioinformatics

EDGE bioinformatics was **developed to help biologists process Next Generation Sequencing data** (in the form of **raw FASTQ files**), even if they have little to no bioinformatics expertise. EDGE is a **highly integrated and interactive web-based platform** that is capable of running many of the standard analyses that biologists require for viral, bacterial/archaeal, and metagenomic samples. EDGE provides the following analytical workflows: **pre-processing, assembly and annotation, reference-based analysis, taxonomy classification, phylogenetic analysis, and PCR analysis**. EDGE provides an intuitive web-based interface for user input, allows users to visualize and interact with selected results (e.g. JBrowse genome browser), and generates a final detailed PDF report. Results in the form of tables, text files, graphic files, and PDFs can be downloaded. A user management system allows tracking of an individual's EDGE runs, along with the ability to share, post publicly, delete, or archive their results.

While EDGE was intentionally designed to be as simple as possible for the user, there is still no single ‘tool’ or algorithm that fits all use-cases in the bioinformatics field. Our intent is to provide a detailed panoramic view of your sample from various analytical standpoints, but users are encouraged to have some knowledge of how each tool/algorithm workflow functions, and some insight into how the results should best be interpreted.

### 1.2 Bioinformatics overview

#### 1.2.1 Inputs:

The input to the EDGE workflows begins with one or more **illumina FASTQ files** for a single sample. (There is currently limited capability of incorporating PacBio and Oxford Nanopore data into the Assembly module.) The user can also enter SRA/ENA accessions to allow processing of publically available datasets. Comparison among samples is not yet supported but development is underway to accommodate such a function for assembly and taxonomy profile comparisons.

## 1.2.2 Workflows:

### Pre-Processing

Assessment of quality control is performed by [FAQCS](#). The host removal step requires the input of one or more reference genomes as FASTA. Several common references are available for selection. Trimmed and host-screened FASTQ files are used for input to the other workflows.

### Assembly and Annotation

We provide the [IDBA](#), [Spades](#), and [MegaHit](#) (in the development version) assembly tools to accommodate a range of sample types and data sizes. When the user selects to perform an assembly, all subsequent workflows can execute analysis with either the reads, the contigs, or both (default).

### Reference-Based Analysis

For comparative reference-based analysis with reads and/or contigs, users must input one or more references (as FASTA or multi-FASTA if there are more than one replicon) and/or select from a drop-down list of RefSeq complete genomes. Results include lists of missing regions (gaps), inserted regions (with input contigs if assembly was performed), SNPs (and coding sequence changes), as well as genome coverage plots and interactive access via JBrowse.

### Taxonomy Classification

For taxonomy classification with reads, multiple tools are used and the results are summarized in heat map and radar plots. Individual tool results are also presented with taxonomy dendograms and Krona plots. Contig classification occurs by assigning taxonomies to all possible portions of contigs. For each contig, the longest and best match (using BWA-MEM) is kept for any region within the contig and the region covered is assigned to the taxonomy of the hit. The next best match to a region of the contig not covered by prior hits is then assigned to that taxonomy. The contig results can be viewed by length of assembly coverage per taxa or by number of contigs per taxa.

### Phylogenetic Analysis

For phylogenetic analysis, the user must select datasets from near neighbor isolates for which the user desires a phylogeny. A minimum of three additional datasets are required to draw a tree. At least one dataset must be an assembly or complete genome. [RefSeq genomes \(Bacteria, Archaea, Viruses\)](#) are available from a dropdown menu, SRA and FASTA entries are allowed, and previously built databases for some select groups of bacteria are provided. This workflow (see [PhaME](#)) is a whole genome SNP-based analysis that uses one reference assembly to which both reads and contigs are mapped. Because this analysis is based on read alignments and/or contig alignments to the reference genome(s), we **strongly recommend only selecting genomes that can be adequately aligned at the nucleotide level (i.e. ~90% identity or better)**. The number of ‘core’ nucleotides able to be aligned among all genomes, and the number of SNPs within the core, are what determine the resolution of the phylogenetic tree. Output phylogenies are presented along with text files outlining the SNPs discovered.

### Primer Analysis

For primer analysis, if the user would like to validate known PCR primers in silico, a FASTA file of primer sequences must be input. New primers can be generated from an assembly as well.

**All commands and tool parameters are recorded in log files to make sure the results are repeatable and traceable.** The main output is an integrated interactive web page that includes summaries of all the workflows run and features tables, graphical plots, and links to genome (if assembled, or of a selected reference) browsers and to access unprocessed results and log files. Most of these summaries, including plots and tables are included within a final PDF report.

## 1.2.3 Limitations

### Pre-processing

For host removal/screening, not all genomes are available from a drop-down list, however

## Assembly and Taxonomy Classification

EDGE has been primarily designed to **analyze microbial (bacterial, archaeal, viral) isolates or (shotgun) metagenome samples**. Due to the complexity and computational resources required for eukaryotic genome assembly, and the fact that the current taxonomy classification tools do not support eukaryotic classification, EDGE does not fully support eukaryotic samples. The combination of large NGS data files and complex metagenomes may also run into computational memory constraints.

### Reference-based analysis

We recommend only aligning against (a limited number of) most closely related genome(s). If this is unknown, the Taxonomy Classification module is recommended as an alternative. If the user selects too many references, this may affect runtimes or require more computational resources than may be available on the user's system.

### Phylogenetic Analysis

Because this pipeline provides SNP-based trees derived from whole genome (and contig) alignments or read mapping, **we recommend selecting genomes within the same species or at least within the same genus**.

## 1.3 Computational Environment

### 1.3.1 EDGE source code, images, and webservers

EDGE was designed to be installed and implemented from within any institute that provides sequencing services or that produces or hosts NGS data. When installed locally, EDGE can access the raw FASTQ files from within the institute, thereby providing immediate access by the biologist for analysis. EDGE is available in a variety of packages to fit various institute needs. **EDGE source code** can be obtained via our [GitHub](#) page. To simplify installation, a [VM in OVF](#) or a [Docker image](#) can also be obtained. A **demonstration version of EDGE** is currently available at <https://bioedge.lanl.gov> with example data sets available to the public to view and/or re-run. This webserver has 24 cores, 512GB ram with Ubuntu 14.04.3 LTS, and also allows EDGE runs of SRA/ENA data. This webserver does not currently support upload of data (due in part to LANL security regulations), however local installations are meant to be fully functional.

# CHAPTER 2

---

## Introduction

---

### 2.1 What is EDGE?

EDGE is a highly adaptable bioinformatics platform that allows laboratories to quickly analyze and interpret genomic sequence data. The bioinformatics platform allows users to address a wide range of use cases including assay validation and the characterization of novel biological threats, clinical samples, and complex environmental samples. EDGE is designed to:

- Align to real world use cases
- Make use of open source (free) software tools
- Run analyses on small, relatively inexpensive hardware
- Provide remote assistance from bioinformatics specialists

### 2.2 Why create EDGE?

EDGE bioinformatics was **developed to help biologists process Next Generation Sequencing data** (in the form of **raw FASTQ files**), even if they have little to no bioinformatics expertise. EDGE is a **highly integrated and interactive web-based platform** that is capable of running many of the standard analyses that biologists require for viral, bacterial/archaeal, and metagenomic samples. EDGE provides the following analytical workflows: **quality trimming and host removal, assembly and annotation, comparisons against known references, taxonomy classification of reads and contigs, whole genome SNP-based phylogenetic analysis, and PCR analysis**. EDGE provides an intuitive web-based interface for user input, allows users to visualize and interact with selected results (e.g. JBrowse genome browser), and generates a final detailed PDF report. Results in the form of tables, text files, graphic files, and PDFs can be downloaded. A user management system allows tracking of an individual's EDGE runs, along with the ability to share, post publicly, delete, or archive their results.

While the design of EDGE was intentionally done to be as simple as possible for the user, there is still no single 'tool' or algorithm that fits all use-cases in the bioinformatics field. Our intent is to provide a detailed panoramic view of your sample from various analytical standpoints, but users are encouraged to have some insight into how each tool or workflow functions, and how the results should best be interpreted.

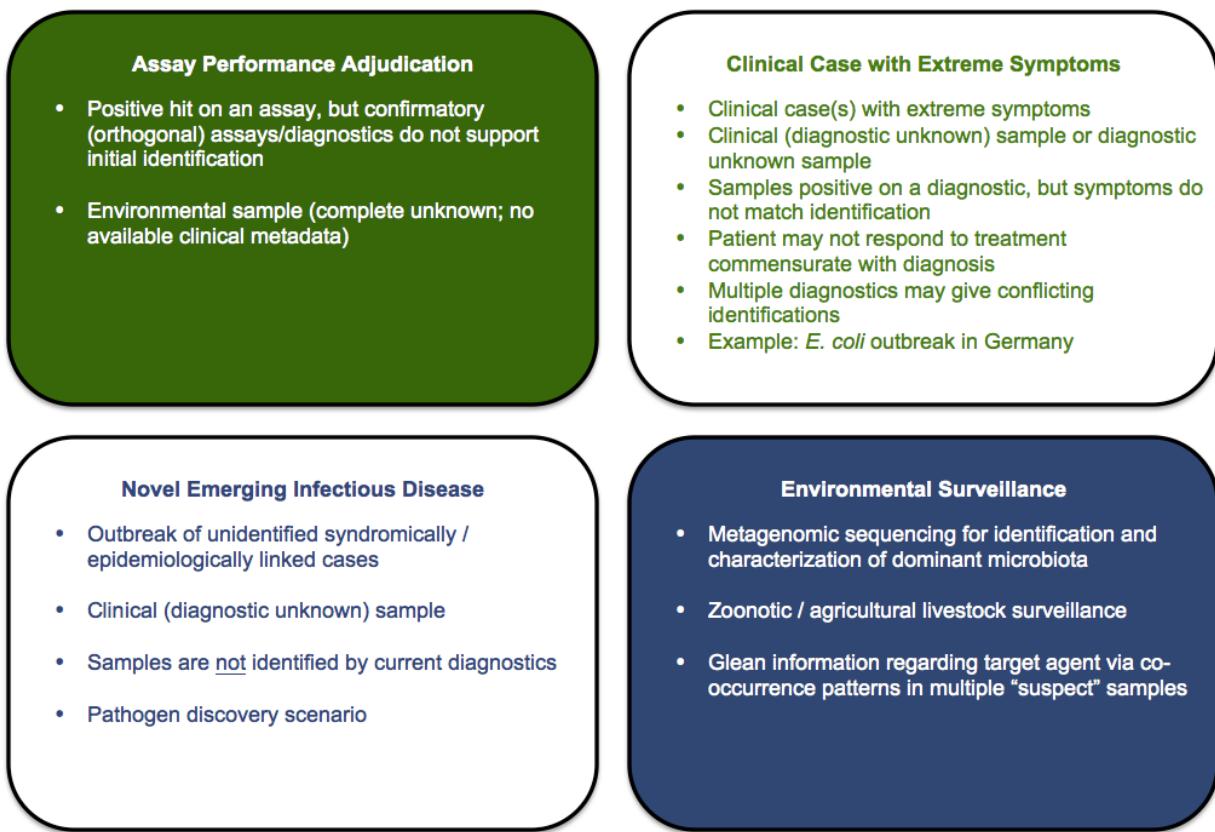


Fig. 1: Four common Use Cases guided initial EDGE Bioinformatic Software development.

# CHAPTER 3

---

## System requirements

---

NOTE: The web-based online version of EDGE, found on [https://bioedge.lanl.gov/edge\\_ui/](https://bioedge.lanl.gov/edge_ui/) is run on our own internal servers and is our recommended mode of usage for EDGE. It does not require any particular hardware or software other than a web browser. This segment and the installation segment only apply if you want to run EDGE through Python or Apache 2, or through the CLI.

The current version of the EDGE pipeline has been extensively tested on a Linux Server with Ubuntu 14.04 and Centos 6.5 and 7.0 operating system and will work on 64bit Linux environments. Perl v5.8 or above is required. Python 2.7 is required. Due to the involvement of several memory/time consuming steps, it requires at least 16Gb memory and at least 8 computing CPUs. A higher computer spec is recommended: 128Gb memory and 16 computing CPUs.

Please ensure that your system has the essential software building packages installed properly before running the installing script.

The following are required installed by system administrator.

---

**Note:** If your system OS is neither Ubuntu 14.04 or Centos 6.5 or 7.0, it may have differnt packages/libraries name and the newer complier (gcc5) on newer OS (ex: Ubuntu 16.04) may fail on coupling some of thirdparty bioinformatics tools. We would suggest to use EDGE [VMware image](#) or [Docker container](#).

---

### 3.1 Ubuntu 14.04

1. Install build essential libraries and dependancies:

```
sudo apt-get install build-essential  
sudo apt-get install libreadline-gplv2-dev  
sudo apt-get install libx11-dev  
sudo apt-get install libxt-dev libgs10-dev  
sudo apt-get install libncurses5-dev  
sudo apt-get install gfortran  
sudo apt-get install inkscape  
sudo apt-get install libwww-perl libxml-libxml-perl libperlio-gzip-perl
```

(continues on next page)

(continued from previous page)

```
sudo apt-get install zlib1g-dev zip unzip libjson-perl
sudo apt-get install libpng-dev
sudo apt-get install cpanminus
sudo apt-get install default-jre
sudo apt-get install firefox
sudo apt-get install wget curl csh
```

2. Install python packages for Metaphlan (Taxonomy assignment software):

```
sudo apt-get install python-numpy python-matplotlib python-scipy libpython2.7-
˓→stdlib
sudo apt-get install python-pip python-pandas python-sympy python-nose
```

3. Install BioPerl:

```
sudo apt-get install bioperl
˓→or
sudo cpan -i -f CJFIELDS/BioPerl-1.6.923.tar.gz
```

4. Install packages for user management system:

```
sudo apt-get install sendmail mysql-client mysql-server phpMyAdmin tomcat7
```

## 3.2 CentOS 6.7

1. Install dependancies using yum:

```
# add epel repository
sudo yum -y install epel-release
su -c 'yum localinstall -y --nogpgcheck http://download1.rpmfusion.org/free/el/
˓→updates/6/i386/rpmfusion-free-release-6-1.noarch.rpm http://download1.rpmfusion.
˓→org/nonfree/el/updates/6/i386/rpmfusion-nonfree-release-6-1.noarch.rpm'
sudo yum -y update

sudo yum -y install\
csh gcc gcc-c++ make curl binutils gd gsl-devel\
libX11-devel readline-devel libXt-devel ncurses-devel inkscape\
freetype freetype-devel zlib zlib-devel git\
blas-devel atlas-devel lapack-devel libpng libpng-devel\
expat expat-devel graphviz java-1.7.0-openjdk\
perl-Archive-Zip perl-Archive-Tar perl-CGI perl-CGI-Session \
perl-DBI perl-GD perl-JSON perl-Module-Build perl-CPAN-Meta-YAML\
perl-XML-LibXML perl-XML-Parser perl-XML-SAX perl-XML-SAX-Writer\
perl-XML-Simple perl-XML-Twig perl-XML-Writer perl-YAML\
perl-Test-Most perl-PerlIO-gzip perl-SOAP-Lite perl-GraphViz
```

2. Install perl cpanm:

```
curl -L http://cpanmin.us | perl - App::cpanminus
```

3. Install perl modules by cpanm:

```
cpanm Graph Time::Piece Data::Dumper IO::Compress::Gzip Data::Stag IO::String
cpanm Algorithm::Munkres Array::Compare Clone Convert::Binary::C_
˓→XML::Parser::PerlSAX
```

(continues on next page)

(continued from previous page)

```
cpanm HTML::Template HTML::TableExtract List::MoreUtils PostScript::TextBlock
cpanm SVG SVG::Graph Set::Scalar Sort::Naturally Spreadsheet::ParseExcel
cpanm -f Bio::Perl
```

#### 4. Install dependent packages for Python:

EDGE requires several packages (NumPy, Matplotlib, SciPy, IPython, Pandas, SymPy and Nose) to work properly. These packages are available at PyPI (<https://pypi.python.org/pypi>) for downloading and installing respectively. Or you can install a Python distribution with dependent packages instead. We suggest users to use Anaconda Python distribution. You can download the installers and find more information at their website (<https://store.continuum.io/cshop/anaconda/>). The installation is interactive. Type in /opt/apps/anaconda when the script asks for the location to install python.:

```
bash Anaconda-2.x.x-Linux-x86.sh
ln -s /opt/apps/anaconda/bin/python /path/to/edge_v1.x/bin/
```

Create symlink anaconda python to edge/bin. So system will use your python over the system's.

