



**omics***pipe* Documentation  
**Release 1.1.3**

**Kathleen Fisch, Ph.D.**

October 10, 2014



<b>1 Omics Pipe: An Automated Framework for Next Generation Sequencing Analysis</b>	<b>1</b>
1.1 Introduction . . . . .	1
1.2 Installation . . . . .	1
1.3 Tutorials . . . . .	1
1.4 Online Resources . . . . .	1
1.5 Site Navigation . . . . .	2
<b>2 About Omics Pipe</b>	<b>3</b>
2.1 Available Pipelines . . . . .	4
2.2 Users . . . . .	4
2.3 Developers . . . . .	5
2.4 Contact . . . . .	5
<b>3 Using Omics Pipe</b>	<b>7</b>
3.1 Requirements . . . . .	7
3.2 Installation . . . . .	7
3.3 Usage . . . . .	8
3.4 Running Omics Pipe on Amazon Web Services (AWS) . . . . .	8
3.5 Tutorial . . . . .	8
3.6 Version history . . . . .	8
3.7 Documentation . . . . .	8
3.8 Questions . . . . .	8
<b>4 OmicsPipe on the Amazon Cloud (AWS EC2) Tutorial</b>	<b>9</b>
4.1 Step 1: Create an AWS Account . . . . .	9
4.2 Step 2: Load the the OmicsPipe on AWS docker image on your machine . . . . .	9
4.3 Step 3: Configure StarCluster . . . . .	9
4.4 Step 4: Create AWS Volumes . . . . .	10
4.5 Step 5: Launch the Cluster . . . . .	11
4.6 Step 6: Upload data to the cluster . . . . .	11
4.7 Installing extra software . . . . .	12
4.8 To build your own docker image . . . . .	12
4.9 Add storage > 1TB to the cluster using LVM (for advanced users) . . . . .	13
4.10 Add storage > 1TB to the cluster using RAID 0 (for advanced users) . . . . .	14
4.11 Backing up your data to S3 . . . . .	16
<b>5 Omics Pipe Tutorial – Installation</b>	<b>17</b>
5.1 Installation . . . . .	17

5.2	Before Running Omics Pipe: Configuring Parameters File . . . . .	17
5.3	Running Omics Pipe . . . . .	17
5.4	Running Omics Pipe with the Test Data and Parameters . . . . .	17
<b>6</b>	<b>Omics Pipe Tutorial – Configuring the Parameter File</b>	<b>19</b>
6.1	Example Omics Pipe Parameter File . . . . .	19
6.2	Explanation of Variables in Omics Pipe Parameter File . . . . .	21
<b>7</b>	<b>Omics Pipe Tutorial – Creating a Custom Pipeline Script</b>	<b>25</b>
7.1	Designing the structure of the pipeline . . . . .	25
7.2	Creating the script . . . . .	25
7.3	Updating your parameters file . . . . .	27
<b>8</b>	<b>Omics Pipe Tutorial – Adding a New Module (Tool)</b>	<b>29</b>
8.1	1. Create a Bash script . . . . .	29
8.2	2. Create a Python module . . . . .	30
8.3	3. Add custom Python module to your custom pipeline . . . . .	31
8.4	4. Add new parameters to parameters file . . . . .	31
<b>9</b>	<b>Omics Pipe Available Pipelines</b>	<b>33</b>
9.1	RNA-seq (Tuxedo) . . . . .	33
9.2	RNA-seq(Anders 2013) . . . . .	33
9.3	Whole Exome Sequencing (GATK) . . . . .	33
9.4	Whole Genome Sequencing (GATK) . . . . .	34
9.5	Whole Genome Sequencing (MUTECT) . . . . .	34
9.6	ChIP-seq (MACS) . . . . .	34
9.7	ChIP-seq (HOMER) . . . . .	34
9.8	Breast Cancer Personalized Genomics Report- RNAseq . . . . .	35
9.9	TCGA Reanalysis Pipeline - RNAseq . . . . .	35
9.10	TCGA Reanalysis Pipeline - RNAseq Counts . . . . .	35
9.11	miRNAseq Counts (Anders 2013) . . . . .	36
9.12	miRNAseq (Tuxedo) . . . . .	36
9.13	All Available Modules . . . . .	36
<b>10</b>	<b>Reference Databases Needed</b>	<b>37</b>
10.1	All Pipelines . . . . .	37
10.2	Reference Data for Cancer Reporting Scripts (RNAseq cancer, TCGA pipelines) . . . . .	37
10.3	References for Variants (RNA-seq cancer, RNA-seq cancer TCGA, WES and WGS pipelines) . . . . .	38
10.4	WES Pipeline . . . . .	39
10.5	ChIP-seq Pipelines . . . . .	39
10.6	SNPiR Pipelines (RNA-seq cancer and RNA-seq cancer TCGA pipelines) . . . . .	39
<b>11</b>	<b>Third Party Software Dependencies</b>	<b>41</b>
11.1	R Packages Needed . . . . .	41
11.2	RNA-seq (Tuxedo) . . . . .	41
11.3	RNA-seq (Anders 2013) . . . . .	41
11.4	Whole Exome Sequencing (GATK) . . . . .	41
11.5	Whole Genome Sequencing (GATK) . . . . .	42
11.6	Whole Genome Sequencing (MUTECT) . . . . .	42
11.7	ChIP-seq (MACS) . . . . .	42
11.8	ChIP-seq (HOMER) . . . . .	42
11.9	Breast Cancer Personalized Genomics Report- RNAseq . . . . .	43
11.10	TCGA Reanalysis Pipeline - RNAseq . . . . .	43
11.11	TCGA Reanalysis Pipeline - RNAseq Counts . . . . .	43
11.12	miRNAseq Counts (Anders 2013) . . . . .	44

11.13 miRNAseq (Tuxedo) . . . . .	44
<b>12 System Requirements</b>	<b>45</b>
<b>13 RNA-seq Tuxedo Modules</b>	<b>47</b>
13.1 FASTQC . . . . .	47
13.2 TopHat . . . . .	47
13.3 Cufflinks . . . . .	48
13.4 Cuffmerge . . . . .	48
13.5 Cuffmergetocompare . . . . .	49
13.6 Cuffdiff . . . . .	49
13.7 R Summary Report . . . . .	49
<b>14 RNA-seq Count Based Modules</b>	<b>51</b>
14.1 FASTQC . . . . .	51
14.2 STAR Aligner . . . . .	51
14.3 HTSEQ-count . . . . .	52
14.4 R Summary Report - DESEQ2 . . . . .	52
<b>15 Breast Cancer Personalized Genomics Report- RNAseq</b>	<b>53</b>
15.1 FASTQC . . . . .	53
15.2 STAR Aligner . . . . .	53
15.3 HTSEQ-count . . . . .	54
15.4 RSEQC . . . . .	54
15.5 Fusion Catcher . . . . .	55
15.6 BWA/SNPiR . . . . .	55
15.7 Intogen . . . . .	58
15.8 OncoRep Cancer Report . . . . .	59
<b>16 TCGA Reanalysis Pipeline - RNAseq</b>	<b>61</b>
16.1 TCGA Download . . . . .	61
16.2 FASTQC . . . . .	61
16.3 STAR Aligner . . . . .	62
16.4 HTSEQ-count . . . . .	62
16.5 RSEQC . . . . .	63
16.6 Fusion Catcher . . . . .	63
16.7 BWA/SNPiR . . . . .	64
16.8 Intogen . . . . .	67
16.9 OncoRep Cancer Report . . . . .	67
<b>17 RNA-seq Count Based Modules- TCGA</b>	<b>69</b>
17.1 TCGA Download . . . . .	69
17.2 FASTQC . . . . .	69
17.3 STAR Aligner . . . . .	70
17.4 HTSEQ-count . . . . .	70
17.5 R Summary Report - DESEQ2 . . . . .	71
<b>18 miRNA-seq Tuxedo Modules</b>	<b>73</b>
18.1 CutAdapt . . . . .	73
18.2 Fastq Length Filter . . . . .	73
18.3 FASTQC . . . . .	73
18.4 TopHat . . . . .	74
18.5 Cufflinks . . . . .	74
18.6 Cuffmerge . . . . .	75
18.7 Cuffmergetocompare . . . . .	75

18.8	Cuffdiff . . . . .	76
18.9	R Summary Report . . . . .	76
<b>19</b>	<b>miRNA-seq Count Based Modules</b>	<b>77</b>
19.1	CutAdapt . . . . .	77
19.2	Fastq Length Filter . . . . .	77
19.3	FASTQC . . . . .	77
19.4	STAR Aligner . . . . .	78
19.5	HTSEQ . . . . .	78
19.6	R Summary Report - DESEQ2 . . . . .	79
<b>20</b>	<b>ChIP-seq Modules – HOMER</b>	<b>81</b>
20.1	FASTQC . . . . .	81
20.2	ChIP trim . . . . .	81
20.3	Bowtie . . . . .	82
20.4	Read Density -HOMER . . . . .	82
20.5	Peak Detection - HOMER . . . . .	82
20.6	Peak Annotation & Visualization - HOMER . . . . .	83
20.7	Find Motifs - HOMER . . . . .	84
<b>21</b>	<b>ChIP-seq Modules – MACS</b>	<b>85</b>
21.1	FASTQC . . . . .	85
21.2	ChIP trim . . . . .	85
21.3	Bowtie . . . . .	86
21.4	MACS . . . . .	86
<b>22</b>	<b>Whole Genome and Whole Exome Sequencing Modules</b>	<b>87</b>
22.1	FASTQC . . . . .	87
22.2	BWA-MEM . . . . .	87
22.3	PICARD Mark Duplicates . . . . .	88
22.4	GATK Preprocessing . . . . .	88
22.5	GATK Variant Discovery . . . . .	89
22.6	GATK Variant Filtering . . . . .	90
<b>23</b>	<b>Whole Genome Sequencing (MUTECT)</b>	<b>93</b>
23.1	FASTQC . . . . .	93
23.2	BWA-MEM . . . . .	93
23.3	MUTECT . . . . .	94
<b>24</b>	<b>All Available Modules</b>	<b>95</b>
<b>25</b>	<b>Version History</b>	<b>117</b>
25.1	<b>1.1.3</b> (2014/08/22) . . . . .	117
25.2	<b>1.1.2b</b> (2014/08/05) . . . . .	117
25.3	<b>1.1.0</b> (2014/07/09) . . . . .	117
25.4	<b>1.0.16</b> (2014/07/01) . . . . .	117
25.5	<b>1.0.15</b> (2014/06/20) . . . . .	118
<b>26</b>	<b>Copyright &amp; License</b>	<b>119</b>
26.1	Omics Pipe . . . . .	119
	<b>Python Module Index</b>	<b>121</b>

---

# Omics Pipe: An Automated Framework for Next Generation Sequencing Analysis

---

## 1.1 Introduction

Welcome to the documentation for Omics Pipe! Omics pipe is an open-source, modular computational platform that automates ‘best practice’ multi-omics data analysis pipelines published in Nature Protocols and other commonly used pipelines, such as GATK. It currently automates and provides summary reports for two RNA-seq pipelines, variant calling from whole exome sequencing (WES), variant calling and copy number variation analysis from whole genome sequencing (WGS), two ChIP-seq pipelines and a custom RNA-seq pipeline for personalized genomic medicine reporting. It also provides automated support for interacting with the The Cancer Genome Atlas (TCGA) datasets, including automatic download and processing of the samples in this database.

*About Omics Pipe*

## 1.2 Installation

*Local Cluster*

*Amazon Web Services (AWS)*

## 1.3 Tutorials

*Step-by-Step Tutorial*

*Creating a Custom Pipeline*

*Creating Custom Modules*

## 1.4 Online Resources

Homepage:	<a href="http://sulab.org/tools/omics-pipe/">http://sulab.org/tools/omics-pipe/</a>
Repository:	<a href="https://bitbucket.org/sulab/omics_pipe">https://bitbucket.org/sulab/omics_pipe</a>
Online Docs:	<a href="http://packages.python.org/omics_pipe">http://packages.python.org/omics_pipe</a>
Download & PyPI:	<a href="http://pypi.python.org/pypi/omics_pipe">http://pypi.python.org/pypi/omics_pipe</a>

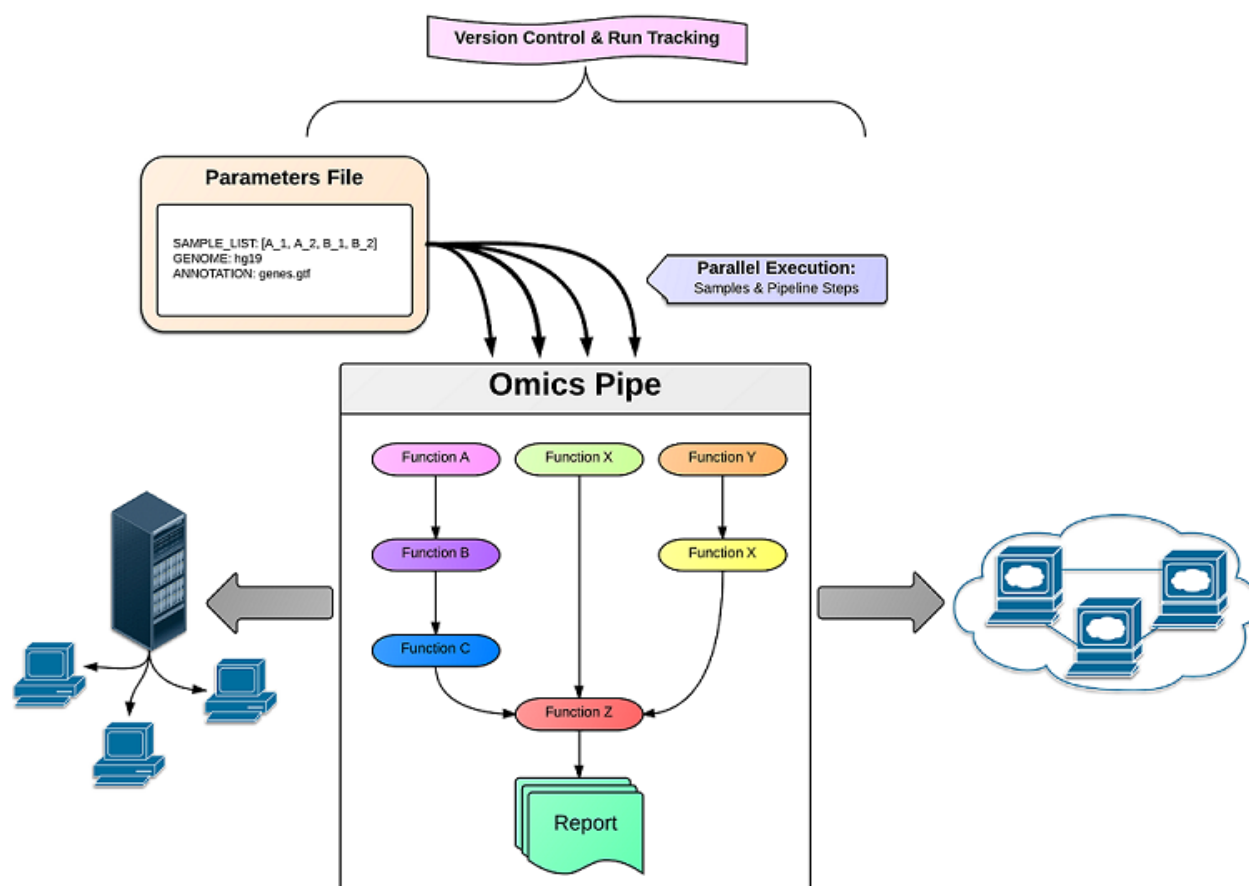
## 1.5 Site Navigation

- [genindex](#)
- [modindex](#)
- [search](#)