#### 5. Install packages for user management system:

```
sudo yum -y install sendmail mysql mysql-server phpmyadmin tomcat
```

## 3.3 CentOS 7



#### 1. Install libraries and dependencies by yum:

```
# add epel repository
sudo yum -y install epel-release

sudo yum install -y libX11-devel readline-devel libXt-devel ncurses-devel \
    inkscape \
    scipy expat expat-devel freetype freetype-devel zlib zlib-devel perl-App-
    cpanminus \
    perl-Test-Most python-pip blas-devel atlas-devel lapack-devel numpy numpy-
    f2py \
    libpng12 libpng12-devel perl-XML-Simple perl-JSON csh gcc gcc-c++ make \
    binutils \
    gd gsl-devel git graphviz java-1.7.0-openjdk perl-Archive-Zip perl-CGI \
    perl-CGI-Session perl-CPAN-Meta-YAML perl-DBI perl-Data-Dumper perl-GD perl-
    IO-Compress \
    perl-Module-Build perl-XML-LibXML perl-XML-Parser perl-XML-SAX perl-XML-SAX-
    Writer \
    perl-XML-Twig perl-XML-Writer perl-YAML perl-PerlIO-gzip python-matplotlib \
    python-six
```

#### 2. Update existing python and perl tools:

```
sudo pip install --upgrade six scipy matplotlib
sudo cpanm App::cpanoutdated
sudo su -
```

(continues on next page)

(continued from previous page)

```
cpan-outdated -p | cpanm  
exit
```

3. Install perl modules by cpanm:

```
cpanm Graph Time::Piece Bio::Perl  
cpanm Algorithm::Munkres Archive::Tar Array::Compare Clone Convert::Binary::C  
cpanm HTML::Template HTML::TableExtract List::MoreUtils PostScript::TextBlock  
cpanm SOAP::Lite SVG SVG::Graph Set::Scalar Sort::Naturally  
    ↳ Spreadsheet::ParseExcel  
cpanm CGI CGI::Simple GD Graph GraphViz XML::Parser::PerlSAX XML::SAX  
    ↳ XML::SAX::Writer XML::Simple XML::Twig XML::Writer
```

4. Install packages for user management system:

```
sudo yum -y install sendmail mariadb-server mariadb phpMyAdmin tomcat
```

5. Configure firewall for ssh, http, https, and smtp:

```
sudo firewall-cmd --permanent --add-service=ssh  
sudo firewall-cmd --permanent --add-service=http  
sudo firewall-cmd --permanent --add-service=https  
sudo firewall-cmd --permanent --add-service=smtp
```

---

**Note:** You may need to turn the SELinux into Permissive mode.

```
sudo setenforce 0
```

---

# CHAPTER 4

---

## Installation

---

### 4.1 EDGE Installation

---

**Note:** A base install is ~8GB for the code base and ~177GB for the databases.

---

1. Please ensure that your system has the *essential software building packages* (page 6). installed properly before proceeding following installation.
2. Download the codebase, databases and third party tools.

```
## Codebase is ~68Mb and contains all the scripts and HTML needed to make EDGE run
wget -c https://edge-dl.lanl.gov/EDGE/1.1/edge_main_v1.1.tgz

## Third party tools is ~1.9Gb and contains the underlying programs needed to do the analysis
wget -c https://edge-dl.lanl.gov/EDGE/1.1/edge_v1.1_thirdParty_softwares.tgz

## Pipeline database is ~7.9Gb and contains the other databases needed for EDGE
wget -c https://edge-dl.lanl.gov/EDGE/1.1/edge_pipeline_v1.1.databases.tgz

## GOTTCCHA database is ~14Gb and contains the custom databases for the GOTTCCHA taxonomic identification pipeline
wget -c https://edge-dl.lanl.gov/EDGE/1.1/GOTTCCHA_db_for_edge_v1.1.tgz

## BWA index is ~41Gb and contains the databases for bwa taxonomic identification pipeline
wget -c https://edge-dl.lanl.gov/EDGE/1.1/bwa_index1.1.tgz

## NCBI Genomes is ~8Gb and contain the full genomes for prokaryotes and some viruses
wget -c https://edge-dl.lanl.gov/EDGE/1.1/NCBI_genomes_for_edge_v1.1.tar.gz
```

**Warning:** Be patient; the database files are huge.

3. Unpack main archive:

```
tar -xvzf edge_main_v1.1.1.tgz
```

**Note:** The main directory, edge\_v1.1.1, will be created.

4. Move the database and third party archives into main directory (edge\_v1.1.1):

```
mv edge_v1.1_thirdParty_softwares.tgz edge_v1.1.1/
mv edge_pipeline_v1.1.databases.tgz edge_v1.1.1/
mv GOTTCHA_db_for_edge_v1.1.tgz edge_v1.1.1/
mv bwa_index1.1.tgz edge_v1.1.1/
mv NCBI_genomes_for_edge_v1.1.tar.gz edge_v1.1.1/
```

5. Change directory to main directory and unpack databases and third party tools archive:

```
cd edge_v1.1.1

# unpack third party tools
tar -xvzf edge_v1.1_thirdParty_softwares.tgz

# unpack databases
tar -xvzf edge_pipeline_v1.1.databases.tgz
tar -xvzf GOTTCHA_db_for_edge_v1.1.tgz
tar -xzvf bwa_index1.1.tgz
tar -xvzf NCBI_genomes_for_edge_v1.1.tar.gz
```

**Note:** To this point, you should see a database directory and a thirdParty directory in the main directory

6. Installing pipeline:

```
./INSTALL.sh
```

It will install the following depended *tools* (page 62).

- Assembly
  - idba
  - spades
- Annotation
  - prokka
  - RATT
  - tRNAscan
  - barrnap
  - BLAST+
  - blastall
  - phageFinder

- glimmer
- aragorn
- prodigal
- tbl2asn
- Alignment
  - hmmer
  - infernal
  - bowtie2
  - bwa
  - mummer
- Taxonomy
  - kraken
  - metaphlan
  - kronatools
  - gottcha
- Phylogeny
  - FastTree
  - RAxML
- Utility
  - bedtools
  - R
  - GNU\_parallel
  - tabix
  - JBrowse
  - primer3
  - samtools
  - sratoolkit
- Perl\_Modules
  - perl\_parallel\_forkmanager
  - perl\_excel\_writer
  - perl\_archive\_zip
  - perl\_string\_approx
  - perl\_pdf\_api2
  - perl\_html\_template
  - perl\_html\_parser
  - perl\_JSON

- perl\_bio\_phylo
  - perl\_xml\_twig
  - perl\_cgi\_session
7. Restart the Terminal Session to allow \$EDGE\_HOME to be exported.

**Note:** After running INSTALL.sh successfully, the binaries and related scripts will be stored in the ./bin and ./scripts directory. It also writes EDGE\_HOME environment variable into .bashrc or .bash\_profile.

#### 4.1.1 Testing the EDGE Installation

After installing the packages above, it is highly recommended to test the installation:

```
> cd $EDGE_HOME/testData
> ./runAllTest.sh
```

```
Working Dir: /data/edge_v1.0/testData
EDGE HOME Dir: /data/edge_v1.0
[00:00:00] [ 6 %]      Test Assembly ..... [OK]
[00:03:02] [13 %]      Test Contigs2Refernece ..... [OK]
[00:08:24] [20 %]      Test ContigsAnnotation ..... [OK]
[00:10:57] [26 %]      Test ContigsTaxonomy ..... [OK]
[00:13:41] [33 %]      Test HostRemoval ..... [OK]
[00:14:21] [40 %]      Test PCRassay ..... [OK]
[00:15:03] [46 %]      Test PhageFinder ..... [OK]
[00:15:16] [53 %]      Test PhylogeneticAnalysis ..... [OK]
[00:15:23] [60 %]      Test QC ..... [OK]
[00:17:40] [66 %]      Test Reads2Contigs ..... [OK]
[00:18:48] [73 %]      Test Reads2Reference ..... [OK]
[00:23:23] [80 %]      Test ReadsTaxonomy ..... [OK]
[00:36:18] [86 %]      Test Report ..... [OK]
[00:36:41] [93 %]      Test SRADownload ..... [OK]
[00:39:08] [100 %]     Test VariantAnalysis ..... [OK]

15/15 test(s) passed

Total Running Time: 00:43:55
```

There are 15 module/unit tests which took around 44 mins in our testing environments. (24 cores 2.60GHz, 512GB ram with Ubuntu 14.04.3 LTS). You will see test output on the terminal indicating test successes and failures. Some tests may fail due to missing external applications/modules/packages or failed installation. These will be noted separately in the \$EDGE\_HOME/testData/runXXXXTest/TestOutput/error.log or log files in each modules. If these are related to features of EDGE that you are not using, this is acceptable. Otherwise, you'll want to ensure that you have the EDGE installed correctly. If the output doesn't indicate any failures, you are now ready to use EDGE through command line. To take advantage of the user friendly GUI, please follow the section below to configure the EDGE Web server.

## 4.1.2 Apache Web Server Configuration

### 1. Install apache2

```
For Ubuntu
> sudo apt-get install apache2

For CentOS
> sudo yum -y install httpd
```

### 2. Enable apache cgid, proxy, headers modules:

```
For Ubuntu
> sudo a2enmod cgid proxy proxy_http headers
```

### 3. Modify/Check sample apache configuration file:

```
Double check $EDGE_HOME/edge_ui/apache_conf/edge_apache.conf alias directories to ↵
↪match EDGE
installation path at line 2,3,13,14,26,51.
The default is configured as http://localhost/edge_ui/ or http://www.yourdomain. ↵
↪com/edge_ui/
```

### 4. (Optional) If users are behind a corporate proxy for internet:

```
Please add proxy info into $EDGE_HOME/edge_ui/apache_conf/edge_apache.conf or ↵
↪$EDGE_HOME/edge_ui/apache_conf/edge_httpd.conf

# Add following proxy env
SetEnv http_proxy http://yourproxy:port
SetEnv https_proxy http://yourproxy:port
SetEnv ftp_proxy http://yourproxy:port
```

### 5. Copy modified edge\_apache.conf to the apache or Insert content into httpd.conf

```
For Ubuntu
> cp $EDGE_HOME/edge_ui/apache_conf/edge_apache.conf /etc/apache2/conf-available/
> ln -s /etc/apache2/conf-available/edge_apache.conf /etc/apache2/conf-enabled/
For CentOS
> cp $EDGE_HOME/edge_ui/apache_conf/edge_apache.conf /etc/httpd/conf.d/
```

### 6. Modify permissions: modify permissions on installed directory to match apache user

```
For Ubuntu 14, the user can be edited at /etc/apache2/envvars and the variable ↵
↪are APACHE_RUN_USER and APACHE_RUN_GROUP.

For CentOS, the user can be edited at /etc/httpd/conf/httpd.conf and the variable ↵
↪are User and Group.

> chown -R xxxxx $EDGE_HOME/edge_ui $EDGE_HOME/edge_ui/JBrowse/data # (xxxxx is ↵
↪the APACHE_RUN_USER value)
```

(continues on next page)

(continued from previous page)

```
> chgrp -R xxxxx $EDGE_HOME/edge_ui $EDGE_HOME/edge_ui/JBrowse/data # (xxxxx is
  ↳the APACHE_RUN_GROUP value)
```

## 7. Restart the apache2 to activate the new configuration

```
For Ubuntu
>sudo service apache2 restart

For CentOS
>sudo httpd -k restart
```

### 4.1.3 User Management system installation

#### 1. Create database: userManagement:

```
> cd $EDGE_HOME/userManagement
> mysql -p -u root
mysql> create database userManagement;
mysql> use userManagement;
```

**Note:** make sure mysql is running. If not, run “sudo service mysqld start”; ;

for CentOS7: “sudo systemctl start mariadb.service && sudo systemctl enable mariadb.service”

#### 2. Load userManagement\_schema.sql:

```
mysql> source userManagement_schema.sql;
```

#### 3. Load userManagement\_constraints.sql:

```
mysql> source userManagement_constraints.sql;
```

#### 4. Create an user account

```
username: yourDBUsername
password: yourDBPassword
(also modify the username/password in userManagementWS.xml file)
and grant all privileges on database "userManagement" to user yourDBUsername

mysql> CREATE USER 'yourDBUsername'@'localhost' IDENTIFIED BY 'yourDBPassword';

mysql> GRANT ALL PRIVILEGES ON userManagement.* to 'yourDBUsername'@'localhost';

mysql>exit;
```

#### 5. Configure tomcat:

```
* Copy mysql-connector-java-5.1.34-bin.jar to /usr/share/tomcat/lib/
```

For Ubuntu and CentOS6

(continues on next page)

(continued from previous page)

```

> cp mysql-connector-java-5.1.34-bin.jar /usr/share/tomcat7/lib/
For CentOS7
> cp mariadb-java-client-1.2.0.jar /usr/share/tomcat/lib/

* Configure tomcat basic auth to secure /user/admin/register web service
  add lines below to /var/lib/tomcat7/conf/tomcat-users.xml of Ubuntu or /etc/
  ↵tomcat/tomcat-users.xml of CentOS

<role rolename="admin"/>
<user username="yourAdminName" password="yourAdminPassword" roles="admin"/>

(also modify the username and password in createAdminAccount.pl file)

* Inactive timeout in /var/lib/tomcat7/conf/web.xml or /etc/tomcat/web.xml_
  ↵(default is 30mins)

<!-- <session-config>
      <session-timeout>30</session-timeout>
    </session-config> -->

* add the line below to tomcat /usr/share/tomcat7/bin/catalina.sh of Ubuntu or /
  ↵etc/tomcat/tomcat.conf of CentOS to increase PermSize:

JAVA_OPTS="-Xms256M -Xmx1024M -XX:PermSize=256m -XX:MaxPermSize=512m"

* Restart tomcat server

for Ubuntu
> sudo service tomcat7 restart
for CentOS6
> sudo service tomcat restart
for CentOS7
> sudo systemctl restart tomcat.service

* Deploy userManagementWS to tomcat server

for Ubuntu
> cp userManagementWS.war /var/lib/tomcat7/webapps/
> cp userManagementWS.xml /var/lib/tomcat7/conf/Catalina/localhost/
for CentOS
> cp userManagementWS.war /var/lib/tomcat/webapps/
> cp userManagementWS.xml /etc/tomcat/Catalina/localhost/

(for CentOS7. The userManagementWS.xml needs to modify the sql connector_
  ↵where driverClassName="org.mariadb.jdbc.Driver")

* Deploy userManagement to tomcat server

for Ubuntu
> cp userManagement.war /var/lib/tomcat7/webapps
for CentOS
> cp userManagement.war /var/lib/tomcat/webapps

* Change settings in /var/lib/tomcat7/webapps/userManagement/WEB-INF/classes/sys.
  ↵properties of Ubuntu.
  ↵          /var/lib/tomcat/webapps/userManagement/WEB-INF/classes/sys.
  ↵properties of CentOS.

```

(continues on next page)

(continued from previous page)

```
host_url=http://www.yourdomain.com:8080/userManagement
email_sender=admin@yourdomain.com
email_host=mail.yourdomain.com
```

**Note:**

**tomcat files in /var/lib/tomcat7 & /usr/share/tomcat7 for Ubuntu** in /var/lib/tomcat & /usr/share/tomcat & /etc/tomcat for CentOS

The tomcat server will automatically decompress the userManagementWS.war and userManagement.war ;

---

6. Setup admin user:

```
* run script createAdminAccount.pl to add admin account with encrypted password to database
> perl createAdminAccount.pl -e admin@my.com -p admin -fn <first name> -ln
<last name>
```

7. Configure the EDGE to use the user management system

- edit \$EDGE\_HOME/edge\_ui/cgi-bin/edge\_config tmpl where user\_management=1

**Note:** If user management system is not in the same domain with edge. ex: <http://www.someother.com/> userManagement The parameter: edge\_user\_management\_url=http://www.someother.com/userManagement

---

8. Enable social (facebook,google,windows live, Linkedin) login function

- edit \$EDGE\_HOME/edge\_ui/cgi-bin/edge\_config tmpl where user\_social\_login=1
- modify \$EDGE\_HOME/edge\_ui/cgi-bin/edge\_user\_management.cgi at line 108/109 of the admin\_email and password according to #6 above.
- modify \$EDGE\_HOME/edge\_ui/javascript/social.js, change apps id you created on each social media.

**Note:** You need to register your EDGE's domain on each social media to get apps id. e.g.: A FACEBOOK app needs to be created and configured for the domain and website set up by EDGE. see <https://developers.facebook.com/> and StackOverflow Q&A

Google+

Windows

LinkedIn

---

9. Optional: configure sendmail to use SMTP to email out of local domain:

```
* edit /etc/mail/sendmail.cf and edit this line:
# "Smart" relay host (may be null)
DS
* and append the correct server right next to DS (no spaces);
```

(continues on next page)

(continued from previous page)

```
# "Smart" relay host (may be null)
DSmail.yourdomain.com

* Then, restart the sendmail service

> sudo service sendmail restart
```

## 4.2 EDGE Docker image

EDGE has a lot of dependencies and can (but doesn't have to) be very challenging to install. The EDGE docker gets around the difficulty of installation by providing a functioning EDGE full install on top of offical Ubuntu 14.04.3 LTS. You can find the image and usage at [docker hub](#).