## About Omics Pipe

Omics pipe is an open-source, modular computational platform that automates ‘best practice’ multi-omics data analysis pipelines published in [Nature Protocols](#) and other commonly used pipelines, such as [GATK](#). It currently automates and provides summary reports for two RNA-seq pipelines, two miRNA-seq pipelines, variant calling from whole exome sequencing (WES), variant calling and copy number variation analysis from whole genome sequencing (WGS), two CHIP-seq pipelines and a custom RNA-seq pipeline for personalized genomic medicine reporting. It also provides automated support for interacting with the The Cancer Genome Atlas ([TCGA](#)) datasets, including automatic download and processing of the samples in this database.

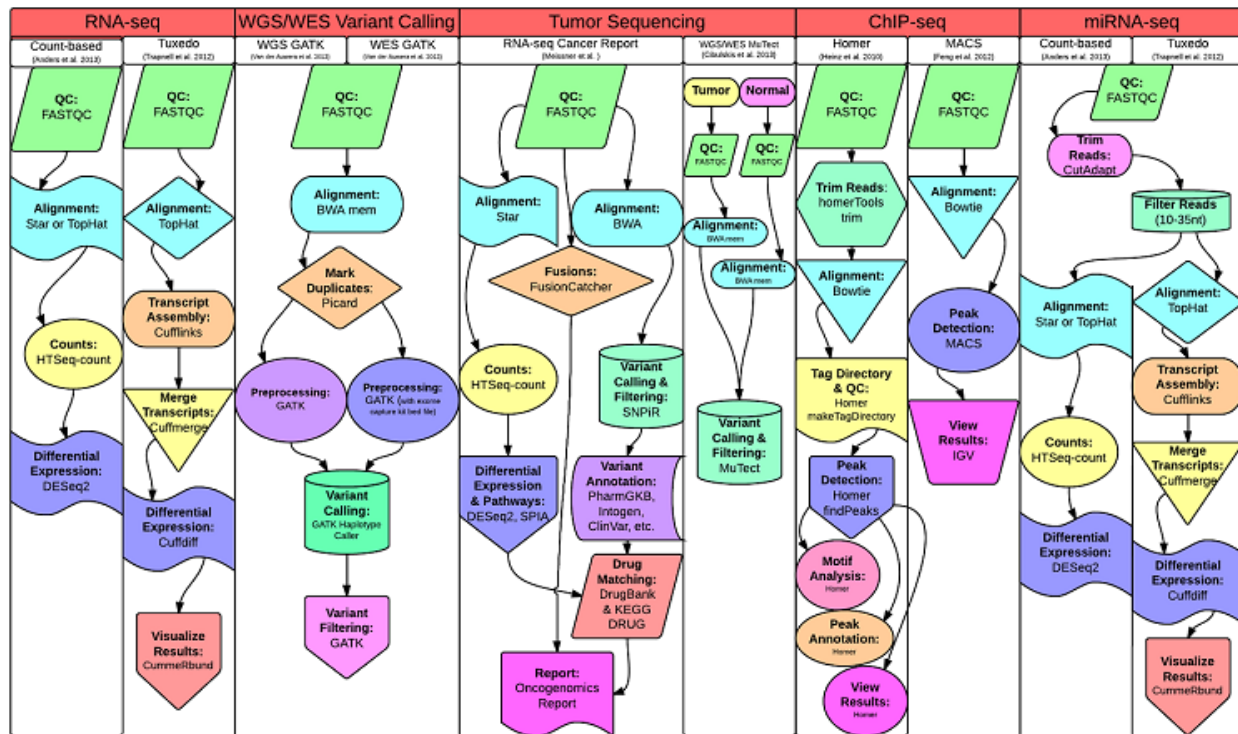


Omics pipe is a Python package that can be installed on a compute cluster, a local installation or in the cloud. It can be downloaded directly from the Omics pipe website for local and cluster installation, or can be used on [AWS](#)

in Amazon EC2. The modular nature of Omics pipe allows researchers to easily and efficiently add new analysis tools with Bash scripts in the form of modules that can then be used to assemble a new analysis pipeline. Omics pipe uses *Ruffus* to pipeline the various analysis modules together into a parallel, automated pipeline. The dependence of Omics pipe on *Ruffus* also allows for the restarting of only the steps in the pipeline that need updating in the event of an error. In addition, *Sumatra* is built into Omics pipe, which provides version control for each run of the pipeline, increasing the reproducibility and documentation of your analyses. Omics pipe interacts with the Distributed Resource Management Application API (*DRMAA*), which automatically submits, controls and monitors jobs to a Distributed Resource Management system, such as a compute cluster or Grid computing infrastructure. This allows you to run samples and steps in the pipeline in parallel in a computationally efficient distributed fashion, without the need to individually schedule and monitor individual jobs. For each supported pipeline in Omics pipe, results files from each step in the pipeline are generated, and an analysis summary report is generated as an HTML report using the R package *knitr*. The summary report provides quality control metrics and visualizations of the results for each sample to enable researchers to quickly and easily interpret the results of the pipeline.

## 2.1 Available Pipelines

*Omics Pipe Available Pipelines* Pipelines supported by this version of omics pipe.



## 2.2 Users

Projects that have used Omics Pipe for solving biological problems. Please submit your story if you would like to share how you use the pipeline for your own research.

- **The Scripps Research Institute, Lotz Lab:** The Lotz Lab in the Department of Molecular and Experimental Medicine at TSRI uses Omics Pipe to perform RNA-seq and miRNA-seq analyses on human articular cartilage samples to elucidate molecular pathways dysregulated in Osteoarthritis.

- **Avera Health:** Researchers working in collaboration with Avera Health use Omics Pipe to analyze sequence data from multiple platforms to provide personalized medicine to breast cancer patients.
- **Scripps Laboratories for tRNA Synthetase Research:** Researchers working in collaboration with Scripps Laboratories for tRNA Synthetase Research use Omics Pipe to analyze ChIP-seq data to explore transcription factor binding sites under experimental conditions.
- **Dorris Neuroscience Center:** Researchers working in collaboration with The Maximov Lab in the Department of Molecular and Cellular Neuroscience at The Scripps Research Institute use Omics Pipe to analyze RNA-seq data to determine how extensive activity-dependent alternative mRNA splicing occurs in the transcriptome of a mouse model that is born and develops to adulthood without synaptic transmission in the forebrain.
- **Sanford Burnham Medical Research Institute:** Researchers working in collaboration with The Peterson Lab in the Bioinformatics and Structural Biology Program at Sanford Burnham Medical Research Institute are using Omics Pipe to perform RNA-seq based global gene expression analysis of dental plaque microbiota derived from twin pairs to identify functional networks of the dental microbiome in relation to dental health and disease.

## 2.3 Developers

Omics Pipe is developed by [Kathleen Fisch](#), [Tobias Meissner](#) and [Louis Gioia](#) at [The Su Lab](#) in the Department of Molecular and Experimental Medicine at [The Scripps Research Institute](#) in beautiful La Jolla, CA.

## 2.4 Contact

Feedback, questions, bug reports, contributions, collaborations, etc. welcome!

[Katie Fisch](#), Ph.D.

Email: [kfisch@scripps.edu](mailto:kfisch@scripps.edu)

Twitter: [@kathleenfisch](#)



---

## Using Omics Pipe

---

Omics Pipe is a Python framework for automating ‘best practice’ next generation sequencing pipelines. Omics Pipe can be run from the command-line by providing it with a YAML parameter file specifying your directory structure and software specific parameters. This executes a parallel automated pipeline on a Distributed Resource Management system (local cluster or Amazon Web Services (AWS)) that efficiently handles job resource allocation, monitoring and restarting. The goals of Omics Pipe are to provide researchers with an open-source computational solution to implement ‘best practice’ pipelines with minimal development overhead and providing visual outputs to aid the researcher in biological interpretation.

To install Omics Pipe, first determine if you are going to be using it on a local compute cluster or on AWS. If you are going to be installing it on your local cluster, follow the directions below (or have your system administrator install it globally). If you are going to create a local installation in your home directory on your cluster but you do not have administrative permissions, you can create a [Python Virtual Environment](#) and then follow the instructions below within the virtual environment.

### 3.1 Requirements

- HPC Cluster or AWS Star Cluster (*Resource Requirements*)
- Python  $\geq 2.6$
- Modules
- **R  $\geq 2.15$** 
  - R Packages (*Third Party Software Dependencies*)
- Third Party Software Dependencies (*Third Party Software Dependencies*)
- Reference Databases (*Reference Databases Needed*)

### 3.2 Installation

- **Option 1:** Install from pypi using pip:

```
pip install omics_pipe
```
- **Option 2:** Install from pypi using easy\_install:

```
easy_install omics_pipe
```

- **Option 3:** Install from source: Download/extract the source code and run:

```
python setup.py install
```

- **Option 4:** Install the latest code directly from the repository:

```
pip install -e hg+https://bitbucket.org/sulab/omics_pipe#egg=omics_pipe
```

- **Option 5:** If you do not have administrator privileges on your system:

Step 1: Set up a 'Python Virtual Environment'

Step 2: Use one of the Options (1-4) above to install Omics Pipe within your virtual environment

## 3.3 Usage

Once you have successfully installed Omics Pipe, you can run a pipeline by typing the command:

```
omics_pipe [-h] [--custom_script_path CUSTOM_SCRIPT_PATH]
           [--custom_script_name CUSTOM_SCRIPT_NAME]
           [--compression {gzip, bzip}]
           {RNAseq_Tuxedo, RNAseq_count_based, RNAseq_cancer_report, RNAseq_TCGA, RNAseq_TCGA_counts,
           parameter_file
```

## 3.4 Running Omics Pipe on Amazon Web Services (AWS)

*AWS Installation Instructions* Installation instructions for setting up the AWS Omics Pipe AMI

## 3.5 Tutorial

*Tutorial* Step-by-step tutorial for running Omics Pipe

*Creating a custom pipeline* Tutorial for creating and running a custom pipeline in Omics Pipe using existing modules

*Adding new modules/tools* Tutorial for adding new modules to Omics Pipe be used in a custom pipeline

## 3.6 Version history

*Version History*

## 3.7 Documentation

The latest copy of this documentation should always be available at: [http://packages.python.org/omics\\_pipe](http://packages.python.org/omics_pipe)

## 3.8 Questions

Email: [kfish@scripps.edu](mailto:kfish@scripps.edu)

Twitter: @kathleenfish

---

## OmicsPipe on the Amazon Cloud (AWS EC2) Tutorial

---

OmicsPipe on AWS uses a custom [StarCluster](#) image, created with [docker.io](#) (which installs [docker.io](#), [environment-modules](#), and [easybuild](#) on an AWS EC2 cluster). All you have to do is get the docker image, upload your data, launch the Amazon cluster and run a single command to analyze all of your data according to published, best-practice methods.

### 4.1 Step 1: Create an AWS Account

1. Create an AWS account by following the instructions at [Amazon-AWS](#)
2. Note your AWS ACCESS KEY ID, AWS SECRET ACCESS KEY and AWS USER ID

### 4.2 Step 2: Load the the OmicsPipe on AWS docker image on your machine

1. Download [docker.io](#) following the instructions for your operating system at [Get-Docker](#)
2. From inside the Docker environment, run the command:

```
docker run -i -t omicspipe/aws_readymade /bin/bash
```

---

**Note:** If you want to share a file from your local computer with the docker container, follow the instructions for [Docker Folder Sharing](#), put your desired file within the shared folder and run the command below (this is recommended for saving your `./starcluster/config` file from the next step:

```
docker run -it --volumes-from NameofSharedDataFolder omicspipe/aws_readymade /bin/bash
```

---

- If you are on a local Ubuntu installation, skip this step and [install](#) the StarCluster client directly.
- If you are using Windows, it might be necessary to update your BIOS to [enable virtualization](#) before installing Docker

### 4.3 Step 3: Configure StarCluster

1. After running the `omicspipe/aws_readymade` Docker container, run the command below to edit the StarCluster configuration file:

```
nano ~/.starcluster/config
```

Or if you prefer vim::

```
vim ~/.starcluster/config
```

2. Enter your “AWS ACCESS KEY ID”, “AWS SECRET ACCESS KEY”, and “AWS USER ID”
3. Change the AWS REGION NAME and AWS REGION HOST variables if you do not live in the AWS us-west region to the appropriate region [AWS Regions](#).
4. Select your desired pre-configured cluster by editing the “DEFAULT\_TEMPLATE” variable or creating a custom cluster. The default is a test cluster with 5 c3.large nodes.
5. Save the edited file (Instructions for [Nano](#) and for [Vim](#))
6. Create your starcluster SSH key by running the command:

```
starcluster createkey omicspipe -o ~/.ssh/omicspipe.rsa
```

- To remove a key from the AWS registry, run the command:

```
starcluster removekey omicspipe
```

- For more information on editing the StarCluster configuration file, see the [StarCluster](#) website.

## 4.4 Step 4: Create AWS Volumes

1. Create AWS volumes to store the raw data and results of your analyses. From within the Docker environment, run:

```
starcluster createvolume --name=data -i ami-52112317 -d -s <volume size in GB> us-west-1a
```

```
starcluster createvolume --name=results -i ami-52112317 -d -s <volume size in GB> us-west-1a
```

- Specify the <volume size in GB> as a number large enough to accommodate all of your raw data and ~4x that size for your results folder
  - Change us-west-1b to your region as described in [AWS Regions](#).
2. Make a volume from the provided snapshot of reference databases (currently only supports *H. sapiens*)
    - Go to the [AWS-Console](#)
    - Click on the [EC2 option](#)
    - Click on Volumes
    - Click on “Create Volume”
    - Set availability zone
    - In Snapshot ID search for “omicspipe\_db” and click on the resulting Snapshot ID
    - Click Create
    - From the Volumes tab, note the “VOLUME\_ID” of the database snapshot
  3. Edit your StarCluster configuration file to add your volume IDs. Run the command below and edit the VOLUME\_ID variables for data, results, and database:



```
nano ~/.starcluster/config
```

Edit the fields below:

```
[volume results]
VOLUME_ID =
MOUNT_PATH = /data/results
```

```
[volume data]
VOLUME_ID =
MOUNT_PATH = /data/data
```

```
[volume database]
VOLUME_ID =
MOUNT_PATH = /data/database
```

4. Save your StarCluster configuration file to `~/.starcluster/config`

## 4.5 Step 5: Launch the Cluster

1. From the Docker container, run the command below to start a new cluster with the name “mypipe”:

```
starcluster start mypipe
```

2. Optional but Recommended: To load balance the cluster, type the command below:

```
starcluster loadbalance mypipe
```

(see ‘**load balance**’ for configuration options. Note: make sure to keep at least one worker node attached to the cluster)

3. SSH into the cluster by running the command below:

```
starcluster sshmaster mypipe
```

## 4.6 Step 6: Upload data to the cluster

Now that you are in your cluster, you can use it like any other cluster. Before running omics pipe on your own data (you can skip this step if you are running the test data, you will want to upload your data. There are two options to upload your data:

1. Upload data from your local machine or cluster using [StarCluster put](#):

```
starcluster put mypipe <myfile> /data/raw
```

2. Use [Webmin](#) to transfer files from your local system to the cluster (recommended for small files only, like parameter files).