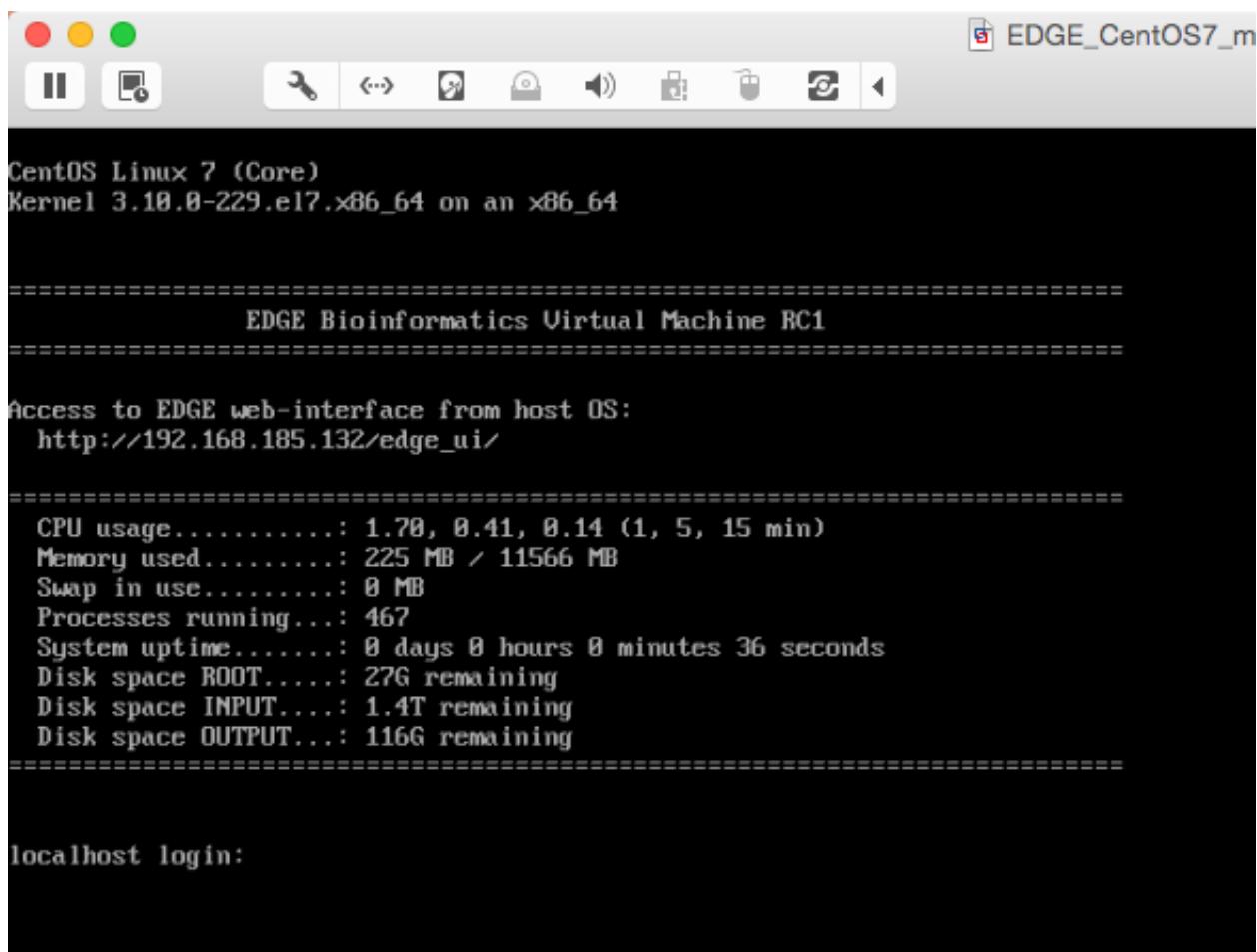
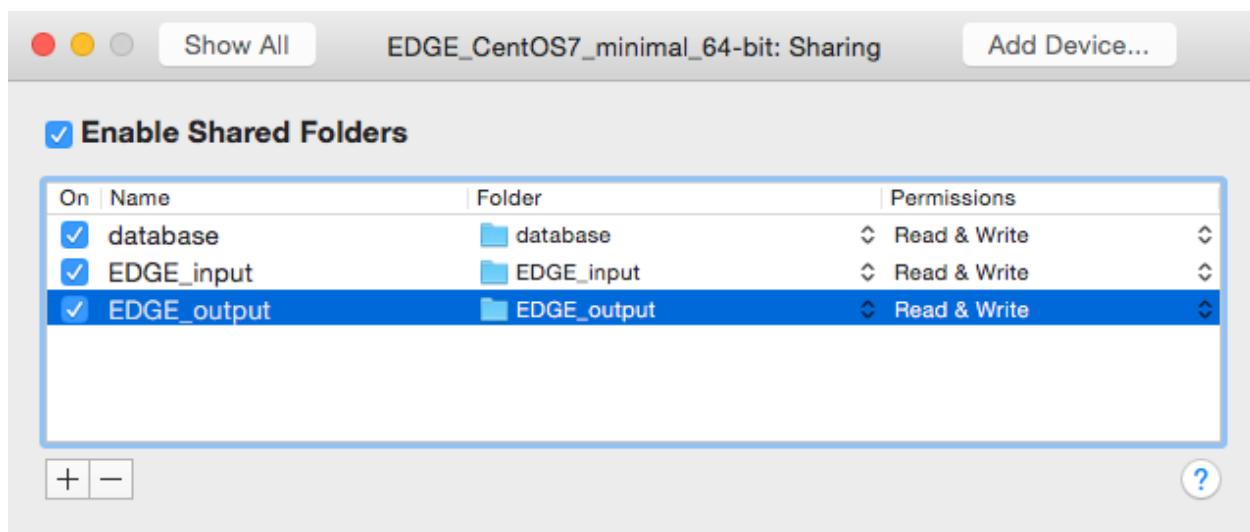
## 4.3 EDGE VMware/OVF Image

You can start using EDGE by launching a local instance of the EDGE VM. The image is built by [VMware Fusion v8.0](#). The pre-built EDGE VM is provided in [Open Virtualization Format \(OVA/OVF\)](#) which is supported by major virtualization players, such as VMware / VirtualBox / Red Hat Enterprise Virtualization, etc. Unfortunately, this may not always work perfectly, as each VM technology seems to use slightly different OVA/OVF implementations that aren't entirely compatible. For example, the [auto-deploy](#) feature and the [path of auto-mount shared folders](#) between host and guest which are used in the EDGE VMware image may not be compatible with other VM technologies (or may need advanced tweaks). Therefore, we highly recommended using [VMware Workstation Player](#) which is free for non-commercial, personal and home use. The [EDGE databases](#) are not included in the image. You will need to download and mount the databases, input and output directories after you launch the VM. Below are instructions to run EDGE VM on your local server:

1. Install [VMware Workstation player](#) .
2. Download VM image (EDGE\_vm\_RC1.ova) from [LANL FTP site](#).
3. Download the [EDGE databases](#) and follow instruction to unpack them.
4. Configure your VM
  - Allocate at least 10GB memory to the VM
  - Share the database, input and output directory to the “database”, “EDGE\_input” and “EDGE\_output” directory in the VM guest OS. If you use VMware, the “Sharing settings” should look like:
5. Start EDGE VM.
6. Access EDGE VM using host browser ([http://<IP\\_OF\\_VM>/edge\\_ui/](http://<IP_OF_VM>/edge_ui/)).

Note that the IP address will also be provided when the instance starts up.

7. Control EDGE VM with default credentials.
  - OS Login: edge/edge
  - EDGE user: [admin@my.edge/admin](mailto:admin@my.edge/admin)
  - MariaDB root: root/edge



# CHAPTER 5

---

## Graphic User Interface (GUI)

---

The User Interface was mainly implemented in [JQuery Mobile](#), CSS, javascript and perl CGI. It is a HTML5-based user interface system designed to make responsive web sites and apps that are accessible on all smartphone, tablet and desktop devices.

See [GUI page](#)

### 5.1 User Login

A user management system has been implemented to provide a level of privacy/security for a user's submitted projects. When this system is activated, any user can view projects that have been made public, but other projects can only be accessed by logging into the system using a registered local EDGE account or via an existing social media account (Facebook, Google+, Windows, or LinkedIn). The users can then run new jobs and view their own previously run projects or those that have been shared with them. Click on the upper-right user icon will pop up an user login window.

The screenshot shows the main landing page of the EDGE bioinformatics website. At the top, there's a navigation bar with links for Home, Upload Files, Run EDGE, and Projects. The main content area features a heading "Empowering the Development of Genomics Expertise". Below this, a paragraph explains the purpose of EDGE bioinformatics, mentioning its user-friendly interface and support for many cutting-edge tools. It also notes that the website is available for use with publicly available data, while local stand-alone implementations can accommodate on-site data analysis. Logos for Los Alamos National Laboratory and NIST are visible. To the right, a large diagram titled "EDGE Workflow" illustrates the analytical pipeline, divided into several stages: 1. Raw sequencing and assembly, 2. Pre-processing, 3. Post-processing, 4. Analysis, and 5. Reporting. Each stage is further detailed with specific steps like "Raw sequencing and assembly", "Quality control", "Assembly", "Annotation", and "Reporting". A central box for "Login to EDGE" includes social media logins for Facebook, Google+, Windows Live, and LinkedIn, along with fields for Email Address and password, and a "Submit" button. Below the login form are links for password reset and sign-up.

## 5.2 Upload Files

For LANL security policy, the function is not implemented at [https://bioedge.lanl.gov/edge\\_ui/](https://bioedge.lanl.gov/edge_ui/).

EDGE supports input from NCBI Sequence Reads Archive (SRA) and select files from the EDGE server. To analyze users' own data, EDGE allows user to upload fastq, fasta and genbank (which can be in gzip format) and text (txt). Max file size is '5gb' and files will be kept for 7 days. Choose "Upload files" from the navigation bar on the left side of the screen. Add users files by clicking "Add Files" button or drag files to the upload feature window. Then, click "Start Upload" button to upload files to EDGE server.

The screenshot shows the "Upload files" page. The left sidebar has links for Home, Upload Files (with a red arrow pointing to it), Run EDGE, and Projects. The main content area has a heading "Upload files" and a note about max file size (5gb) and allowed file types (fastq, fasta, genbank, txt). Below this is a large box with a "Drag files here." placeholder. At the bottom of the box are two buttons: "Add Files" and "Start Upload". Above the "Start Upload" button, file size and progress information (0 b, 0%) are displayed. The top right corner shows a user profile for "Chienchi" and a search icon.

## 5.3 Initiating an analysis job

Choose “Run EDGE” from the navigation bar on the left side of the screen.

The screenshot shows the main interface of the EDGE bioinformatics platform. At the top, there's a green header bar with the logo and name "EDGE bioinformatics" and the URL "@bioedge.lanl.gov". Below the header is a navigation menu with links for "Home", "Upload Files", "Run EDGE", and "Projects". The main content area has a title "Empowering the Development of Genomics Expertise". To the right of the title are logos for Los Alamos National Laboratory and NIST. Below the title is a bulleted list of features:

- No need for high-level bioinformaticists at the site of sequencing.
- No need to move large datasets.
- Edge Bioinformatics enables sequencing as a solution in the biological facilities where human-resources, space, bandwidth, and time are limited
- A highly adaptable bioinformatics platform that allows laboratories to quickly analyze and interpret genomic sequence data.
- Allow users to address a wide range of use cases including assay validation and the characterization of novel biological threats, clinical samples, and complex environmental samples.

Below this list is a complex flowchart titled "The EDGE Environment" showing the data pipeline from raw sequencing data to final reports. The flowchart is organized into several main sections: QC Statistics and Graphics, Raw Sequencing Data FASTQ files, Bio Processing, Clean Data FASTQ files, Quality FASTQ files, Assembly Validation and Annotation, Reference-based Analysis, and finally Interactive Mapping Graphics. Arrows indicate the flow of data between these stages, with many paths merging back into a central hub before reaching the final reporting stage.

This will cause a section to appear called “Input Raw Reads.” Here, you may browse the EDGE Input Directory and select FASTQ files containing the reads to be analyzed. EDGE supports gzip compressed fastq files. At minimum, EDGE will accept two FASTQ files containing paired reads and/or one FASTQ file containing single reads as initial input. Alternatively, rather than providing files through the EDGE Input Directory, you may decide to use as input reads from the Sequence Read Archive (SRA). In this case, select the “yes” option next to “Input from NCBI Sequence Reads Archive” and a field will appear where you can type in an SRA accession number.

The screenshot shows the "Input Your Sample" page with the "Input Raw Reads" section highlighted. The left sidebar contains links for "Home", "Upload Files", "Run EDGE", and "Projects". The main content area has a title "Input Raw Reads". Below the title is a note: "EDGE requires sequence data files in FASTQ format. EDGE allows both paired-end and single-end sequences." The "Input Raw Reads" form includes fields for "Project name" (with a note "(required, at 3 but less than 30 characters)"), "Description" (with a note "(optional)"), and a "Yes" or "No" radio button for "Input from NCBI Short Reads Archive(SRA)". There are three main file upload fields: "Pair-1 FASTQ file", "Pair-2 FASTQ file", and "Single-end FASTQ file", each with a "absolute file path/select file" input field and a circular "additional options" icon. Below these fields is a link "[ additional options ]". At the bottom of the form is a section titled "Batch Project Submission" with a circular icon.

In addition to the input read files, you have to specify a project name. The project name is restricted to only alphanumeric characters and underscores and requires a minimum of three characters. For example, a project name of “E.coli. Project” is not acceptable, but a project name of “E.coli\_project” could be used instead. In the “Description” fields you may enter free text that describes your project. If you would like, you may use as input more reads files than the minimum of 2 paired read files or one file of single reads. To do so, click “additional options” to expose more fields, including two buttons for “Add Paired-end Input” and “Add Single-end Input”.

The screenshot shows the 'Input Raw Reads' section of the EDGE interface. It includes fields for 'Project name' (required, at 3 but less than 30 characters), 'Description' (optional), and a toggle switch for 'Input from NCBI Short Reads Archive(SRA)' where 'Yes' is selected. Below this, there are sections for 'Paired-end reads:' and 'Single-end FASTQ file'. Under 'Paired-end reads:', there are fields for 'Pair-1 FASTQ file' and 'Pair-2 FASTQ file', both with 'absolute file path/select file' buttons. An 'and/or' button is present between these fields and the 'Single-end FASTQ file' field, which also has its own 'absolute file path/select file' button. At the bottom, there is a 'Specify Output Path' field (optional), a 'Use # of CPUs' field set to '8', and a 'Config file' field with an 'absolute file path/select file' button. A note below states: 'Your customized parameters can be used again. You can utilize the file selector above to upload a standard config file generated by EDGE bioinformatics.' At the bottom right, there are 'Add Paired-end Input' and 'Add Single-end Input' buttons with '+' icons.

In the “additional options”, there are several more options, for output path, number of CPUs, and config file. In most cases, you can ignore these options, but they are described briefly below.

### 5.3.1 Output path

You may specify the output path if you would like your results to be output to a specific location. In most cases, you can leave this field blank and the results will be automatically written to a standard location, \$EDGE\_HOME/edge\_ui/EDGE\_output. In most cases, it is sufficient to leave these options to the default settings.

### 5.3.2 Number of CPUs

Additionally, you may specify the number of CPUs to be used. The default and minimum value is one-fourth of total number of server CPUs. You may adjust this value if you wish. Assuming your hardware has 64 CPUs, the default is 16 and the maximum you should choose is 62 CPUs. Otherwise, if the jobs currently in progress use the maximum number of CPUs, the new submitted job will be queued (and colored in grey. Color-coding see [Checking the status of an analysis job](#) (page 31)). For instance, if you have only one job running, you may choose 62 CPUs. However, if you are planning to run 6 different jobs simultaneously, you should divide the computing resources (in this case, 10 CPUs per each job, totaling 60 CPUs for 6 jobs).

### 5.3.3 Config file

Below the “Use # of CPUs” field is a field where you may select a configuration file. A configuration file is automatically generated for each job when you click “Submit.” This field could be used if you wanted to restart a job that hadn’t finished for some reason (e.g. due to power interruption, etc.). This option ensures that your submission will be run exactly the same way as previously, with all the same options.

**See also:**

[Example of config file](#) (page 38)

### 5.3.4 Batch project submission

The “Batch project submission” section is toggled off by default. Clicking on it will open it up and toggle off the “Input Sequence” section at the same time. When you have many samples in “EDGE Input Directory” and would like to run them with the same configuration, instead of submitting several times, you can compile a text file with project name, fastq inputs and optional project descriptions (upload or paste it) and submit through the “Batch project submission” section

**Input Your Sample**

EDGE requires sequence data files in FASTQ format. EDGE allows both paired-end and single-end sequences.

**Input Raw Reads**

**Batch Project Submission**

Run EDGE with Multiple projects using a tools set configuration. Click [\[Sample Input\]](#) to see the example.

Batch text file	<input type="button" value="Browse..."/> No file selected.
or enter your batch sample descriptions here	
Batch sample input:	

## 5.4 Choosing processes/analyses

Once you have selected the input files and assigned a project name and description, you may either click “Submit” to submit an analysis job using the default parameters, or you may change various parameters prior to submitting the job.

The default settings include quality filter and trimming, assembly, annotation, and community profiling. Therefore, if you choose to use default parameters, the analysis will provide an assessment of what organism(s) your sample is composed of, but will not include host removal, primer design, etc. Below the “Input Your Sample” section is a section called “Choose Processes / Analyses”. It is in this section that you may modify parameters if you would like to use settings other than the default settings for your analysis (discussed in detail below).

**Choose Processes / Analyses**

EDGE provides many modules to do various analyses. You can choose to run or skip a specific process. Parameters/options are provided for most of the analyses. You can click here to [turn all on](#), [expand all sections](#) or [close all sections](#).

<input checked="" type="checkbox"/> Pre-processing	<input checked="" type="button"/> On
<input checked="" type="checkbox"/> Assembly and Annotation	<input checked="" type="button"/> On
<input checked="" type="checkbox"/> Reference-based Analysis	<input type="button"/> Off
<input checked="" type="checkbox"/> Taxonomy Classification	<input checked="" type="button"/> On
<input checked="" type="checkbox"/> Phylogenetic Analysis	<input type="button"/> Off
<input checked="" type="checkbox"/> PCR Primer Tools	<input type="button"/> Off

[Submit](#)   [Reset](#)

#### 5.4.1 Pre-processing

Pre-processing is by default on, but can be turned off via the toggle switch on the right hand side. The default parameters should be sufficient for most cases. However, if your experiment involves specialized adapter sequences that need to be trimmed, you may do so in the Quality Trim and Filter subsection. There are two options for adapter trimming. You may either supply a FASTA file containing the adapter sequences to be trimmed, or you may specify N number of bases to be trimmed from either end of each read.

**Pre-processing**

**a. Quality Trim and Filter**

Run Quality Trim and Filter  Yes  No

Trim Quality Level: 5

Average Quality Cutoff: 0

Minimum Read Length: 50

"N" Base Cutoff: 0

Low Complexity Filter Ratio: 0.85

Adapter FASTA: (optional) absolute file path/select file

Cut #bp from 5'-end: 0

Cut #bp from 3'-end: 0

**b. Host Removal**

Run Host Removal  Yes  No

Select Genome(s): Select host genome(s)...

and/or

Host FASTA file: absolute file path/select file

Similarity (%): 90

---

**Note:** Trim Quality Level can be used to trim reads from both ends with defined quality. "N" base cutoff can be used to filter reads which have more than this number of continuous base "N". Low complexity is defined by the fraction of mono-/di-nucleotide sequence. Ref: FaQCs.

---

The host removal subsection allows you to subtract host-derived reads from your dataset, which can be useful for metagenomic (complex) samples such as clinical samples (blood, tissue), or environmental samples like insects. In order to enable host removal, within the "Host Removal" subsection of the "Choose Processes / Analyses" section, switch the toggle box to "On" and select either from the pre-build host list ( Human , Invertebrate Vectors of Human Pathogens , PhiX , RefSeq Bacteria and RefSeq Viruses .) or the appropriate host FASTA file for your experiment from the navigation field. The Similarity (%) can be varied if desired, but the default is 90 and we would not recommend using a value less than 90.

## 5.4.2 Assembly And Annotation

The Assembly option by default is turned on. It can be turned off via the toggle button. EDGE performs iterative kmers de novo assembly by IDBA-UD . It performs well on isolates as well as metagenomes but it may not work well on very large genomes. By default, it starts from kmer=31 and iterative step by adding 20 to maximum kmer=121. When the maximum k value is larger than the input average reads length, it will automatically adjust the maximum value to average reads length minus 1. User can set the minimum cutoff value on the final contigs. By default, it will filter out all contigs with size smaller than 200 bp.

**Assembly and Annotation**

On

Bypass assembly and use pre-assembled contigs  Yes  No

Assembler  IDBA\_UD  SPAdes

IDBA\_UD performs well on isolates as well as metagenomes but it may not work well on very large genomes.