- In the AWS Management Console go to “Security Groups”
- Select the “StarCluster-0\_95\_5” group associated with your cluster’s name
- On the Inbound tab click on “Edit”
- Click on “Add Rule” and a new “Custom TCP Rule” will appear. On “Port Range” enter “10000” and on “Source” select “My IP”

- Hit “Save”
- Select Instances in the AWS management console and note the “Public IP” of your instance
- In a Web browser, enter [https://the\\_public\\_ip:10000](https://the_public_ip:10000). Type in the Login info when prompted: user: root password: sulab
- This will take a few seconds to load, and once it does, you can navigate your cluster’s file structure with the tabs on the left
- To upload a file from your local file system, click “upload” and specify the directory /data/data to upload your data.

## 4.7 Installing extra software

Both the [GATK](#) and [MuTect](#) software are used by OmicsPipe, but they require licenses from The Broad Institute and cannot be distributed with the OmicsPipe software. GATK and MuTect are free to download after accepting the license agreement.

To install GATK:

1. [Download GATK](#)
2. Upload the GenomeAnalysisTK.jar file to the /root/.local/easybuild/software/gatk/3.2-2 using either [Webmin](#) or [StarCluster put](#)
3. Make the jar file executable by running the command:

```
chmod +x /root/.local/easybuild/software/gatk/3.2-2/GenomeAnalysisTK.jar
```

To install MuTect:

1. [Download MuTect](#)
2. Upload the muTect-1.1.4.jar file to the /root/.local/easybuild/software/mutect/1.1.4 using either [Webmin](#) or [StarCluster put](#)
3. Make the jar file executable by running the command:

```
chmod +x /root/.local/easybuild/software/mutect/1.1.4/muTect-1.1.4.jar
```

Adding software that OmicsPipe was not built with might require a little more configuration, but OmicsPipe is designed as a foundation to which new software can be added. New software can obviously be added in any manner that the user prefers, but to follow the structure that was used to build OmicsPipe, please refer to the “custombuild” scripts.

---

### Important:

- If you configure software that you think extends the functionality of OmicsPipe, please create a pull request on our [Bitbucket](#) page.
- 

## 4.8 To build your own docker image

1. Download docker.io following the instructions at [Get-Docker](#)
2. Run the command:

```
docker build -t <Repository Name> https://bitbucket.org/sulab/omics_pipe/downloads/Dockerfile_AW
```

This will store the dockercluster image in the Repository Name of your choice.

There is also an [AWS\\_custombuild Dockerfile](#), which can be used to build an Amazon Machine Image from scratch

## 4.9 Add storage > 1TB to the cluster using LVM (for advanced users)

1. Within StarCluster create x new volumes by running:

```
nvolumes=2 #number of volumes
vsize=1000 #in gb
instance='curl -s http://169.254.169.254/latest/meta-data/instance-id'
akey=<AWS KEY>
skey=<AWS SECRET KEY>
region=us-west-1
zone=us-west-1a

for x in $(seq 1 $nvolumes)
do
    ec2-create-volume \
        --aws-access-key $akey \
        --aws-secret-key $skey \
        --size $vsize \
        --region $region \
        --availability-zone $zone
done > /tmp/vols.txt
```

2. Attach the volumes to the head node:

```
i=0
for vol in $(awk '{print $2}' /tmp/vols.txt)
do
    i=$(( i + 1 ))
    ec2-attach-volume $vol \
        -O $akey \
        -W $skey \
        -i $instance \
        --region $region \
        -d /dev/sdh${i}
done > /tmp/attach.txt
```

3. Mark the EBS volumes as physical volumes:

```
for i in $(find /dev/xvdh*)
do
    pvcreate $i
done
```

4. Create a volume group:

```
vgcreate vg /dev/xvdh*
```

5. Create a logical volume:

```
lvcreate -l100%VG -n lv vg
```

6. Create the file system:

```
mkfs -t xfs /dev/vg/lv
```

8. Mount the file system:

```
mount /dev/vg/lv /data/data_large
```

9. Create mount point and mount the device:

```
mkdir /data/data_large
mount /dev/md0 /data/data_large
```

10. Add new mountpoint to /etc/exports:

```
for x in $(qconf -sh | tail -n +2)
do
    echo '/data/data_large' ${x}' (async,no_root_squash,no_subtree_check,rw)' >> /etc/exports
done
```

11. Reload /etc/exports:

```
exportfs -a
```

12. Mount the new folder on all nodes:

```
for x in $(qconf -sh | tail -n +2)
do
    ssh $x 'mkdir /data/data_large'
    ssh $x 'mount -t nfs master:/data/data_large /data/data_large'
done
```

### How to increase volume size?

1. Create and attach EBS volumes as described in steps 1. & 2. and then create the additional physical volumes:

```
for i in $(cat /tmp/attach.txt | cut -f 4 | sed 's/[^0-9]*//g')
do
    pvcreate /dev/xvdh${i}
done
```

2. Add new volumes to the volume group:

```
for i in $(cat /tmp/attach.txt | cut -f 4 | sed 's/[^0-9]*//g')
do
    vgextend vg /dev/xvdh${i}
done

lvextend -l100%VG /dev/mapper/vg-lv
```

3. Grow the file system to the new size:

```
xfs_growfs /data/data_large
```

## 4.10 Add storage > 1TB to the cluster using RAID 0 (for advanced users)

1. Within StarCluster create x new volumes by running:

```

nvolumes=2 #number of volumes
vsize=1000 #in gb
instance=`curl -s http://169.254.169.254/latest/meta-data/instance-id`
akey=<AWS KEY>
skey=<AWS SECRET KEY>
region=us-west-1
zone=us-west-1a

for x in $(seq 1 $nvolumes)
do
    ec2-create-volume \
        --aws-access-key $akey \
        --aws-secret-key $skey \
        --size $vsize \
        --region $region \
        --availability-zone $zone
done > /tmp/vols.txt

```

## 2. Attach the volumes to the head node:

```

i=0
for vol in $(awk '{print $2}' /tmp/vols.txt)
do
    i=$(( i + 1 ))
    ec2-attach-volume $vol \
        -O $akey \
        -W $skey \
        -i $instance \
        --region $region \
        -d /dev/sdh${i}
done

```

## 3. Create a raid 0 volume:

```
mdadm --create -l 0 -n $nvolumes /dev/md0 /dev/xvdh*
```

## 4. Create a file system:

```
mkfs -t ext4 /dev/md0
```

## 5. Create mount point and mount the device:

```
mkdir /data/data_large
mount /dev/md0 /data/data_large
```

## 6. Add new mountpoint to /etc/exports:

```

for x in $(qconf -sh | tail -n +2)
do
    echo '/data/data_large' ${x}' (async,no_root_squash,no_subtree_check,rw)' >> /etc/exports
done

```

## 7. Reload /etc/exports:

```
exportfs -a
```

## 8. Mount the new folder on all nodes:

```

for x in $(qconf -sh | tail -n +2)
do

```

```
ssh $x 'mkdir /data/data_large'  
ssh $x 'mount -t nfs master:/data/data_large /data/data_large'  
done
```

## 4.11 Backing up your data to S3

1. Run:

```
s3cmd --configure
```

and follow the instructions

2. Create a S3 bucket:

```
s3cmd mb s3://backup
```

3. Copy data to the bucket:

```
s3cmd put -r /data/results s3://backup
```

More info on s3cmd here: <https://github.com/s3tools/s3cmd>

---

## Omics Pipe Tutorial – Installation

---

### 5.1 Installation

*Installation instructions*

### 5.2 Before Running Omics Pipe: Configuring Parameters File

---

**Note:** Before running `omics_pipe`, you must configure the parameters file, which is a YAML document. Follow the instructions here: *Configuring the parameters file*

---

### 5.3 Running Omics Pipe

When you are ready to run `omics pipe`, simply type the command:

```
omics_pipe RNAseq_count_based /path/to/parameter_file.yaml
```

To run the basic `RNAseq_count_based` pipeline with your parameter file. Additional usage instructions below and are available by typing `omics_pipe -h`:

```
omics_pipe [-h] [--custom_script_path CUSTOM_SCRIPT_PATH]
           [--custom_script_name CUSTOM_SCRIPT_NAME]
           [--compression {gzip, bzip}]
           {RNAseq_Tuxedo, RNAseq_count_based, RNAseq_cancer_report, RNAseq_TCGA, RNAseq_TCGA_counts,
           Tumorseq_MUTECT, miRNAseq_count_based, miRNAseq_tuxedo, WES_GATK, WGS_GATK,
           parameter_file
```

If your `.fastq` files are compressed, please use the compression option and indicate the type of compression used for your files. Currently supported compression types are `gzip` and `bzip`.

### 5.4 Running Omics Pipe with the Test Data and Parameters

To run Omics Pipe with the test parameter files and data, type the commands below to run each pipeline.

---

**Note:** Replace the `~` with the path to your Omics Pipe installation.

---

**RNA-seq (Tuxedo):**

```
omics_pipe RNAseq_Tuxedo ~/tests/test_params_RNAseq_Tuxedo.yaml
```

**RNA-seq(Anders 2013):**

```
omics_pipe RNAseq_count_based ~/tests/test_params_RNAseq_counts.yaml
```

**Whole Exome Sequencing (GATK):**

```
omics_pipe WES_GATK ~/tests/test_params_WES_GATK.yaml
```

**Whole Genome Sequencing (GATK):**

```
omics_pipe WGS_GATK ~/tests/test_params_WGS_GATK.yaml
```

**Whole Genome Sequencing (MUTECT):**

```
omics_pipe Tumorseq_MUTECT ~/tests/test_params_MUTECT.yaml
```

**ChIP-seq (MACS):**

```
omics_pipe ChIPseq_MACS ~/tests/test_params_MACS.yaml
```

**ChIP-seq (HOMER):**

```
omics_pipe ChIPseq_HOMER ~/tests/test_params_HOMER.yaml
```

**Breast Cancer Personalized Genomics Report- RNAseq:**

```
omics_pipe RNAseq_cancer_report ~/tests/test_params_RNAseq_cancer.yaml
```

**TCGA Reanalysis Pipeline - RNAseq:**

```
omics_pipe RNAseq_TCGA ~/tests/test_params_RNAseq_TCGA.yaml
```

**TCGA Reanalysis Pipeline - RNAseq Counts:**

```
omics_pipe RNAseq_TCGA_counts ~/tests/test_params_RNAseq_TCGA_counts.yaml
```

**miRNAseq Counts (Anders 2013):**

```
omics_pipe miRNAseq_count_based ~/tests/test_params_miRNAseq_counts.yaml
```

**miRNAseq (Tuxedo):**

```
omics_pipe miRNAseq_tuxedo ~/tests/test_params_miRNAseq_Tuxedo.yaml
```



---

## Omics Pipe Tutorial – Configuring the Parameter File

---

Before running Omics Pipe, you must configure the parameters file, which is a [YAML](#) document. Example parameters files are located within the `omics_pipe/test` folder for each pipeline. Copy one of these parameters files into your working directory, and edit the parameters to work with your sample names, directory structure, software options and software versions. Make sure to keep the formatting and parameter names exactly the same as in the example parameters files.

---

**Note:** Make sure to follow the YAML format exactly. Ensure that there is only one space after each colon.

---

**Note:** For parameters in quotes in the test parameters file, please make sure to keep them in quotes in your custom parameter file.

---

The `STEP` parameter should be the function name of the last step in the pipeline that you want to run (e.g. `run_tophat`). To run the pre-installed pipelines all the way through, this should be `“last_function.”`

**Warning:** Do not change the `STEPS` or `STEPS_DE` parameters for a pre-installed pipeline.

---

**Note:** Fastq files: paired end: 2 files, `“Name_1.fastq”` and `“Name_2.fastq”` representing read 1 and read 2. Have all fastq files in same raw data folder

---

### 6.1 Example Omics Pipe Parameter File

test\_params.yaml in `omics_pipe/tests`:

```
SAMPLE_LIST: [test1, test2, test3]
```

```
STEP: last_function
```

```
STEPS: [fastqc, star, htseq, last_function]
```

```
RAW_DATA_DIR: /gpfs/home/kfisch/scripts/omics_pipeline-devel/tests
```

```
FLAG_PATH: /gpfs/home/kfisch/scripts/omics_pipeline-devel/tests/test_run/flags
```

```
HTSEQ_RESULTS: /gpfs/home/kfisch/scripts/omics_pipeline-devel/tests/test_run/counts
```

```
LOG_PATH: /gpfs/home/kfisch/scripts/omics_pipeline-devel/tests/test_run/logs
```

```
QC_PATH: /gpfs/home/kfisch/scripts/omics_pipeline-devel/tests/test_run
RESULTS_PATH: /gpfs/home/kfisch/test
STAR_RESULTS: /gpfs/home/kfisch/scripts/omics_pipeline-devel/tests/test_run/star
WORKING_DIR: /gpfs/home/kfisch/scripts/omics_pipeline-devel/omics_pipe/scripts
REPORT_RESULTS: /gpfs/home/kfisch/scripts/omics_pipeline-devel/tests/test_run
ENDS: SE
FASTQC_VERSION: '0.10.1'
GENOME: /gpfs/group/databases/Homo_sapiens/UCSC/hg19/Sequence/WholeGenomeFasta/genome.fa
HTSEQ_OPTIONS: -m intersection-nonempty -s no -t exon
PIPE_MULTIPROCESS: 100
PIPE_REBUILD: 'True'
PIPE_VERBOSE: 5
REF_GENES: /gpfs/group/databases/Homo_sapiens/UCSC/hg19/Annotation/Genes/genes.gtf
RESULTS_EMAIL: kfisch@scripps.edu
STAR_INDEX: /gpfs/group/databases/Homo_sapiens/UCSC/hg19/star_genome
STAR_OPTIONS: --readFilesCommand cat --runThreadN 8 --outSAMstrandField intronMotif --outFilterIntron
STAR_VERSION: '2.3.0'
TEMP_DIR: /scratch/kfisch
QUEUE: workq
USERNAME: kfisch
DRMAA_PATH: /opt/applications/pbs-drmaa/current/gnu/lib/libdrmaa.so
DPS_VERSION: '1.3.1111'
BAM_FILE_NAME: Aligned.out.bam
PARAMS_FILE: '/gpfs/home/kfisch/scripts/omics_pipeline-devel/tests/test_params_RNAseq_counts.yaml'
DESEQ_META: /gpfs/home/kfisch/scripts/omics_pipeline-devel/tests/counts_meta.csv
DESIGN: '~ condition'
PVAL: '0.05'
DESEQ_RESULTS: /gpfs/home/kfisch/scripts/omics_pipeline-devel/tests/test_run/DESEQ
SUMATRA_DB_PATH: /gpfs/home/kfisch/sumatra
```

```
SUMATRA_RUN_NAME: test_counts_sumatra_project
REPOSITORY: https://kfish@bitbucket.org/sulab/omics_pipe
HG_USERNAME: Kathleen Fisch <kfish@scripps.edu>
```