Minimum Kmer Length: 31

Maximum Kmer Length: 121

Step Size: 20

Minimum Contig Length: 200

Annotation  Yes  No

Minimum Contig Length for Annotation: 700

Annotation Tool  Prokka  RATT

Specify Kingdom  Archaea  Bacteria  Mitochondria  Viruses  Others

Please choose the genome type you would like to annotate for Prokka to do genome annotation.

The Annotation module will be performed only if the assembly option is turned on and reads were successfully assembled. EDGE has the option of using Prokka or RATT to do genome annotation. For most cases, Prokka is the appropriate tool to use, however, if your input is a viral genome with attached reference annotation (GenBank file), RATT is the preferred method. If for some reason the assembly fails (ex: run out of Memory), EDGE will bypass any modules requiring a contigs file including the annotation analysis.

## 5.4.3 Reference-based Analysis

The reference-based analysis section allows you to map reads/contigs to the provided references, which can be useful for known isolated species such as cultured samples, to get the coverage information and validate the assembled contigs. In order to enable reference-based analysis, switch the toggle box to “On” and select either from the pre-

build Reference list ( [Ebola virus genomes](#) (page 61) , [E.coli 55989](#) , [E.coli O104H4](#) , [E.coli O127H6](#) and [E.coli K12 MG1655](#) .) or the appropriate FASTA/GenBank file for your experiment from the navigation field.

**Reference-based Analysis**  On

Given one or multiple reference genome FASTA files, EDGE will turn on the analysis of the reads/contigs mapping to reference and JBrowse reference track generation. Given a reference genome genbank file, EDGE will also turn on variant analysis.

Select Genome(s)

and/or

Reference genome

Identify Unmapped Reads  Yes  No

Identify Unmapped Contigs  Yes  No

EDGE will try to classify reads and contigs that are unmapped to references by mapping them to NCBI RefSeq database.

Given a reference genome fasta file, EDGE will turn on the analysis of the reads/contigs mapping to reference and JBrowse reference track generation. If a GenBank file is provided, EDGE will also turn on variant analysis.

#### 5.4.4 Taxonomy Classification

Taxonomic profiling is performed via the “Taxonomy Classification” feature. This is a useful feature not only for complex samples, but also for purified microbial samples (to detect contamination). In the “Community profiling” subsection in the “Choose Processes / Analyses section,” community profiling can be turned on or off via the toggle button.

**Taxonomy Classification**  On

a. Read-based Taxonomy Classification

EDGE will use all reads by default. You can change the behavior to use reads that are unmapped to the reference if Reference-based Analysis is on.

Always use all reads  Yes  No

EDGE uses multiple tools for taxonomy classification including GOTTCHA (bacterial & viral databases), MetaPhlAn, MetaPhyler (short read version), Kraken, MetaScope and reads mapping to NCBI RefSeq using BWA.

Classification Tools

b. Contig-based Taxonomy Classification

EDGE will map contigs to NCBI genomes and make taxonomy inference to each contigs.

Contigs Classification  Yes  No

There is an option to “Always use all reads” or not. If “Always use all reads” is not selected, then only those reads that do not map to the user-supplied reference will be shown in downstream analyses (i.e. the results will only include what is different from the reference). Additionally, the user can use different profiling tools with checkbox selection menu. EDGE uses multiple tools for taxonomy classification including [GOTTCHA](#) (bacterial & viral databases) , [MetaPhlAn](#) , [Kraken](#) and reads mapping to NCBI RefSeq using [BWA](#) .

Turning on the “Contig-Based Taxonomy Classification” section will initiate mapping contigs against NCBI databases for taxonomy and functional annotations.

#### 5.4.5 Phylogenomic Analysis

EDGE supports 5 pre-computed pathogen databases ( *E.coli*, *Yersinia*, *Francisella*, *Brucella*, *Bacillus* (page 54)) for SNP phylogeny analysis. You can also choose to build your own database by first selecting a build method (either FastTree or RAxML), then selecting a pathogen from the “Search Genomes” search function. You can also add FASTA files or SRA Accessions.

**Phylogenetic Analysis**

On

EDGE supports 5 pre-computed databases for SNP phylogeny analysis and two tree builders. FastTree is faster and RAxML is slower but more accurate.

Tree Build Method

FastTree RAxML

Pathogen SNP DB

Select a Pathogen...

or

Select Genome(s)

Search genomes...

Add Genome(s)

absolute FASTA file path/select file

SRA Accessions

ex:SRR2133399,SRR576632

#### 5.4.6 PCR Primer Tools

EDGE includes PCR-related tools for use by those who want to use PCR data for their projects.

The screenshot shows the 'PCR Primer Tools' interface. At the top right is a blue button labeled 'On'. Below it, under 'a. Primer Validation', there is a 'Run Primer Validation' toggle switch set to 'Yes' (blue), and a 'Maximum Mismatch' slider set to 1. A note says: 'Given a primer file, EDGE will run validation of the primer pair to the reference and/or assembled contigs, as available.' Under 'Primer Fasta Sequences', there is a text input field with placeholder 'absolute file path/select file' and a browse icon. Below that, 'Maximum Mismatch' is set to 1. Under 'b. Primer Design', there is another 'Run Primer Design' toggle switch set to 'Yes' (blue). A note says: 'EDGE will design primers based on the assembled contigs.' Below are several input fields with sliders: 'Tm Optimum (C)' at 59, 'Tm Range (C)' from 57 to 63, 'Length Optimum (bp)' at 20, 'Length Range (bp)' from 18 to 27, 'Background Tm Differential (C)' at 5, and 'Number of Primer Pairs' at 5.

#### • Primer Validation

The “Primer Validation” tool can be used to verify whether and where given primer sequences would align to the genome of the sequenced organism. Prior to initiating the analysis, primer sequences in FASTA format must be deposited in the folder on the desktop in the directory entitled “EDGE Input Directory.”

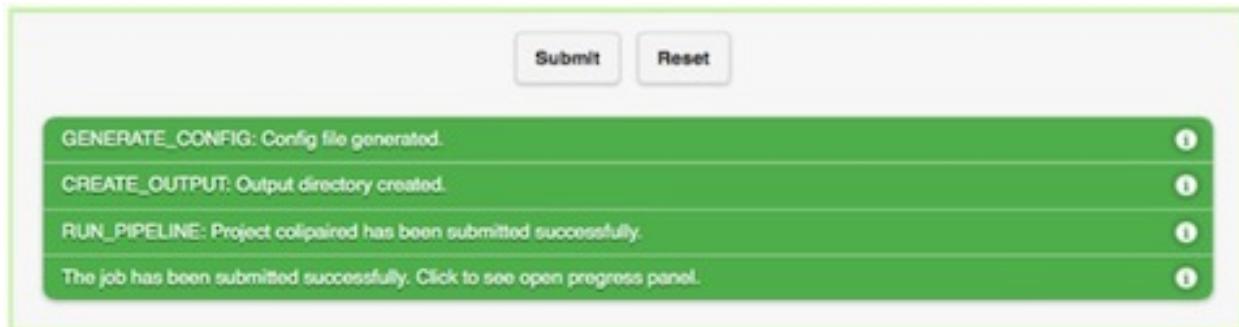
In order to initiate primer validation, within the “Primer Validation” subsection switch the “Run Primer Validation” toggle button to “On”. Then, within the “Primer FASTA Sequences” navigation field, select your file containing the primer sequences of interest. Next, in the “Maximum Mismatch” field, choose the maximum number of mismatches you wish to allow per primer sequence. The available options are 0, 1, 2, 3, or 4.

#### • Primer Design

If you would like to design new primers that will differentiate a sequenced microorganism from all other bacteria and viruses in NCBI, you can do so using the “Primer Design” tool. To initiate primer design switch the “Run Primer Design” toggle button to “On”. There are default settings supplied for Melting Temperature, Primer Length, Tm Differential, and Number of Primer Pairs, but you can change these settings if desired.

## 5.5 Submission of a job

When you have selected the appropriate input files and desired analysis options, and you are ready to submit the analysis job, click on the “Submit” button at the bottom of the page. Immediately you will see indicators of successful job submission and job status below the submit button, in green. If there is something wrong with the input, it will stop the submission and show the message in red, highlighting the sections with issues.

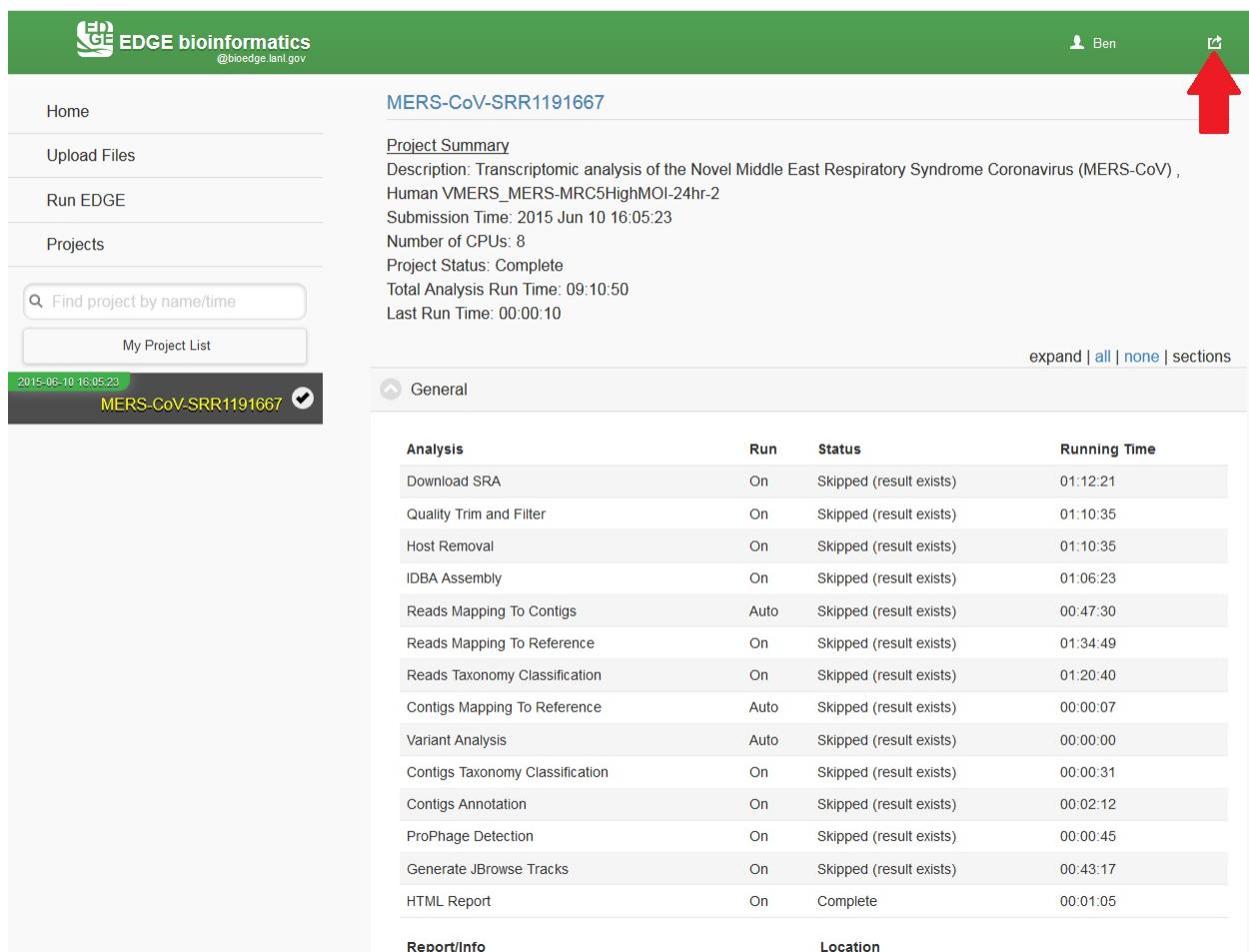


## 5.6 Checking the status of an analysis job

Once an analysis job has been submitted, it will become visible in the left navigation bar. There is a grey, red, orange, green color-coding system that indicates job status as follow:

Status	Not yet begun	Error	In progress (running)	Completed
Color	Grey	Red	Orange	Green

While the job is in progress, clicking on the project in the left navigation bar will allow you to see which individual steps have been completed or are in progress, and results that have already been produced. Clicking the job progress widget at top right opens up a more concise view of progress.



The screenshot shows the EDGE bioinformatics project status page for project MERS-CoV-SRR1191667. The left sidebar includes links for Home, Upload Files, Run EDGE, Projects, and a search bar for 'Find project by name/time'. The main content area displays project details: Description (Transcriptomic analysis of the Novel Middle East Respiratory Syndrome Coronavirus (MERS-CoV)), Human VMERS\_MERS-MRC5HighMOI-24hr-2, Submission Time (2015 Jun 10 16:05:23), Number of CPUs (8), Project Status (Complete), Total Analysis Run Time (09:10:50), and Last Run Time (00:00:10). Below this is a table of analysis steps with columns for Analysis, Run, Status, and Running Time. The table shows various steps like Download SRA, Quality Trim and Filter, Host Removal, IDBA Assembly, etc., with their respective statuses and execution times. At the bottom, there are links for Report/Info and Location.

Analysis	Run	Status	Running Time
Download SRA	On	Skipped (result exists)	01:12:21
Quality Trim and Filter	On	Skipped (result exists)	01:10:35
Host Removal	On	Skipped (result exists)	01:10:35
IDBA Assembly	On	Skipped (result exists)	01:06:23
Reads Mapping To Contigs	Auto	Skipped (result exists)	00:47:30
Reads Mapping To Reference	On	Skipped (result exists)	01:34:49
Reads Taxonomy Classification	On	Skipped (result exists)	01:20:40
Contigs Mapping To Reference	Auto	Skipped (result exists)	00:00:07
Variant Analysis	Auto	Skipped (result exists)	00:00:00
Contigs Taxonomy Classification	On	Skipped (result exists)	00:00:31
Contigs Annotation	On	Skipped (result exists)	00:02:12
ProPhage Detection	On	Skipped (result exists)	00:00:45
Generate JBrowse Tracks	On	Skipped (result exists)	00:43:17
HTML Report	On	Complete	00:01:05

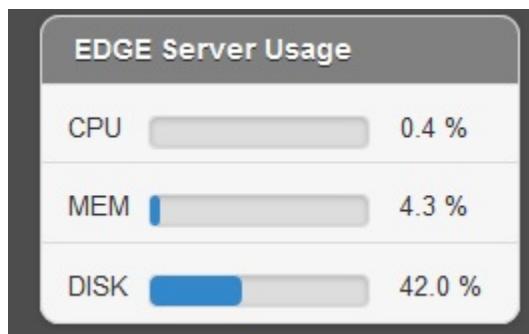
The screenshot shows the EDGE bioinformatics web interface. On the left, there's a sidebar with links for Home, Upload Files, Run EDGE, Projects, a search bar for 'Find project by name/time', and a 'My Project List' button. The main area displays a project summary for 'MERS-CoV-SRR1191667'. It includes a 'Project Summary' section with a description of the transcriptomic analysis of the Novel Middle East Respiratory Syndrome Coronavirus (Human VMERS\_MERS-MRC5HighMOL-24hr-2), submission time (2015 Jun 10 16:05:23), number of CPUs (8), project status (Complete), total analysis run time (09:10:50), and last run time (00:00:10). Below this is a 'General' section with a table of analysis steps and their status:

Analysis	Run	Status
Download SRA	On	Skipped (result exists)
Quality Trim and Filter	On	Skipped (result exists)
Host Removal	On	Skipped (result exists)
IDBA Assembly	On	Skipped (result exists)
Reads Mapping To Contigs	Auto	Skipped (result exists)
Reads Mapping To Reference	On	Skipped (result exists)
Reads Taxonomy Classification	On	Skipped (result exists)
Contigs Mapping To Reference	Auto	Skipped (result exists)
Variant Analysis	Auto	Skipped (result exists)
Contigs Taxonomy Classification	On	Skipped (result exists)
Contigs Annotation	On	Skipped (result exists)
ProPhage Detection	On	Skipped (result exists)
Generate JBrowse Tracks	On	Skipped (result exists)
HTML Report	On	Complete

Below the table is a 'Report/Info' section with a 'Location' column. On the right side of the interface, there are three panels: 'Job Progress' (listing completed steps like 'Download SRA' with green checkmarks), 'EDGE Server Usage' (showing CPU at 0.4%, MEM at 4.3%, and DISK at 42.0%), and 'Action' (with buttons for 'View live log', 'Force to rerun this project', and 'Interrupt running project').

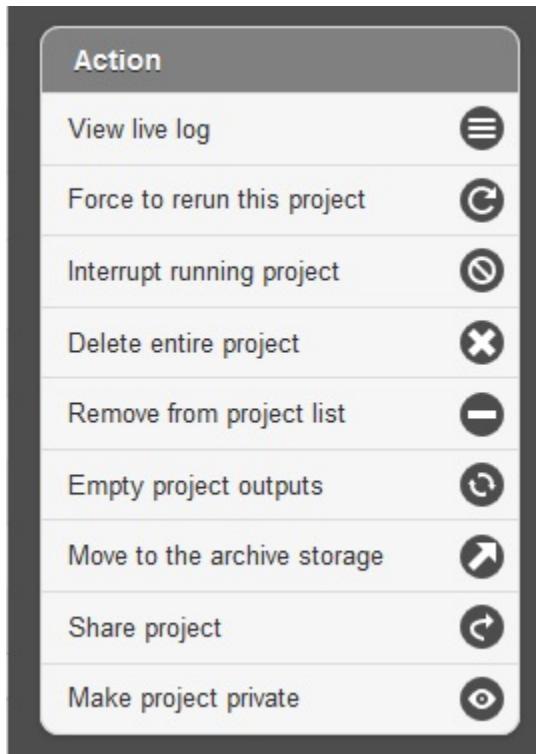
## 5.7 Monitoring the Resource Usage

In the job project sidebar, you can see there is an “EDGE Server Usage” widget that dynamically monitors the server resource usage for %CPU, %MEMORY and %DISK space. If there is not enough available disk space, you may consider deleting or archiving the submitted job with the Action tool described below.



## 5.8 Management of Jobs

Below the resource monitor is the “Action” tool, used for managing jobs in progress or existing projects.