## 6.2 Explanation of Variables in Omics Pipe Parameter File

These parameters are for the RNAseq Count Based Pipeline. Parameters vary by pipeline. See examples in the tests/ folder.

test\_params.yaml in omics\_pipe/tests:

```
#sample names ie "Name" for paired and single end reads. So, "Name" for paired-end would expect two
SAMPLE_LIST: [test1, test2, test3]

#Function to be run within pipeline. If you want to run the whole pipeline, leave this as last_function
STEP: last_function

#All steps within the pipeline. Do not change this parameter for pre-installed pipelines. If you create
STEPS: [fastqc, star, htseq, last_function]

#Directory where your raw .fastq files are located.
RAW_DATA_DIR: /gpfs/home/kfish/scripts/omics_pipeline-devel/tests

#Directory where you would like to have the flag files created. Flag files are empty files that indicate
FLAG_PATH: /gpfs/home/kfish/scripts/omics_pipeline-devel/tests/test_run/flags

#Directory for HTSEQ results
HTSEQ_RESULTS: /gpfs/home/kfish/scripts/omics_pipeline-devel/tests/test_run/counts

#Directory where log files will be written
LOG_PATH: /gpfs/home/kfish/scripts/omics_pipeline-devel/tests/test_run/logs

#Directory for QC results
QC_PATH: /gpfs/home/kfish/scripts/omics_pipeline-devel/tests/test_run

#Upper level results directory. Sumatra will check all subfolders of this directory for new files to
RESULTS_PATH: /gpfs/home/kfish/test

#Directory where STAR results will be written
STAR_RESULTS: /gpfs/home/kfish/scripts/omics_pipeline-devel/tests/test_run/star

#Where omics_pipe is installed, this path will be pointing to ~/omics_pipe/scripts.
WORKING_DIR: /gpfs/home/kfish/scripts/omics_pipeline-devel/omics_pipe/scripts

#Directory for the summary report
REPORT_RESULTS: /gpfs/home/kfish/scripts/omics_pipeline-devel/tests/test_run

#SE is single end, PE is paired-end sequencing reads
ENDS: SE

#Version number of FASTQC
FASTQC_VERSION: '0.10.1'

#Full path to Genome fasta file
```

```
GENOME: /gpfs/group/databases/Homo_sapiens/UCSC/hg19/Sequence/WholeGenomeFasta/genome.fa

#Options for HTSEQ
HTSEQ_OPTIONS: -m intersection-nonempty -s no -t exon

#Number of multiple processes you want Ruffus to spawn at once
PIPE_MULTIPROCESS: 100

#Ruffus parameter. No need to change.
PIPE_REBUILD: 'True'

#Ruffus parameter. No need to change.
PIPE_VERBOSE: 5

#Full path to reference gene annotations
REF_GENES: /gpfs/group/databases/Homo_sapiens/UCSC/hg19/Annotation/Genes/genes.gtf

#Your email.
RESULTS_EMAIL: kfisch@scripps.edu

#Directory pointing to STAR_INDEX (you may have to create this)
STAR_INDEX: /gpfs/group/databases/Homo_sapiens/UCSC/hg19/star_genome

#Options for STAR
STAR_OPTIONS: --readFilesCommand cat --runThreadN 8 --outSAMstrandField intronMotif --outFilterIntron

#Version number of STAR
STAR_VERSION: '2.3.0'

#Path to temporary directory
TEMP_DIR: /scratch/kfisch

#Name of the queue on your local cluster you wish to use
QUEUE: workq

#Username for local cluster
USERNAME: kfisch

#Path to your local cluster installation of DRMAA (ask your sys admin for this)
DRMAA_PATH: /opt/applications/pbs-drmaa/current/gnu/lib/libdrmaa.so

#Version number of Drug Pair Seeker
DPS_VERSION: '1.3.1111'

#Name of create Bam file. Will be Aligned.out.bam if you are using STAR and accepted_hits.bam if you
BAM_FILE_NAME: Aligned.out.bam

#Full path to your parameter file. Make sure to include the single quotes.
PARAMS_FILE: '/gpfs/home/kfisch/scripts/omics_pipeline-devel/tests/test_params_RNAseq_counts.yaml'

#Location of the meta data csv file for DESEQ. See tests/counts_meta.csv for an example.
DESEQ_META: /gpfs/home/kfisch/scripts/omics_pipeline-devel/tests/counts_meta.csv

#Design for DESEQ differential expression. Leave as is if you use the exact design as in the counts_r
DESIGN: '~ condition'

#P-value threshold
PVAL: '0.05'
```

```
#Directory for DESEQ results
DESEQ_RESULTS: /gpfs/home/kfisch/scripts/omics_pipeline-devel/tests/test_run/DESEQ

#Directory where you want to store your Sumatra database. Once you run this once, you do not have to
SUMATRA_DB_PATH: /gpfs/home/kfisch/sumatra

#Name of your project. You do not need to change this for subsequent runs of the pipeline, but you can
SUMATRA_RUN_NAME: test_counts_sumatra_project

#Location of omics pipe repository (you can leave this)
REPOSITORY: https://kfisch@bitbucket.org/sulab/omics_pipe

#Your Mercurial username
HG_USERNAME: Kathleen Fisch <kfisch@scripps.edu>
```



---

## Omics Pipe Tutorial – Creating a Custom Pipeline Script

---

A pipeline script is a .py file that has the steps in the pipeline that you want to run in your analysis. To create a custom pipeline, you will create a new Python script (\*.py) file and place it in your working directory (or wherever you want). The available analysis steps built-in to omics\_pipe are available in the (List of currently available omics\_pipe analysis steps).

You may add new modules directly to the module directory (see *Adding Custom Modules*), and they will become available steps that you can use in your custom pipeline.

**There are three steps to creating a custom pipeline:**

1. Designing the structure of your pipeline
2. Creating the script
3. Updating your parameters file

The section below details each of these steps.

### 7.1 Designing the structure of the pipeline

Omics\_pipe depends upon the pipelining module Ruffus to handle the automation. Please read the documentation at the Ruffus website <http://www.ruffus.org.uk/> for more information. To design your pipeline, you need to decide

- What analysis modules you want to run,
- What order you want the analysis modules to run in,
- Which, if any, of the analysis modules depend upon the results from another analysis module.

For example, we will create a custom pipeline that runs fastqc, star and htseq (depends on output from star).