The available actions are:

- **View live log** A terminal-like screen showing all the command lines and progress log information. This is useful for troubleshooting or if you want to repeat certain functions through command line at edge server.
- **Force to rerun this project** Rerun a project with the same inputs and configuration. No additional input needs.
- **Interrupt running project** Immediately stop a running project.
- **Delete entire project** Delete the entire output directory of the project.
- **Remove from project list** Keep the output but remove project name from the project list
- **Empty project outputs** Clean all the results but keep the config file. User can use this function to do a clean rerun.
- **Move to an archive directory** For performance reasons, the output directory will be put in local storage. User can use this function to move projects from local storage to a slower but larger network storage, which are configured when the edge server is installed.
- **Share Project** Allow guests and other users to view the project.
- **Make project Private** Restrict access to viewing the project to only yourself.

## 5.9 Other Methods of Accessing EDGE

### 5.9.1 Internal Python Web Server

EDGE includes a simple web server for single-user applications or other testing. It is not robust enough for production usage, but it is simple enough that it can be run on practically any system.

To run gui, type:

```
$EDGE_HOME/start_edge_ui.sh
```

This will start a localhost and the GUI html page will be opened by your default browser.

### 5.9.2 Apache Web Server

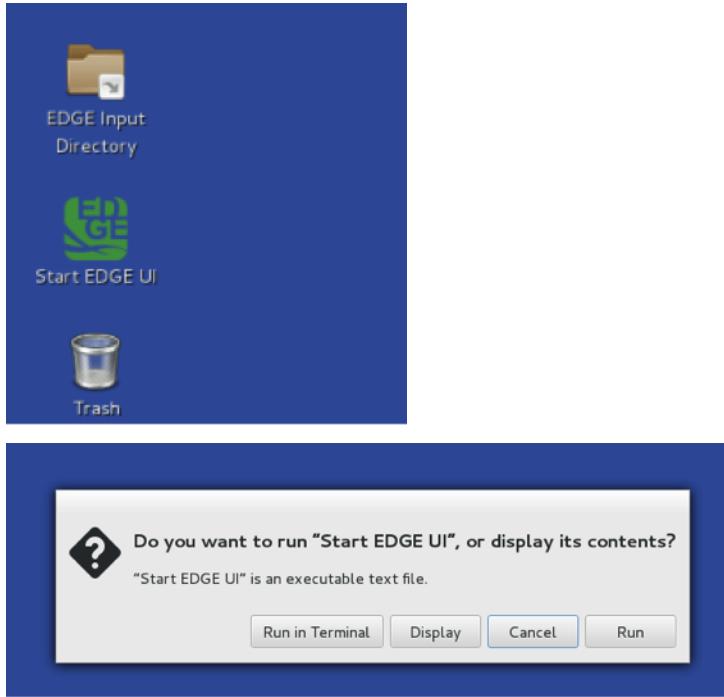
The preferred installation of EDGE uses Apache 2 (See *Apache Web Server Configuration* (page 14)), and serves the application as a proper system service. A sample httpd.conf (or apache2.conf, depending on your operating system) is provided in the root directory of your installation. If this configuration is used, EDGE will be available on any IP or hostname registered to the machine, on ports 80 and 8080.

You can access EDGE by opening either the desktop link (below), or your browser, and entering <http://localhost:80> in the address bar.

---

**Note:** If the desktop environment is available, after installation, a “Start EDGE UI” icon should be on the desktop. Click on the green icon and choose “Run in Terminal.” Results should be the same as those obtained by the above method to start the GUI.

---



The URL address is 127.0.0.1:8080/index.html. It may not be that powerful, as it is hosted by Apache HTTP Server, but it works. With system administrator help, the Apache HTTP Server is the suggested method to host the gui interface.

---

**Note:** You may need to configure the edge\_wwwroot and input and output in the edge\_ui/edge\_config.tmpl file while configuring the Apache HTTP Server and link to external drive or network drive if needed.

---

A Terminal window will display messages and errors as you run EDGE. Under normal operating conditions you can minimize this window. Should an error/problem arise, you may maximize this window to view the error.

```
EDGE
Turning on localhost
webdir: "/Users/218819/Projects/edge-run/edge-ui", port 8888
bash-3.2$ 127.0.0.1 - - [24/Nov/2014 16:58:06] "GET /index.html HTTP/1.1" 200 -
127.0.0.1 - - [24/Nov/2014 16:58:07] "GET /css/edge-output.css HTTP/1.1" 200 -
127.0.0.1 - - [24/Nov/2014 16:58:07] "GET /css/jquery.mobile.1.4.3.min.css HTTP/1.1" 200 -
127.0.0.1 - - [24/Nov/2014 16:58:07] "GET /css/edge.css HTTP/1.1" 200 -
127.0.0.1 - - [24/Nov/2014 16:58:07] "GET /css/jquery.mobile.icons.min.css HTTP/1.1" 200 -
127.0.0.1 - - [24/Nov/2014 16:58:07] "GET /css/jqueryFileTree.css HTTP/1.1" 200 -
127.0.0.1 - - [24/Nov/2014 16:58:07] "GET /css/edge-theme.min.css HTTP/1.1" 200 -
127.0.0.1 - - [24/Nov/2014 16:58:07] "GET /css/tooltipster.css HTTP/1.1" 200 -
127.0.0.1 - - [24/Nov/2014 16:58:07] "GET /css/jquery.lazyloadxt.spinner.min.css HTTP/1.1" 200 -
127.0.0.1 - - [24/Nov/2014 16:58:07] "GET /css/tablesorter.css HTTP/1.1" 200 -
127.0.0.1 - - [24/Nov/2014 16:58:07] "GET /javascript/jquery.js HTTP/1.1" 200 -
127.0.0.1 - - [24/Nov/2014 16:58:07] "GET /javascript/jquery.mobile-1.4.3.min.js HTTP/1.1" 200 -
127.0.0.1 - - [24/Nov/2014 16:58:07] "GET /javascript/edge.js HTTP/1.1" 200 -
127.0.0.1 - - [24/Nov/2014 16:58:07] "GET /javascript/jqueryFileTree.js HTTP/1.1" 200 -
127.0.0.1 - - [24/Nov/2014 16:58:07] "GET /javascript/raphael-min.js HTTP/1.1" 200 -
127.0.0.1 - - [24/Nov/2014 16:58:07] "GET /javascript/jspолосvg-min.js HTTP/1.1" 200 -
127.0.0.1 - - [24/Nov/2014 16:58:07] "GET /javascript/jquery.tooltipster.min.js HTTP/1.1" 200 -
127.0.0.1 - - [24/Nov/2014 16:58:07] "GET /javascript/jquery.lazyloadxt.extra.min.js HTTP/1.1" 200 -
127.0.0.1 - - [24/Nov/2014 16:58:07] "GET /javascript/jquery.tablesorter.min.js HTTP/1.1" 200 -
127.0.0.1 - - [24/Nov/2014 16:58:07] "GET /images/edge_logo.svg HTTP/1.1" 200 -
```

**Warning:** IMPORTANT: Do not close this window!

The Browser window is the window in which you will interact with EDGE.

# CHAPTER 6

---

## Command Line Interface (CLI)

---

The command line usage is as follows:

```
Usage: perl runPipeline.pl [options] -c config.txt -p 'reads1.fastq reads2.fastq' -o ↵
      ↵out_directory
Version 1.1
Input File:
  -u          Unpaired reads, Single end reads in fastq
  -p          Paired reads in two fastq files and separate by space in quote
  -c          Config File
Output:
  -o          Output directory.
Options:
  -ref         Reference genome file in fasta
  -primer      A pair of Primers sequences in strict fasta format
  -cpu         number of CPUs (default: 8)
  -version     print verison
```

A config file (example in the below section, the *Graphic User Interface (GUI)* (page 20) will generate config automatically), reads Files in fastq format, and a output directory are required when run by command line. Based on the configuration file, if all modules are turned on, EDGE will run the following steps. Each step contains at least one command line scripts/programs.

1. Data QC
2. Host Removal QC
3. *De novo* Assembling
4. Reads Mapping To Contig
5. Reads Mapping To Reference Genomes

6. Taxonomy Classification on All Reads or unMapped to Reference Reads
7. Map Contigs To Reference Genomes
8. Variant Analysis
9. Contigs Taxonomy Classification
10. Contigs Annotation
11. ProPhage detection
12. PCR Assay Validation
13. PCR Assay Adjudication
14. Phylogenetic Analysis
15. Generate JBrowse Tracks
16. HTML report

## 6.1 Configuration File

The config file is a text file with the following information. If you are going to do host removal, you need to [build host index](#) (page 54) for it and change the fasta file path in the config file.

```
[Count Fastq]
DoCountFastq=auto

[Quality Trim and Filter]
## boolean, 1=yes, 0=no
DoQC=1
##Targets quality level for trimming
q=5
##Trimmed sequence length will have at least minimum length
min_L=50
##Average quality cutoff
avg_q=0
##"N" base cutoff. Trimmed read has more than this number of continuous base "N" ↴
will be discarded.
n=1
##Low complexity filter ratio, Maximum fraction of mono-/di-nucleotide sequence
lc=0.85
## Trim reads with adapters or contamination sequences
adapter=/PATH/adapter.fasta
## phiX filter, boolean, 1=yes, 0=no
phiX=0
## Cut # bp from 5 end before quality trimming/filtering
5end=0
## Cut # bp from 3 end before quality trimming/filtering
3end=0

[Host Removal]
## boolean, 1=yes, 0=no
DoHostRemoval=1
## Use more Host= to remove multiple host reads
Host=/PATH/all_chromosome.fasta
similarity=90
```

(continues on next page)

(continued from previous page)

```
[Assembly]
## boolean, 1=yes, 0=no
DoAssembly=1
##Bypass assembly and use pre-assembled contigs
assembledContigs=
minContigSize=200
## spades or idba_ud
assembler=idba_ud
idbaOptions="--pre_correction --mink 31"
## for spades
singleCellMode=
pacbioFile=
nanoporeFile=

[Reads Mapping To Contigs]
# Reads mapping to contigs
DoReadsMappingContigs=auto

[Reads Mapping To Reference]
# Reads mapping to reference
DoReadsMappingReference=0
bowtieOptions=
# reference genbank or fasta file
reference=
MapUnmappedReads=0

[Reads Taxonomy Classification]
## boolean, 1=yes, 0=no
DoReadsTaxonomy=1
## If reference genome exists, only use unmapped reads to do Taxonomy Classification.
# Turn on AllReads=1 will use all reads instead.
AllReads=0
enabledTools=gottcha-genDB-b,gottcha-speDB-b,gottcha-strDB-b,gottcha-genDB-v,gottcha-
# speDB-v,gottcha-strDB-v,metaphlan,bwa,kraken_mini

[Contigs Mapping To Reference]
# Contig mapping to reference
DoContigMapping=auto
## identity cutoff
identity=85
MapUnmappedContigs=0

[Variant Analysis]
DoVariantAnalysis=auto

[Contigs Taxonomy Classification]
DoContigsTaxonomy=1

[Contigs Annotation]
## boolean, 1=yes, 0=no
DoAnnotation=1
# kingdom: Archaea Bacteria Mitochondria Viruses
kingdom=Bacteria
contig_size_cut_for_annotation=700
## support tools: Prokka or RATT
annotateProgram=Prokka
```

(continues on next page)

(continued from previous page)

```

annotateSourceGBK=

[ProPhage Detection]
DoProPhageDetection=1

[Phylogenetic Analysis]
DoSNPtree=1
## Available choices are Ecoli, Yersinia, Francisella, Brucella, Bacillus
SNPdbName=Ecoli
## FastTree or RAxML
treeMaker=FastTree
## SRA accessions ByRun, ByExp, BySample, ByStudy
SNP_SRA_ids=

[Primer Validation]
DoPrimerValidation=1
maxMismatch=1
primer=

[Primer Adjudication]
## boolean, 1=yes, 0=no
DoPrimerDesign=0
## desired primer tm
tm_opt=59
tm_min=57
tm_max=63
## desired primer length
len_opt=18
len_min=20
len_max=27
## reject primer having Tm < tm_diff difference with background Tm
tm_diff=5
## display # top results for each target
top=5

[Generate JBrowse Tracks]
DoJBrowse=1

[HTML Report]
DoHTMLReport=1

```

## 6.2 Test Run

EDGE provides an example data set which is an E. coli MiSeq dataset and has been subsampled to ~10x fold coverage reads.

In the EDGE home directory,

```

cd testData
sh runTest.sh

```

See *Output* (page 50)

```
Project Start: 2015 Oct 15 11:26:30
Version: 1.1
The Output Directory path exists
  If you use different input, it may mess up the result with existing files.
[Quality Trim and Filter]
Quality Trim and Filter Finished
[Assembly]
IDBA Assembly Finished
[Reads Mapping To Contigs]
Reads Mapping to Contigs Finished
[Reads Mapping To Reference]
Reads Mapping to Reference Finished
  Unmapped reads retrieved
[Reads Taxonomy Classification]
Reads Taxonomy Classification Finished
[Contigs Mapping To Reference]
Contigs Mapping to Reference Finished
[Variant Analysis]
GFF3 file not exists. Skip Variant Analysis
[Contigs Taxonomy Classification]
Contigs Taxonomy Classification Finished
[Contigs Annotation]
Contig Annotation Finished
[ProPhage Detection]
ProPhage Detection Finished
[Phylogenetic Analysis]
Phylogenetic Analysis Finished
[Primer Validation]
Primer Validation Finished
[Generate JBrowse Tracks]
Generate JBrowse Tracks Finished
Produce Final PDF Report
  Running time: 00:00:02

[HTML Report]
HTML Report Finished
Total Running time: 00:00:02

All Done.
```

Fig. 1: Snapshot from the terminal.

## 6.3 Descriptions of each module

Each module comes with default parameters and user can see the optional parameters by entering the program name with -h or -help flag without any other arguments.

### 1. Data QC

- Required step? No
- Command example

```
perl $EDGE_HOME/scripts/illumina_fastq_QC.pl -p 'Ecoli_10x.1.fastq Ecoli_10x.2.
˓fastq' -q 5 -min_L 50 -avg_q 5 -n 0 -lc 0.85 -d QcReads -t 10
```

- What it does
  - Quality control
  - Read filtering
  - Read trimming
- Expected input
  - Paired-end/Single-end reads in FASTQ format
- Expected output
  - QC.1.trimmed.fastq
  - QC.2.trimmed.fastq
  - QC.unpaired.trimmed.fastq
  - QC.stats.txt
  - QC\_qc\_report.pdf

### 2. Host Removal QC

- Required step? No
- Command example

```
perl $EDGE_HOME/scripts/host_reads_removal_by_mapping.pl -p 'QC.1.trimmed.fastq
˓QC.2.trimmed.fastq' -u QC.unpaired.trimmed.fastq -ref human_chromosomes.fasta -
˓o QcReads -cpu 10
```

- What it does
  - Read filtering
- Expected input
  - Paired-end/Single-end reads in FASTQ format
- Expected output
  - host\_clean.1.fastq
  - host\_clean.2.fastq
  - host\_clean.mapping.log
  - host\_clean.unpaired.fastq
  - host\_clean.stats.txt

### 3. IDBA Assembling

- Required step? No
- Command example

```
fq2fa --merge host_clean.1.fastq host_clean.2.fastq pairedForAssembly.fasta
idba_ud --num_threads 10 -o AssemblyBasedAnalysis/idba --pre_correction_
↳pairedForAssembly.fasta
```

- What it does
  - Iterative kmers de novo Assembly, it performs well on isolates as well as metagenomes. It may not work well on very large genomes.
- Expected input
  - Paired-end/Single-end reads in FASTA format
- Expected output
  - contig.fa
  - scaffold.fa (input paired end)

### 4. Reads Mapping To Contig

- Required step? No
- Command example

```
perl $EDGE_HOME/scripts/runReadsToContig.pl -p 'host_clean.1.fastq host_clean.2.
↳fastq' -d AssemblyBasedAnalysis/readsMappingToContig -pre readsToContigs -ref_
↳AssemblyBasedAnalysis/contigs.fa
```

- What it does
  - Mapping reads to assembled contigs
- Expected input
  - Paired-end/Single-end reads in FASTQ format
  - Assembled Contigs in Fasta format
  - Output Directory
  - Output prefix
- Expected output
  - readsToContigs.alnstats.txt
  - readsToContigs\_coverage.table
  - readsToContigs\_plots.pdf
  - readsToContigs.sort.bam
  - readsToContigs.sort.bam.bai

### 5. Reads Mapping To Reference Genomes

- Required step? No
- Command example:

```
perl $EDGE_HOME/scripts/runReadsToGenome.pl -p 'host_clean.1.fastq host_clean.2.  
-fastq' -d ReadsBasedAnalysis -pre readsToRef -ref Reference.fna
```

- What it does
  - Mapping reads to reference genomes
  - SNPs/Indels calling
- Expected input
  - Paired-end/Single-end reads in FASTQ format
  - Reference genomes in Fasta format
  - Output Directory
  - Output prefix
- Expected output
  - readsToRef.alnstats.txt
  - readsToRef\_plots.pdf
  - readsToRef\_refID.coverage
  - readsToRef\_refID.gap.coords
  - readsToRef\_refID.window\_size\_coverage
  - readsToRef.ref\_windows\_gc.txt
  - readsToRef.raw.bcf
  - readsToRef.sort.bam
  - readsToRef.sort.bam.bai
  - readsToRef.vcf

## 6. Taxonomy Classification on All Reads or unMapped to Reference Reads

- Required step? No
- Command example:

```
perl $EDGE_HOME/scripts/microbial_profiling/microbial_profiling_configure.pl  
-$EDGE_HOME/scripts/microbial_profiling/microbial_profiling.settings.tmpl  
-gottcha-speDB-b > microbial_profiling.settings.ini  
perl $EDGE_HOME/scripts/microbial_profiling/microbial_profiling.pl -o Taxonomy -  
-s microbial_profiling.settings.ini -c 10 UnmappedReads.fastq
```

- What it does
  - Taxonomy Classification using multiple tools, including BWA mapping to NCBI Refseq, metaphlan, kraken, GOTTCCHA.
  - Unify varies output format and generate reports
- Expected input
  - Reads in FASTQ format
  - Configuration text file (generated by microbial\_profiling\_configure.pl)
- Expected output

- Summary EXCEL and text files.
- Heatmaps tools comparison
- Radarchart tools comparison
- Krona and tree-style plots for each tool.

## 7. Map Contigs To Reference Genomes

- Required step? **No**
- Command example:

```
perl $EDGE_HOME/scripts/nucmer_genome_coverage.pl -e 1 -i 85 -p contigsToRef_
  ↵Reference.fna contigs.fa
```

- What it does
  - Mapping assembled contigs to reference genomes
  - SNPs/Indels calling
- Expected input
  - Reference genome in Fasta Format
  - Assembled contigs in Fasta Format
  - Output prefix
- Expected output
  - contigsToRef\_avg\_coverage.table
  - contigsToRef.delta
  - contigsToRef\_query\_unUsed.fasta
  - contigsToRef.snps
  - contigsToRef.coords
  - contigsToRef.log
  - contigsToRef\_query\_novel\_region\_coord.txt
  - contigsToRef\_ref\_zero\_cov\_coord.txt

## 8. Variant Analysis

- Required step? **No**
- Command example:

```
perl $EDGE_HOME/scripts/SNP_analysis.pl -genbank Reference.gbk -SNP contigsToRef_
  ↵snps -format nucmer
perl $EDGE_HOME/scripts/gap_analysis.pl -genbank Reference.gbk -gap contigsToRef_
  ↵ref_zero_cov_coord.txt
```

- What it does
  - Analyze variants and gaps regions using annotation file.
- Expected input
  - Reference in GenBank format
  - SNPs/INDELS/Gaps files from “Map Contigs To Reference Genomes“

- Expected output
  - contigsToRef.SNPs\_report.txt
  - contigsToRef.Indels\_report.txt
  - GapVSReference.report.txt

## 9. Contigs Taxonomy Classification

- Required step? No
- Command example:

```
perl $EDGE_HOME/scripts/contig_classifier_by_bwa/contig_classifier_by_bwa.pl --db
  ↵$EDGE_HOME/database/bwa_index/NCBI-Bacteria-Virus.fna --threads 10 --prefix_
  ↵OuputCT --input contigs.fa
```

- What it does
  - Taxonomy Classification on contigs using BWA mapping to NCBI Refseq
- Expected input
  - Contigs in Fasta format
  - NCBI Refseq genomes bwa index
  - Output prefix
- Expected output
  - prefix.assembly\_class.csv
  - prefix.assembly\_class.top.csv
  - prefix.ctg\_class.csv
  - prefix.ctg\_class.LCA.csv
  - prefix.ctg\_class.top.csv
  - prefix.unclassified.fasta

## 10. Contig Annotation

- Required step? No
- Command example:

```
prokka --force --prefix PROKKA --outdir Annotation contigs.fa
```

- What it does
  - The rapid annotation of prokaryotic genomes.
- Expected input
  - Assembled Contigs in Fasta format
  - Output Directory
  - Output prefix
- Expected output
  - It produces GFF3, GBK and SQN files that are ready for editing in Sequin and ultimately submitted to Genbank/DDJB/ENA.

## 11. ProPhage detection

- Required step? **No**
- Command example:

```
perl $EDGE_HOME/scripts/phageFinder_prepare.pl -o Prophage -p Assembly Annotation/
  ↵PROKKA.gff Annotation/PROKKA.fna
$EDGE_HOME/thirdParty/phage_finder_v2.1/bin/phage_finder_v2.1.sh Assembly
```

- What it does
  - Identify and classify prophages within prokaryotic genomes.
- Expected input
  - Annotated Contigs GenBank file
  - Output Directory
  - Output prefix
- Expected output
  - phageFinder\_summary.txt

## 12. PCR Assay Validation

- Required step? **No**
- Command example:

```
perl $EDGE_HOME/scripts/pcrValidation/validate_primers.pl -ref contigs.fa -primer_
  ↵primers.fa -mismatch 1 -output AssayCheck
```

- What it does
  - In silico PCR primer validation by sequence alignment.
- Expected input
  - Assembled Contigs/Reference in Fasta format
  - Output Directory
  - Output prefix
- Expected output
  - pcrContigValidation.log
  - pcrContigValidation.bam

## 13. PCR Assay Adjudication

- Required step? **No**
- Command example:

```
perl $EDGE_HOME/scripts/pcrAdjudication/pcrUniquePrimer.pl --input contigs.fa --
  ↵gff3 PCR.Adjudication.primers.gff3
```

- What it does
  - Design unique primer pairs for input contigs.
- Expected input

- Assembled Contigs in Fasta format
- Output gff3 file name
- Expected output
  - PCR.Adjudication.primers.gff3
  - PCR.Adjudication.primers.txt

#### 14. Phylogenetic Analysis

- Required step? No
- Command example:

```
perl $EDGE_HOME/scripts/prepare_SNP_phylogeny.pl -o output/SNP_Phylogeny/Ecoli -
  ↪tree FastTree -db Ecoli -n output -cpu 10 -p QC.1.trimmed.fastq QC.2.trimmed.
  ↪fastq -c contigs.fa -s QC.unpaired.trimmed.fastq
perl $EDGE_HOME/scripts/SNPphy/runSNPphylogeny.pl output/SNP_Phylogeny/Ecoli/
  ↪SNPphy.ctrl
```

- What it does
  - Perform SNP identification against selected pre-built SNPdb or selected genomes
  - Build SNP based multiple sequence alignment for all and CDS regions
  - Generate Tree file in newick/PhyloXML format
- Expected input
  - SNPdb path or genomesList
  - Fastq reads files
  - Contig files
- Expected output
  - SNP based phylogenetic multiple sequence alignment
  - SNP based phylogenetic tree in newick/PhyloXML format.
  - SNP information table

#### 15. Generate JBrowse Tracks

- Required step? No
- Command example:

```
perl $EDGE_HOME/scripts/edge2jbrowse_converter.pl --in-ref-fa Reference.fna --in-
  ↪ref-gff3 Reference.gff --proj_outdir EDGE_project_dir
```

- What it does
  - Convert several EDGE outputs into JBrowse tracks for visualization for contigs and reference, respectively.
- Expected input
  - EDGE project output Directory
- Expected output
  - EDGE post-processed files for JBrowse tracks in the JBrowse directory.
  - Tracks configuration files in the JBrowse directory.

## 16. HTML Report

- Required step? No
- Command example:

```
perl $EDGE_HOME/scripts/munger/outputMunger_w_temp.pl EDGE_project_dir
```

- What it does
  - Generate statistical numbers and plots in an interactive html report page.
- Expected input
  - EDGE project output Directory
- Expected output
  - report.html

## 6.4 Other command-line utility scripts

1. To extract certain taxa fasta from contig classification result:

```
cd /home/edge_install/edge_ui/EDGE_output/41/AssemblyBasedAnalysis/Taxonomy
perl /home/edge_install/scripts/contig_classifier_by_bwa/extract_fasta_by_taxa.pl
  ↵-fasta ../contigs.fa -csv ProjectName.ctg_class.top.csv -taxa "Enterobacter
  ↵cloacae" > Ecloacae.contigs.fa
```

2. To extract unmapped/mapped reads fastq from the bam file:

```
cd /home/edge_install/edge_ui/EDGE_output/41/AssemblyBasedAnalysis/
  ↵readsMappingToContig
# extract unmapped reads
perl /home/edge_install/scripts/bam_to_fastq.pl -unmapped readsToContigs.sort.bam
# extract mapped reads
perl /home/edge_install/scripts/bam_to_fastq.pl -mapped readsToContigs.sort.bam
```

3. To extract mapped reads fastq of a specific contig/reference from the bam file:

```
cd /home/edge_install/edge_ui/EDGE_output/41/AssemblyBasedAnalysis/
  ↵readsMappingToContig
perl /home/edge_install/scripts/bam_to_fastq.pl -id ProjectName_00001 -mapped_
  ↵readsToContigs.sort.bam
```

# CHAPTER 7

---

## Output

---

The output directory structure contains ten major sub-directories when all modules are turned on. In addition to the main directories, EDGE will generate a [final report](#) in portable document file format (pdf), process log and error log file in the project main directory.

- AssayCheck
- AssemblyBasedAnalysis
- HostRemoval
- HTML\_Report
- JBrowse
- QcReads
- ReadsBasedAnalysis
- ReferenceBasedAnalysis
- Reference
- SNP\_Phylogeny

In the graphic user interface, EDGE generates an interactive output webpage which includes summary statistics and taxonomic information, etc. The easiest way to interact with the results is through the web interface. If a project run finished through the command line, user can open the report html file in the HTML\_report subdirectory off-line. When a project run is finished, user can click on the project id from the menu and it will generate the interactive html report on the fly. User can browse the data structure by clicking the project link and visualize the result by JBrowse links, download the pdf files, etc.

The screenshot shows the EDGE bioinformatics web interface. At the top, there's a green header bar with the 'EDGE bioinformatics' logo and a user icon. Below the header, on the left, is a sidebar with links for 'Home', 'Run EDGE', 'Projects', a search bar ('Find project by name/time'), and a 'Public Project List' button. On the right, the main content area displays a project summary for 'Ebola\_Virus\_SRX674271'. It includes a 'Project Summary' section with details like 'Description: SRX674271 EBOV sequencing from human serum, RNAseq, 2014 outbreak in Sierra Leone', 'Submission Time: 2015 May 26 17:53:26', 'Number of CPUs: 8', 'Project Status: Complete', 'Total Analysis Run Time: 00:19:43', and 'Last Run Time: 00:19:22'. Below this is a 'Project Details' section with expandable sections for 'General', 'Pre-processing', 'Assembly and Annotation', 'Reference-Based Analysis', 'Taxonomy Classification', and 'PCR Primer Analysis'. At the bottom of the page, there's a footer bar with the text 'EDGE version 1.1' and logos for Los Alamos National Laboratory, the US Department of Energy, and the National Nuclear Security Administration.

## 7.1 Example Output

See [http://lanl-bioinformatics.github.io/EDGE/example\\_output/report.html](http://lanl-bioinformatics.github.io/EDGE/example_output/report.html)

---

**Note:** The example link is just an example of graphic output. The JBrowse and links are not accessible in the example links.

---

# CHAPTER 8

---

## Databases

---

### 8.1 EDGE provided databases

#### 8.1.1 MvirDB

A Microbial database of protein toxins, virulence factors and antibiotic resistance genes for bio-defense applications

- paper: <http://www.ncbi.nlm.nih.gov/pubmed/?term=17090593>
- website: <http://mvirdb.llnl.gov/>

#### 8.1.2 NCBI Refseq

EDGE prebuilt blast db and bwa\_index of NCBI RefSeq genomes.

- Bacteria: <ftp://ftp.ncbi.nih.gov/genomes/Bacteria/all.fna.tar.gz>
  - Version: NCBI 2015 Aug 11
  - 2786 genomes
- Virus: [NCBI Virus](#)
  - Version: NCBI 2015 Aug 11
  - 4834 RefSeq + Neighbor Nucleotides (51300 sequences)

see \$EDGE\_HOME/database/bwa\_index/id\_mapping.txt for all gi/accession to genome name lookup table.

#### 8.1.3 Krona taxonomy

- paper: <http://www.ncbi.nlm.nih.gov/pubmed/?term=21961884>
- website: <http://sourceforge.net/p/krona/home/krona/>

## Update Krona taxonomy db

Download these files from <ftp://ftp.ncbi.nih.gov/pub/taxonomy>:

```
wget ftp://ftp.ncbi.nih.gov/pub/taxonomy/gi_taxid_nucl.dmp.gz
wget ftp://ftp.ncbi.nih.gov/pub/taxonomy/gi_taxid_prot.dmp.gz
wget ftp://ftp.ncbi.nih.gov/pub/taxonomy/taxdump.tar.gz
```

Transfer the files to the taxonomy folder in the standalone KronaTools installation and run:

```
$EDGE_HOME/thirdParty/KronaTools-2.4/updateTaxonomy.sh --local.
```

## 8.1.4 Metaphlan database

MetaPhlAn relies on unique clade-specific marker genes identified from 3,000 reference genomes.

- paper: <http://www.ncbi.nlm.nih.gov/pubmed/?term=22688413>
- website: <http://huttenhower.sph.harvard.edu/metaphlan>

## 8.1.5 Human Genome

The bwa index is prebuilt in the EDGE. The human hs\_ref\_GRCh38 sequences from NCBI ftp site.

- website [ftp://ftp.ncbi.nlm.nih.gov/genomes/H\\_sapiens/Assembled\\_chromosomes/seq/](ftp://ftp.ncbi.nlm.nih.gov/genomes/H_sapiens/Assembled_chromosomes/seq/)

## 8.1.6 MiniKraken DB

Kraken is a system for assigning taxonomic labels to short DNA sequences, usually obtained through metagenomic studies. MiniKraken is a pre-built 4 GB database constructed from complete bacterial, archaeal, and viral genomes in RefSeq (as of Mar. 30, 2014).

- paper: <http://www.ncbi.nlm.nih.gov/pubmed/?term=24580807>
- website: <http://ccb.jhu.edu/software/kraken/>

## 8.1.7 GOTTCHA DB

A novel, annotation-independent and signature-based metagenomic taxonomic profiling tool. (manuscript in submission)

- website: <https://github.com/LANL-Bioinformatics/GOTTCHA>

## 8.1.8 SNPdb

SNP database based on whole genome comparison. Current available db are *Ecoli*, *Yersinia*, *Francisella*, *Brucella*, *Bacillus* (page 54) .

## 8.1.9 Invertebrate Vectors of Human Pathogens

The bwa index is prebuilt in the EDGE.

- paper: <http://www.ncbi.nlm.nih.gov/pubmed/?term=22135296>
- website: <https://www.vectorbase.org>

Version: 2014 July 24

## 8.1.10 Other optional database

Not in the EDGE but you can download.

- NCBI nr/nt blastDB: <ftp://ftp.ncbi.nih.gov/blast/db/>

## 8.2 Building bwa index

Here take human genome as example.

1. Download the human hs\_ref\_GRCh38 sequences from NCBI ftp site.

Go to [ftp://ftp.ncbi.nlm.nih.gov/genomes/H\\_sapiens/Assembled\\_chromosomes/seq/](ftp://ftp.ncbi.nlm.nih.gov/genomes/H_sapiens/Assembled_chromosomes/seq/) Or use a provided perl script in \$EDGE\_HOME/scripts/

```
perl $EDGE_HOME/scripts/download_human_refseq_genome.pl output_dir
```

2. Gunzip the downloaded fasta file and concatenate them into one human genome multifasta file:

```
gunzip hs_ref_GRCh38.*.fa.gz
cat hs_ref_GRCh38.*.fa > human_ref_GRCh38.all.fasta
```

3. Use the installed bwa to build the index:

```
$EDGE_HOME/bin/bwa index human_ref_GRCh38.all.fasta
```

Now, you can configure the config file with “host=/path/human\_ref\_GRCh38.all.fasta” for host removal step.

## 8.3 SNP database genomes

SNP database was pre-built from the below genomes.

### 8.3.1 Ecoli Genomes

Name	Description	URL
Ecoli_042	Escherichia coli 042, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_11128	Escherichia coli O111:H- str. 11128, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_11368	Escherichia coli O26:H11 str. 11368 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_12009	Escherichia coli O103:H2 str. 12009, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_2009EL2050	Escherichia coli O104:H4 str. 2009EL-2050 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>

Table 1 – continued from previous page

Name	Description	URL
Ecoli_2009EL2071	Escherichia coli O104:H4 str. 2009EL-2071 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_2011C3493	Escherichia coli O104:H4 str. 2011C-3493 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_536	Escherichia coli 536, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_55989	Escherichia coli 55989 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_ABU_83972	Escherichia coli ABU 83972 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_APEC_O1	Escherichia coli APEC O1 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_ATCC_8739	Escherichia coli ATCC 8739 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_BL21_DE3	Escherichia coli BL21(DE3) chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_BW2952	Escherichia coli BW2952 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_CB9615	Escherichia coli O55:H7 str. CB9615 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_CE10	Escherichia coli O7:K1 str. CE10 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_CFT073	Escherichia coli CFT073 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_DH1	Escherichia coli DH1, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_Di14	Escherichia coli str. 'clone D i14' chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_Di2	Escherichia coli str. 'clone D i2' chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_E2348_69	Escherichia coli O127:H6 str. E2348/69 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_E24377A	Escherichia coli E24377A chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_EC4115	Escherichia coli O157:H7 str. EC4115 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_ED1a	Escherichia coli ED1a chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_EDL933	Escherichia coli O157:H7 str. EDL933 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_ETEC_H10407	Escherichia coli ETEC H10407, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_HS	Escherichia coli HS, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_IAI1	Escherichia coli IAI1 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_IAI39	Escherichia coli IAI39 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_IHE3034	Escherichia coli IHE3034 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_K12_DH10B	Escherichia coli str. K-12 substr. DH10B chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_K12_MG1655	Escherichia coli str. K-12 substr. MG1655 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_K12_W3110	Escherichia coli str. K-12 substr. W3110, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_KO11FL	Escherichia coli KO11FL chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_LF82	Escherichia coli LF82, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_NA114	Escherichia coli NA114 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_NRG_857C	Escherichia coli O83:H1 str. NRG 857C chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_P12b	Escherichia coli P12b chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_REL606	Escherichia coli B str. REL606 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_RM12579	Escherichia coli O55:H7 str. RM12579 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_S88	Escherichia coli S88 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_SE11	Escherichia coli O157:H7 str. Sakai chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_SE15	Escherichia coli SE11 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_SMS35	Escherichia coli SE15, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_Sakai	Escherichia coli SMS-3-5 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_TW14359	Escherichia coli O157:H7 str. TW14359 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli UM146	Escherichia coli UM146 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_UMN026	Escherichia coli UMN026 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_UMNK88	Escherichia coli UMNK88 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_UTI89	Escherichia coli UTI89 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_W	Escherichia coli W chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Ecoli_Xuzhou21	Escherichia coli Xuzhou21 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Sboydii_CDC_3083_94	Shigella boydii CDC 3083-94 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>
Sboydii_Sb227	Shigella boydii Sb227 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov">http://www.ncbi.nlm.nih.gov</a>

Cont.

Table 1 – continued from previous page

Name	Description	URL
Sdysenteriae_Sd197	Shigella dysenteriae Sd197, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/384137007">http://www.ncbi.nlm.nih.gov/nuccore/384137007</a>
Sflexneri_2002017	Shigella flexneri 2002017 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/162418099">http://www.ncbi.nlm.nih.gov/nuccore/162418099</a>
Sflexneri_2a_2457T	Shigella flexneri 2a str. 2457T, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/108805998">http://www.ncbi.nlm.nih.gov/nuccore/108805998</a>
Sflexneri_2a_301	Shigella flexneri 2a str. 301 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/384120592">http://www.ncbi.nlm.nih.gov/nuccore/384120592</a>
Sflexneri_5_8401	Shigella flexneri 5 str. 8401 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/16120353">http://www.ncbi.nlm.nih.gov/nuccore/16120353</a>
Ssonnei_53G	Shigella sonnei 53G, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/108810166">http://www.ncbi.nlm.nih.gov/nuccore/108810166</a>
Ssonnei_Ss046	Shigella sonnei Ss046 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/170022262">http://www.ncbi.nlm.nih.gov/nuccore/170022262</a>

### 8.3.2 Yersinia Genomes

Name	Description	URL
Ypestis_A1122	Yersinia pestis A1122 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/384137007">http://www.ncbi.nlm.nih.gov/nuccore/384137007</a>
Ypestis_Angola	Yersinia pestis Angola chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/162418099">http://www.ncbi.nlm.nih.gov/nuccore/162418099</a>
Ypestis_Antiqua	Yersinia pestis Antiqua chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/108805998">http://www.ncbi.nlm.nih.gov/nuccore/108805998</a>
Ypestis_CO92	Yersinia pestis CO92 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/16120353">http://www.ncbi.nlm.nih.gov/nuccore/16120353</a>
Ypestis_D106004	Yersinia pestis D106004 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/384120592">http://www.ncbi.nlm.nih.gov/nuccore/384120592</a>
Ypestis_D182038	Yersinia pestis D182038 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/384124469">http://www.ncbi.nlm.nih.gov/nuccore/384124469</a>
Ypestis_KIM_10	Yersinia pestis KIM 10 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/22123922">http://www.ncbi.nlm.nih.gov/nuccore/22123922</a>
Ypestis_Medievalis_Har	Yersinia pestis biovar Medievalis str. Harbin 35 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/384412706">http://www.ncbi.nlm.nih.gov/nuccore/384412706</a>
Ypestis_Microtus_91001	Yersinia pestis biovar Microtus str. 91001 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/45439865">http://www.ncbi.nlm.nih.gov/nuccore/45439865</a>
Ypestis_Nepal516	Yersinia pestis Nepal516 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/108810166">http://www.ncbi.nlm.nih.gov/nuccore/108810166</a>
Ypestis_Pestoides_F	Yersinia pestis Pestoides F chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/145597324">http://www.ncbi.nlm.nih.gov/nuccore/145597324</a>
Ypestis_Z176003	Yersinia pestis Z176003 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/294502110">http://www.ncbi.nlm.nih.gov/nuccore/294502110</a>
Ypseudotuberculosis_IP_31758	Yersinia pseudotuberculosis IP 31758 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/153946813">http://www.ncbi.nlm.nih.gov/nuccore/153946813</a>
Ypseudotuberculosis_IP_32953	Yersinia pseudotuberculosis IP 32953 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/51594359">http://www.ncbi.nlm.nih.gov/nuccore/51594359</a>
Ypseudotuberculosis_PB1	Yersinia pseudotuberculosis PB1/+ chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/186893344">http://www.ncbi.nlm.nih.gov/nuccore/186893344</a>
Ypseudotuberculosis_YPIII	Yersinia pseudotuberculosis YPIII chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/170022262">http://www.ncbi.nlm.nih.gov/nuccore/170022262</a>

### 8.3.3 Francisella Genomes

Name	Description	URL
Fnovicida_U112	Francisella novicida U112 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/118496615">http://www.ncbi.nlm.nih.gov/nuccore/118496615</a>
Ftularen-sis_holarctica_F92	Francisella tularensis subsp. holarctica F92 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/423049750">http://www.ncbi.nlm.nih.gov/nuccore/423049750</a>
Ftularen-sis_holarctica_FSC200	Francisella tularensis subsp. holarctica FSC200 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/422937995">http://www.ncbi.nlm.nih.gov/nuccore/422937995</a>
Ftularen-sis_holarctica_FTNF002-00	Francisella tularensis subsp. holarctica FTNF002-00 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/156501369">http://www.ncbi.nlm.nih.gov/nuccore/156501369</a>
Ftularen-sis_holarctica_LVS	Francisella tularensis subsp. holarctica LVS chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/89255449">http://www.ncbi.nlm.nih.gov/nuccore/89255449</a>
Ftularen-sis_holarctica_OSU18	Francisella tularensis subsp. holarctica OSU18 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/115313981">http://www.ncbi.nlm.nih.gov/nuccore/115313981</a>
Ftularen-sis_mediasiatica_FSC147	Francisella tularensis subsp. mediasiatica FSC147 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/187930913">http://www.ncbi.nlm.nih.gov/nuccore/187930913</a>
Ftularensis_TIGB03	Francisella tularensis TIGB03 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/379716390">http://www.ncbi.nlm.nih.gov/nuccore/379716390</a>
Ftularen-sis_tularensis_FSC198	Francisella tularensis subsp. tularensis FSC198 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/110669657">http://www.ncbi.nlm.nih.gov/nuccore/110669657</a>
Ftularen-sis_tularensis_NE061598	Francisella tularensis subsp. tularensis NE061598 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/385793751">http://www.ncbi.nlm.nih.gov/nuccore/385793751</a>
Ftularen-sis_tularensis_SCHU_S4	Francisella tularensis subsp. tularensis SCHU S4 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/255961454">http://www.ncbi.nlm.nih.gov/nuccore/255961454</a>
Ftularen-sis_tularensis_TI0902	Francisella tularensis subsp. tularensis TI0902 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/379725073">http://www.ncbi.nlm.nih.gov/nuccore/379725073</a>
Ftularen-sis_tularensis_WY96-3418	Francisella tularensis subsp. tularensis WY96-3418 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/134301169">http://www.ncbi.nlm.nih.gov/nuccore/134301169</a>

### 8.3.4 Brucella Genomes

Name	Description	URL
Baborts_1_9941	Brucella abortus bv. 1 str. 9-941	<a href="http://www.ncbi.nlm.nih.gov/bioproject/58019">http://www.ncbi.nlm.nih.gov/bioproject/58019</a>
Baborts_A13334	Brucella abortus A13334	<a href="http://www.ncbi.nlm.nih.gov/bioproject/83615">http://www.ncbi.nlm.nih.gov/bioproject/83615</a>
Baborts_S19	Brucella abortus S19	<a href="http://www.ncbi.nlm.nih.gov/bioproject/58873">http://www.ncbi.nlm.nih.gov/bioproject/58873</a>
Bcanis_ATCC_23365	Brucella canis ATCC 23365	<a href="http://www.ncbi.nlm.nih.gov/bioproject/59009">http://www.ncbi.nlm.nih.gov/bioproject/59009</a>
Bcanis_HSK_A52141	Brucella canis HSK A52141	<a href="http://www.ncbi.nlm.nih.gov/bioproject/83613">http://www.ncbi.nlm.nih.gov/bioproject/83613</a>
Bceti_TE10759_12	Brucella ceti TE10759-12	<a href="http://www.ncbi.nlm.nih.gov/bioproject/229880">http://www.ncbi.nlm.nih.gov/bioproject/229880</a>
Bceti_TE28753_12	Brucella ceti TE28753-12	<a href="http://www.ncbi.nlm.nih.gov/bioproject/229879">http://www.ncbi.nlm.nih.gov/bioproject/229879</a>
Bmelitensis_1_16M	Brucella melitensis bv. 1 str. 16M	<a href="http://www.ncbi.nlm.nih.gov/bioproject/200008">http://www.ncbi.nlm.nih.gov/bioproject/200008</a>
Bmeliten-sis_Abortus_2308	Brucella melitensis biovar Abortus 2308	<a href="http://www.ncbi.nlm.nih.gov/bioproject/16203">http://www.ncbi.nlm.nih.gov/bioproject/16203</a>
Bmeliten-sis_ATCC_23457	Brucella melitensis ATCC 23457	<a href="http://www.ncbi.nlm.nih.gov/bioproject/59241">http://www.ncbi.nlm.nih.gov/bioproject/59241</a>
Bmelitensis_M28	Brucella melitensis M28	<a href="http://www.ncbi.nlm.nih.gov/bioproject/158857">http://www.ncbi.nlm.nih.gov/bioproject/158857</a>
Bmelitensis_M590	Brucella melitensis M5-90	<a href="http://www.ncbi.nlm.nih.gov/bioproject/158855">http://www.ncbi.nlm.nih.gov/bioproject/158855</a>
Bmelitensis_NI	Brucella melitensis NI	<a href="http://www.ncbi.nlm.nih.gov/bioproject/158853">http://www.ncbi.nlm.nih.gov/bioproject/158853</a>
Bmicroti_CCM_4915	Brucella microti CCM 4915	<a href="http://www.ncbi.nlm.nih.gov/bioproject/59319">http://www.ncbi.nlm.nih.gov/bioproject/59319</a>
Bovis_ATCC_25840	Brucella ovis ATCC 25840	<a href="http://www.ncbi.nlm.nih.gov/bioproject/58113">http://www.ncbi.nlm.nih.gov/bioproject/58113</a>
Bpinnipediais_B2_94	Brucella pinnipedialis B2/94	<a href="http://www.ncbi.nlm.nih.gov/bioproject/71133">http://www.ncbi.nlm.nih.gov/bioproject/71133</a>
Bsuis_1330	Brucella suis 1330	<a href="http://www.ncbi.nlm.nih.gov/bioproject/159871">http://www.ncbi.nlm.nih.gov/bioproject/159871</a>
Bsuis_ATCC_23445	Brucella suis ATCC 23445	<a href="http://www.ncbi.nlm.nih.gov/bioproject/59015">http://www.ncbi.nlm.nih.gov/bioproject/59015</a>
Bsuis_VBI22	Brucella suis VBI22	<a href="http://www.ncbi.nlm.nih.gov/bioproject/83617">http://www.ncbi.nlm.nih.gov/bioproject/83617</a>



### 8.3.5 Bacillus Genomes

Name	Description	URL
Banthraxis_A0248	Bacillus anthracis str. A0248, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/229599883">http://www.ncbi.nlm.nih.gov/nuccore/229599883</a>
Banthraxis_Ames	Bacillus anthracis str. 'Ames Ancestor' chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/50196905">http://www.ncbi.nlm.nih.gov/nuccore/50196905</a>
Ban-thracis_Ames_Ancestor	Bacillus anthracis str. Ames chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/30260195">http://www.ncbi.nlm.nih.gov/nuccore/30260195</a>
Banthraxis_CDC_684	Bacillus anthracis str. CDC 684 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/227812678">http://www.ncbi.nlm.nih.gov/nuccore/227812678</a>
Banthraxis_H9401	Bacillus anthracis str. H9401 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/386733873">http://www.ncbi.nlm.nih.gov/nuccore/386733873</a>
Banthraxis_Sterne	Bacillus anthracis str. Sterne chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/49183039">http://www.ncbi.nlm.nih.gov/nuccore/49183039</a>
Bcereus_03BB102	Bacillus cereus 03BB102, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/225862057">http://www.ncbi.nlm.nih.gov/nuccore/225862057</a>
Bcereus_AH187	Bacillus cereus AH187 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/217957581">http://www.ncbi.nlm.nih.gov/nuccore/217957581</a>
Bcereus_AH820	Bacillus cereus AH820 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/218901206">http://www.ncbi.nlm.nih.gov/nuccore/218901206</a>
Bcereus_anthraxis_CI	Bacillus cereus biovar anthracis str. CI chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/301051741">http://www.ncbi.nlm.nih.gov/nuccore/301051741</a>
Bcereus_ATCC_10987	Bacillus cereus ATCC 10987 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/42779081">http://www.ncbi.nlm.nih.gov/nuccore/42779081</a>
Bcereus_ATCC_14579	Bacillus cereus ATCC 14579, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/30018278">http://www.ncbi.nlm.nih.gov/nuccore/30018278</a>
Bcereus_B4264	Bacillus cereus B4264 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/218230750">http://www.ncbi.nlm.nih.gov/nuccore/218230750</a>
Bcereus_E33L	Bacillus cereus E33L chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/52140164">http://www.ncbi.nlm.nih.gov/nuccore/52140164</a>
Bcereus_F837_76	Bacillus cereus F837/76 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/376264031">http://www.ncbi.nlm.nih.gov/nuccore/376264031</a>
Bcereus_G9842	Bacillus cereus G9842 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/218895141">http://www.ncbi.nlm.nih.gov/nuccore/218895141</a>
Bcereus_NC7401	Bacillus cereus NC7401, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/375282101">http://www.ncbi.nlm.nih.gov/nuccore/375282101</a>
Bcereus_Q1	Bacillus cereus Q1 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/222093774">http://www.ncbi.nlm.nih.gov/nuccore/222093774</a>
Bthuringiensis_AlHakam	Bacillus thuringiensis str. Al Hakam chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/118475778">http://www.ncbi.nlm.nih.gov/nuccore/118475778</a>
Bthuringiensis_BMB171	Bacillus thuringiensis BMB171 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/296500838">http://www.ncbi.nlm.nih.gov/nuccore/296500838</a>
Bthuringiensis_Bt407	Bacillus thuringiensis Bt407 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/409187965">http://www.ncbi.nlm.nih.gov/nuccore/409187965</a>
Bthuringiensis_chinensis_CT43	Bacillus thuringiensis serovar chinensis CT-43 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/384184088">http://www.ncbi.nlm.nih.gov/nuccore/384184088</a>
Bthuringiensis_finitimus_YBT020	Bacillus thuringiensis serovar finitimus YBT-020 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/384177910">http://www.ncbi.nlm.nih.gov/nuccore/384177910</a>
Bthuringiensis_konukian_9727	Bacillus thuringiensis serovar konukian str. 97-27 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/49476684">http://www.ncbi.nlm.nih.gov/nuccore/49476684</a>
Bthuringiensis_MC28	Bacillus thuringiensis MC28 chromosome, complete genome	<a href="http://www.ncbi.nlm.nih.gov/nuccore/407703236">http://www.ncbi.nlm.nih.gov/nuccore/407703236</a>

## 8.4 Ebola Reference Genomes

Accession	Description	URL
NC_014372	Tai Forest ebolavirus isolate Tai Forest virus H.sapiens-tc/CIV/1994/Pauleoula-CI, complete genome.	<a href="http://www.ncbi.nlm.nih.gov/nuccore/NC_014372">http://www.ncbi.nlm.nih.gov/nuccore/NC_014372</a>
FJ217162	Cote d'Ivoire ebolavirus, complete genome.	<a href="http://www.ncbi.nlm.nih.gov/nuccore/FJ217162">http://www.ncbi.nlm.nih.gov/nuccore/FJ217162</a>
FJ968794	Sudan ebolavirus strain Boniface, complete genome.	<a href="http://www.ncbi.nlm.nih.gov/nuccore/FJ968794">http://www.ncbi.nlm.nih.gov/nuccore/FJ968794</a>
NC_006432	Sudan ebolavirus isolate Sudan virus H.sapiens-tc/UGA/2000/Gulu-808892, complete genome.	<a href="http://www.ncbi.nlm.nih.gov/nuccore/NC_006432">http://www.ncbi.nlm.nih.gov/nuccore/NC_006432</a>
KJ660348	Zaire ebolavirus isolate H.sapiens-wt/GIN/2014/Gueckedou-C05, complete genome.	<a href="http://www.ncbi.nlm.nih.gov/nuccore/KJ660348">http://www.ncbi.nlm.nih.gov/nuccore/KJ660348</a>
KJ660347	Zaire ebolavirus isolate H.sapiens-wt/GIN/2014/Gueckedou-C07, complete genome.	<a href="http://www.ncbi.nlm.nih.gov/nuccore/KJ660347">http://www.ncbi.nlm.nih.gov/nuccore/KJ660347</a>
KJ660346	Zaire ebolavirus isolate H.sapiens-wt/GIN/2014/Kissidougou-C15, complete genome.	<a href="http://www.ncbi.nlm.nih.gov/nuccore/KJ660346">http://www.ncbi.nlm.nih.gov/nuccore/KJ660346</a>
JN638998	Sudan ebolavirus - Nakisamata, complete genome.	<a href="http://www.ncbi.nlm.nih.gov/nuccore/JN638998">http://www.ncbi.nlm.nih.gov/nuccore/JN638998</a>
AY354458	Zaire ebolavirus strain Zaire 1995, complete genome.	<a href="http://www.ncbi.nlm.nih.gov/nuccore/AY354458">http://www.ncbi.nlm.nih.gov/nuccore/AY354458</a>
AY729654	Sudan ebolavirus strain Gulu, complete genome.	<a href="http://www.ncbi.nlm.nih.gov/nuccore/AY729654">http://www.ncbi.nlm.nih.gov/nuccore/AY729654</a>
EU338380	Sudan ebolavirus isolate EBOV-S-2004 from Sudan, complete genome.	<a href="http://www.ncbi.nlm.nih.gov/nuccore/EU338380">http://www.ncbi.nlm.nih.gov/nuccore/EU338380</a>
KM655246	Zaire ebolavirus isolate H.sapiens-tc/COD/1976/Yambuku-Ecran, complete genome.	<a href="http://www.ncbi.nlm.nih.gov/nuccore/KM655246">http://www.ncbi.nlm.nih.gov/nuccore/KM655246</a>
KC242801	Zaire ebolavirus isolate EBOV/H.sapiens-tc/COD/1976/deRover, complete genome.	<a href="http://www.ncbi.nlm.nih.gov/nuccore/KC242801">http://www.ncbi.nlm.nih.gov/nuccore/KC242801</a>
KC242800	Zaire ebolavirus isolate EBOV/H.sapiens-tc/GAB/2002/Ilembe, complete genome.	<a href="http://www.ncbi.nlm.nih.gov/nuccore/KC242800">http://www.ncbi.nlm.nih.gov/nuccore/KC242800</a>
KC242799	Zaire ebolavirus isolate EBOV/H.sapiens-tc/COD/1995/13709 Kikwit, complete genome.	<a href="http://www.ncbi.nlm.nih.gov/nuccore/KC242799">http://www.ncbi.nlm.nih.gov/nuccore/KC242799</a>
KC242798	Zaire ebolavirus isolate EBOV/H.sapiens-tc/GAB/1996/1Ikot, complete genome.	<a href="http://www.ncbi.nlm.nih.gov/nuccore/KC242798">http://www.ncbi.nlm.nih.gov/nuccore/KC242798</a>
KC242797	Zaire ebolavirus isolate EBOV/H.sapiens-tc/GAB/1996/1Oba, complete genome.	<a href="http://www.ncbi.nlm.nih.gov/nuccore/KC242797">http://www.ncbi.nlm.nih.gov/nuccore/KC242797</a>
KC242796	Zaire ebolavirus isolate EBOV/H.sapiens-tc/COD/1995/13625 Kikwit, complete genome.	<a href="http://www.ncbi.nlm.nih.gov/nuccore/KC242796">http://www.ncbi.nlm.nih.gov/nuccore/KC242796</a>
KC242795	Zaire ebolavirus isolate EBOV/H.sapiens-tc/GAB/1996/1Mbie, complete genome.	<a href="http://www.ncbi.nlm.nih.gov/nuccore/KC242795">http://www.ncbi.nlm.nih.gov/nuccore/KC242795</a>
KC242794	Zaire ebolavirus isolate EBOV/H.sapiens-tc/GAB/1996/2Nza, complete genome.	<a href="http://www.ncbi.nlm.nih.gov/nuccore/KC242794">http://www.ncbi.nlm.nih.gov/nuccore/KC242794</a>

# CHAPTER 9

---

## Third Party Tools

---

### 9.1 Assembly

- IDBA-UD
  - Citation: Peng, Y., et al. (2012) IDBA-UD: a de novo assembler for single-cell and metagenomic sequencing data with highly uneven depth, *Bioinformatics*, 28, 1420-1428.
  - Site: [http://i.cs.hku.hk/~alse/hkubrg/projects/idba\\_ud/](http://i.cs.hku.hk/~alse/hkubrg/projects/idba_ud/)
  - Version: 1.1.1
  - License: GPLv2
- SPAdes
  - Citation: Nurk, Bankevich et al. (2013) Assembling single-cell genomes and mini-metagenomes from chimeric MDA products. *J Comput Biol.* 2013 Oct;20(10):714-37
  - Site: <http://bioinf.spbau.ru/spades>
  - Version: 3.5.0
  - License: GPLv2

### 9.2 Annotation

- RATT
  - Citation: Otto, T.D., et al. (2011) RATT: Rapid Annotation Transfer Tool, *Nucleic acids research*, 39, e57.
  - Site: <http://ratt.sourceforge.net/>
  - Version:
  - License:

- Note: **The original RATT program does not deal with reverse complement strain annotations transfer. We edited the source code to fix it.**
- Prokka
  - Citation: Seemann, T. (2014) Prokka: rapid prokaryotic genome annotation, Bioinformatics, 30,2068-2069.
  - Site: <http://www.vicbioinformatics.com/software.prokka.shtml>
  - Version: 1.11
  - License: GPLv2
  - Note: **The NCBI tool tbl2asn included within PROKKA can have very slow runtimes (up to several hours) while it is dealing with numerous contigs, such as when we input metagenomic data. We modified the code to allow parallel processing using tbl2asn.**
- tRNAscan
  - Citation: Lowe, T.M. and Eddy, S.R. (1997) tRNAscan-SE: a program for improved detection of transfer RNA genes in genomic sequence, Nucleic acids research, 25, 955-964.
  - Site: <http://lowelab.ucsc.edu/tRNAscan-SE/>
  - Version: 1.3.1
  - License: GPLv2
- Barrnap
  - Citation:
  - Site: <http://www.vicbioinformatics.com/software.barrnap.shtml>
  - Version: 0.42
  - License: GPLv3
- BLAST+
  - Citation: Camacho, C., et al. (2009) BLAST+: architecture and applications, BMC bioinformatics, 10, 421.
  - Site: <ftp://ftp.ncbi.nlm.nih.gov/blast/executables/blast+/2.2.29/>
  - Version: 2.2.29
  - License: Public domain
- blastall
  - Citation: Altschul, S.F., et al. (1990) Basic local alignment search tool, Journal of molecular biology, 215, 403-410.
  - Site: <ftp://ftp.ncbi.nlm.nih.gov/blast/executables/release/2.2.26/>
  - Version: 2.2.26
  - License: Public domain
- Phage\_Finder
  - Citation: Fouts, D.E. (2006) Phage\_Finder: automated identification and classification of prophage regions in complete bacterial genome sequences, Nucleic acids research, 34, 5839-5851.
  - Site: <http://phage-finder.sourceforge.net/>
  - Version: 2.1

- License: GPLv3
- Glimmer
  - Citation: Delcher, A.L., et al. (2007) Identifying bacterial genes and endosymbiont DNA with Glimmer, *Bioinformatics*, 23, 673-679.
  - Site: <http://ccb.jhu.edu/software/glimmer/index.shtml>
  - Version: 302b
  - License: Artistic License
- ARAGORN
  - Citation: Laslett, D. and Canback, B. (2004) ARAGORN, a program to detect tRNA genes and tmRNA genes in nucleotide sequences, *Nucleic acids research*, 32, 11-16.
  - Site: <http://mbio-serv2.mbioekol.lu.se/ARAGORN/>
  - Version: 1.2.36
  - License:
- Prodigal
  - Citation: Hyatt, D., et al. (2010) Prodigal: prokaryotic gene recognition and translation initiation site identification, *BMC bioinformatics*, 11, 119.
  - Site: <http://prodigal.ornl.gov/>
  - Version: 2\_60
  - License: GPLv3
- tbl2asn
  - Citation:
  - Site: <http://www.ncbi.nlm.nih.gov/genbank/tbl2asn2/>
  - Version: 24.3 (2015 Apr 29th)
  - License:

**Warning:** tbl2asn must be compiled within the past year to function. We attempt to recompile every 6 months or so. Most recent compilation is 26 Feb 2015

## 9.3 Alignment

- HMMER3
  - Citation: Eddy, S.R. (2011) Accelerated Profile HMM Searches, *PLoS computational biology*, 7, e1002195
  - Site: <http://hmmer.janelia.org/>
  - Version: 3.1b1
  - License: GPLv3
- Infernal
  - Citation: Nawrocki, E.P. and Eddy, S.R. (2013) Infernal 1.1: 100-fold faster RNA homology searches, *Bioinformatics*, 29, 2933-2935.

- Site: <http://infern.al.janelia.org/>
- Version: 1.1rc4
- License: GPLv3
- Bowtie 2
  - Citation: Langmead, B. and Salzberg, S.L. (2012) Fast gapped-read alignment with Bowtie 2, *Nature methods*, 9, 357-359.
  - Site: <http://bowtie-bio.sourceforge.net/bowtie2/index.shtml>
  - Version: 2.1.0
  - License: GPLv3
- BWA
  - Citation: Li, H. and Durbin, R. (2009) Fast and accurate short read alignment with Burrows-Wheeler transform, *Bioinformatics*, 25, 1754-1760.
  - Site: <http://bio-bwa.sourceforge.net/>
  - Version: 0.7.12
  - License: GPLv3
- MUMmer3
  - Citation: Kurtz, S., et al. (2004) Versatile and open software for comparing large genomes, *Genome biology*, 5, R12.
  - Site: <http://mummer.sourceforge.net/>
  - Version: 3.23
  - License: GPLv3

## 9.4 Taxonomy Classification

- Kraken
  - Citation: Wood, D.E. and Salzberg, S.L. (2014) Kraken: ultrafast metagenomic sequence classification using exact alignments, *Genome biology*, 15, R46.
  - Site: <http://ccb.jhu.edu/software/kraken/>
  - Version: 0.10.4-beta
  - License: GPLv3
- Metaphlan
  - Citation: Segata, N., et al. (2012) Metagenomic microbial community profiling using unique clade-specific marker genes, *Nature methods*, 9, 811-814.
  - Site: <http://huttenhower.sph.harvard.edu/metaphlan>
  - Version: 1.7.7
  - License: Artistic License
- GOTTCCHA

- Citation: Tracey Allen K. Freitas, Po-E Li, Matthew B. Scholz, Patrick S. G. Chain (2015) Accurate Metagenome characterization using a hierarchical suite of unique signatures. Nucleic Acids Research (DOI: 10.1093/nar/gkv180)
- Site: <https://github.com/LANL-Bioinformatics/GOTTCHA>
- Version: 1.0b
- License: GPLv3

## 9.5 Phylogeny

- FastTree
  - Citation: Morgan N. Price, Paramvir S. Dehal, and Adam P. Arkin. 2009. FastTree: Computing Large Minimum Evolution Trees with Profiles instead of a Distance Matrix. Mol Biol Evol (2009) 26 (7): 1641-1650
  - Site: <http://www.microbesonline.org/fasttree/>
  - Version: 2.1.7
  - License: GPLv2
- RAxML
  - Citation: Stamatakis,A. 2014. RAxML version 8: A tool for phylogenetic analysis and post-analysis of large phylogenies. Bioinformatics, 30:1312-1313
  - Site: <http://sco.h-its.org/exelixis/web/software/raxml/index.html>
  - Version: 8.0.26
  - License: GPLv2
- Bio::Phylo
  - Citation: Rutger A Vos, Jason Caravas, Klaas Hartmann, Mark A Jensen and Chase Miller, (2011). Bio::Phylo - phyloinformatic analysis using Perl. BMC Bioinformatics 12:63.
  - Site: <http://search.cpan.org/~rvosa/Bio-Phylo/>
  - Version: 0.58
  - License: GPLv3

## 9.6 Visualization and Graphic User Interface

- JQuery Mobile
  - Site: <http://jquerymobile.com>
  - Version :1.4.3
  - License: CC0
- jsPhyloSVG
  - Citation: Smits SA, Ouverney CC, (2010) jsPhyloSVG: A Javascript Library for Visualizing Interactive and Vector-Based Phylogenetic Trees on the Web. PLoS ONE 5(8): e12267.
  - Site: <http://www.jsphylosvg.com>

- Version: 1.55
- License: GPL
- JBrowse
  - Citation: Skinner, M.E., et al. (2009) **JBrowse**: a next-generation genome browser, *Genome research*, 19, 1630-1638.
  - Site: <http://jbrowse.org>
  - Version: 1.11.6
  - License: Artistic License 2.0/LGPLv1
- KronaTools
  - Citation: Ondov, B.D., Bergman, N.H. and Phillippy, A.M. (2011) Interactive metagenomic visualization in a Web browser, *BMC bioinformatics*, 12, 385.
  - Site: <http://sourceforge.net/projects/krona/>
  - Version: 2.4
  - License: BSD

## 9.7 Utility

- BEDTools
  - Citation: Quinlan, A.R. and Hall, I.M. (2010) BEDTools: a flexible suite of utilities for comparing genomic features, *Bioinformatics*, 26, 841-842.
  - Site: <https://github.com/arq5x/bedtools2>
  - Version: 2.19.1
  - License: GPLv2
- R
  - Citation: R Core Team (2013). R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. URL <http://www.R-project.org/>.
  - Site: <http://www.r-project.org/>
  - Version: 2.15.3
  - License: GPLv2
- GNU\_parallel
  - Citation: O. Tange (2011): GNU Parallel - The Command-Line Power Tool, ;login: The USENIX Magazine, February 2011:42-47
  - Site: <http://www.gnu.org/software/parallel/>
  - Version: 20140622
  - License: GPLv3
- tabix
  - Citation:
  - Site: <http://sourceforge.net/projects/samtools/files/tabix/>

- Version: 0.2.6
- License:
- Primer3
  - Citation: Untergasser, A., et al. (2012) Primer3—new capabilities and interfaces, Nucleic acids research, 40, e115.
  - Site: <http://primer3.sourceforge.net/>
  - Version: 2.3.5
  - License: GPLv2
- SAMtools
  - Citation: Li, H., et al. (2009) The Sequence Alignment/Map format and SAMtools, Bioinformatics, 25, 2078-2079.
  - Site: <http://samtools.sourceforge.net/>
  - Version: 0.1.19
  - License: MIT
- FaQCs
  - Citation: Chienchi Lo, Patrick S.G. Chain (2014) Rapid evaluation and Quality Control of Next Generation Sequencing Data with FaQCs. BMC Bioinformatics. 2014 Nov 19;15
  - Site: <https://github.com/LANL-Bioinformatics/FaQCs>
  - Version: 1.34
  - License: GPLv3
- wigToBigWig
  - Citation: Kent, W.J., et al. (2010) BigWig and BigBed: enabling browsing of large distributed datasets, Bioinformatics, 26, 2204-2207.
  - Site: <https://genome.ucsc.edu/goldenPath/help/bigWig.html#Ex3>
  - Version: 4
  - License:
- sratoolkit
  - Citation:
  - Site: <https://github.com/ncbi/sra-tools>
  - Version: 2.4.4
  - License:

# CHAPTER 10

---

## FAQs and Troubleshooting

---

### 10.1 FAQs

- Can I speed up the process?

You may increase the number of CPUs to be used from the “additional options” of the input section. The default and minimum value is one-eighth of total number of server CPUs.

- There is no enough disk space for storing projects data. How do I do?

There is an archive project action which will move the whole project directory to the directory path configured in the \$EDGE\_HOME/sys.properties. We also recommend a symbolic link for the \$EDGE\_HOME/edge\_ui/EDGE\_input directory which points to the location where the user’s (or sequencing center’s) raw data are stored, obviating unnecessary data transfer via web protocol and saving local storage.

- How to decide various QC parameters?

The default parameters should be sufficient for most cases. However, if you have very depth coverage of the sequencing data, you may increase the trim quality level and average quality cutoff to only use high quality data.

- How to set K-mer size for IDBA\_UD assembly?

By default, it starts from kmer=31 and iterative step by adding 20 to maximum kmer=121. Larger K-mers would have higher rate of uniqueness in the genome and would make the graph simpler, but it requires deep sequencing depth and longer read length to guarantee the overlap at any genomic location and it is much more sensitive to sequencing errors and heterozygosity. Professor Titus Brown has a [good blog on general k-mer size discussion](#).

- How many reference genomes for Reference-Based Analysis and Phylogenetic Analysis can be used from the EDGE GUI?

The default maximum is 20 and there is a minimum 3 genomes criteria for the Phylogenetic Analysis. But it can be configured when installing EDGE.

## 10.2 Troubleshooting

- In the GUI, if you are trying to enter information into a specific field and it is grayed out or won't let you, try refreshing the page by clicking the icon in the right top of the browser window.
- Process.log and error.log files may help on the troubleshooting.

### 10.2.1 Coverage Issues

- Average Fold Coverage reported in the HTML output and by the output tables generated in {output directory}/AssemblyBasedAnalysis/ReadsMappingToContigs/ are calculated with mpileup using the default options for metagenomes. These settings discount reads that are unpaired within a contig or with an insert size out of the expected bounds. This will result in an underreporting of the average fold coverage based on the generated BAM file, but one that the team feels is more accurate given the intended use of this environment.

### 10.2.2 Data Migration

- The preferred method of transferring data to the EDGE appliance is via SFTP. Using an SFTP client such as FileZilla, connect to port 22 using your system's username and password.
- In the case of very large transfers, you may wish to use a USB hard drive or thumb drive.
- If the data is being transferred from another LINUX machine, the server will recognize partitions that use the FAT, ext2, ext3, or ext4 filesystems.
- **If the data is being transferred from a Windows machine, the partition may use the NTFS filesystem. If this is the case, the following steps can be taken:**
  - Open the command line interface by clicking the Applications menu in the top left corner (or use SSH to connect to the system).
  - Enter the command: “sudo yum install ntfs-3g ntfs-3g-devel -y”
  - Enter your password if required.
- After a reboot, you should be able to connect your Windows hard drive to the system, and it will mount like a normal disk.

## 10.3 Discussions / Bugs Reporting

- We have created a mailing list for EDGE users. If you would like to receive notifications about the updates and join the discussion, please join the mailing list by becoming the member of edge-users groups.

[EDGE user's google group](#)

- We appreciate any feedback or concerns you may have about EDGE. If you encounter any bugs, you can report them to our GitHub issue tracker.

[Github issue tracker](#)

- Any other questions? You are welcome to [Contact Us](#) (page 72)

# CHAPTER 11

---

## Copyright

---

? Copyright 2013-2019 Los Alamos National Security, LLC. All rights reserved.

Copyright (2013). Triad National Security, LLC. All rights reserved.

This program was produced under U.S. Government contract 89233218CNA000001 for Los Alamos National Laboratory (LANL), which is operated by Triad National Security, LLC for the U.S. Department of Energy/National Nuclear Security Administration.

All rights in the program are reserved by Triad National Security, LLC, and the U.S. Department of Energy/National Nuclear Security Administration. The Government is granted for itself and others acting on its behalf a nonexclusive, paid-up, irrevocable worldwide license in this material to reproduce, prepare derivative works, distribute copies to the public, perform publicly and display publicly, and to permit others to do so.

This is open source software; you can redistribute it and/or modify it under the terms of the GPLv3 License. If software is modified to produce derivative works, such modified software should be clearly marked, so as not to confuse it with the version available from LANL. Full text of the [GPLv3 License](#) can be found in the License file in the main development branch of the repository.

# CHAPTER 12

---

## Contact Us

---

Questions? Concerns? Please feel free to email our google group at [edge-users@googlegroups.com](mailto:edge-users@googlegroups.com) or contact a dev team member listed below.

Name	Email
Patrick Chain	pchain@lanl.gov
Chien-Chi Lo	chienchi@lanl.gov
Paul Li	po-e@lanl.gov
Karen Davenport	kwdavenport@lanl.gov
Joe Anderson	joseph.j.anderson2.civ@mail.mil
Kim Bishop-Lilly	kimberly.a.bishop-lilly.ctr@mail.mil

# CHAPTER 13

---

## Citation

---

Enabling the democratization of the genomics revolution with a fully integrated web-based bioinformatics platform  
Po-E Li; Chien-Chi Lo; Joseph J. Anderson; Karen W. Davenport; Kimberly A. Bishop-Lilly; Yan Xu; Sanaa Ahmed;  
Shihai Feng; Vishwesh P. Mokashi; Patrick S.G. Chain  
Nucleic Acids Research 2016;  
doi: [10.1093/nar/gkw1027](https://doi.org/10.1093/nar/gkw1027)