### 7.2 Creating the script

To create the script, create a text file name custom\_script.py (or whatever name you choose). At the top of the file, cut and copy this text:

```
#!/usr/bin/env python
from ruffus import *
import sys
import os
import time
import datetime
```

```
import drmaa
from omics_pipe.utils import *
from omics_pipe.parameters.default_parameters import default_parameters
p = Bunch(default_parameters)
os.chdir(p.WORKING_DIR)
now = datetime.datetime.now()
date = now.strftime("%Y-%m-%d %H:%M")
print p
for step in p.STEPS:
    vars()['inputList_' + step] = []
    for sample in p.SAMPLE_LIST:
        vars()['inputList_' + step].append([sample, "%s/%s_%s_completed.flag" % (p.FLAG_PATH,
```

After this has been completed, you will need to import each of the analysis modules that you will use in your pipeline. For each analysis module, write the line below (fill in analysis\_name with the name of the analysis module):

```
from omics_pipe.modules.analysis_name import analysis_name
```

In our example, you will have three lines (see below):

```
from omics_pipe.modules.fastqc import fastqc
from omics_pipe.modules.star import star
from omics_pipe.modules.htseq import htseq
```

Now you are ready to write the functions to run each of these steps in the analysis. For each step in our analysis pipeline, we will need to write a function. You can cut and copy these from the pre-packaged analysis pipeline scripts (or eventually a function reference) or create them. Each function needs to have two decorators from Ruffus: `@parallel(inputList_analysis_name)` to specify that the pipeline should run in parallel for more than one sample and `@check_if_uptodate(check_file_exists)` to call a function to check if that step in the pipeline is up to date. Name each function with the name of the analysis prefixed by “run\_.” The function input should always be (sample, analysis\_name\_flag). Within the function, you will call the analysis module that you loaded above. If you want an analysis module to run only after a module it depends upon finishes, you must add the `@follows()` Ruffus decorator before the function, with the name of the step that it depends upon. For example, if htseq needs to run after star, you would put `@follows(run_star)` above the `run_htseq` function. If you have steps that do not have functions that are dependent upon them, you can create a more complex pipeline structure by creating a “Last Function” that ties together all steps of your pipeline. The last function below is an example of such a function, and it also produces a PDF diagram of your pipeline when it completes. The functions for our example are below.

```
@parallel(inputList_fastqc)
@check_if_uptodate(check_file_exists)
def run_fastqc(sample, fastqc_flag):
    fastqc(sample, fastqc_flag)
    return

@parallel(inputList_star)
@check_if_uptodate(check_file_exists)
def run_star(sample, star_flag):
    star(sample, star_flag)
    return

@parallel(inputList_htseq)
@check_if_uptodate(check_file_exists)
@follows(run_star)
def run_htseq(sample, htseq_flag):
    htseq(sample, htseq_flag)
    return
```



```

@parallel(inputList_last_function)
@check_if_uptodate(check_file_exists)
@follows(run_fastqc, run_htseq)
def last_function(sample, last_function_flag):
    print "PIPELINE HAS FINISHED SUCCESSFULLY!!! YAY!"
    pipeline_graph_output = p.FLAG_PATH + "/pipeline_" + sample + "_" + str(date) + ".pdf"
    pipeline_printout_graph (pipeline_graph_output, 'pdf', step, no_key_legend=False)
    stage = "last_function"
    flag_file = "%s/%s_%s_completed.flag" % (p.FLAG_PATH, stage, sample)
    open(flag_file, 'w').close()
    return

```

Once you have created all of the functions for each step of your pipeline, cut and copy the code below to the bottom of your script:

```

if __name__ == '__main__':

    pipeline_run(p.STEP, multiprocessing = p.PIPE_MULTIPROCESS, verbose = p.PIPE_VERBOSE, gnu_make_r

```

At this point, please save your script and move on to step 3.

## 7.3 Updating your parameters file

In order for your script to run successfully, you need to configure your parameter file so that each analysis module has the necessary parameters to execute successfully. The full list of parameters for all modules in the current version of omics<sub>pipe</sub> are located in the omics<sub>pipe</sub>/parameters/default\_parameters.py file (and eventually organized somewhere). You can view the list of necessary parameters for each analysis module by importing the analysis module into an interactive python session (from omics<sub>pipe</sub>.modules.analysis\_module import analysis\_module) and typing analysis\_module.\_\_doc\_\_. The parameters necessary for that analysis module will be listed under “parameters from parameters file.” These parameters must be put into your parameters.yaml file and spelled exactly as shown (including all caps). Below is the list of parameters that are necessary to run omics<sub>pipe</sub> in addition to the module specific parameters.

```

SAMPLE_LIST: [test, test1]
STEP: run_last_function
STEPS: [fastqc, star, htseq, last_function]
RAW_DATA_DIR: /gpfs/group/sanford/patient/SSKT/test_patient/RNA/RNA_seq/data
FLAG_PATH: /gpfs/group/sanford/patient/SSKT/test_patient/RNA/RNA_seq/logs/flags
LOG_PATH: /gpfs/group/sanford/patient/SSKT/test_patient/RNA/RNA_seq/logs
WORKING_DIR: /gpfs/home/kfisch/virt_env/virt2/lib/python2.6/site-packages/omics_pipe-1.0.7-py2.6.egg
ENDS: PE
PIPE_MULTIPROCESS: 100
PIPE_REBUILD: 'True'
PIPE_VERBOSE: 5
RESULTS_EMAIL: kfisch@scripps.edu
TEMP_DIR: /scratch/kfisch
DPS_VERSION: '1.3.1111'
QUEUE: bigmem
PARAMS_FILE: /gpfs/home/kfisch/omics_pipe_docs/test_params.yaml
USERNAME: kfisch
GENOME: /gpfs/group/databases/Homo_sapiens/UCSC/hg19/Sequence/WholeGenomeFasta/genome.fa
CHROM: /gpfs/group/databases/Homo_sapiens/UCSC/hg19/Sequence/Chromosomes
REF_GENES: /gpfs/group/databases/Homo_sapiens/UCSC/hg19/Annotation/Genes/genes.gtf
STAR_INDEX: /gpfs/group/databases/Homo_sapiens/UCSC/hg19/star_genome
BOWTIE_INDEX: /gpfs/group/databases/Homo_sapiens/UCSC/hg19/Sequence/Bowtie2Index/genome

```

Once you have all of the necessary parameters in your parameters.yaml file, for your custom script you will need to change the STEP and STEPS parameters. In the STEP parameter, you will write the name of the last function in your pipeline that you want to run, which should be configured so that it captures all steps in the pipeline (as in the example above). Make sure to put **run\_** in front of this, since you are calling the function, not the analysis module. In order for omics\_pipe to know what steps you have in your pipeline, you need to list each analysis module name in the STEPS parameter separated with commas (without **run\_** in the prefix). You are now ready to run your custom script.

Running omics\_pipe with a custom pipeline script When you call the omics\_pipe function, you will specify the path to your custom script using the command

```
omics_pipe custom --custom_script_path ~/path/to/the/script -custom_script_name customscript /path/t
```

This will automatically load your custom script and run through the steps in your pipeline using the default modules available in omics\_pipe.

---

## Omics Pipe Tutorial – Adding a New Module (Tool)

---

Users can easily create new analysis modules for use within `omics_pipe`. The user has two options for creating new analysis modules: - Adding analysis modules directly within the `omics_pipe/scripts` installation directory - Creating a new working directory where all analysis modules scripts are located (this can be changed in the parameters file by changing the `WORKING_DIR` parameter to the desired location). If you want to use option 2, in order to use pre-installed analysis modules, for the time being you must copy these analysis modules to your new working directory. If you choose option 1, you can simply add additional analysis modules and they will be accessible along with the pre-installed analysis modules.

To create a new analysis module, you need to perform four steps: 1. Create a Bash script with the command to be sent to the cluster 2. Create a Python module that calls the Bash script 3. Add your module to your custom pipeline 4. Add new module parameters to parameters file

The section below details each of these steps.

### 8.1 1. Create a Bash script

The first step in creating your custom module is to create the Bash script with the command you would like to run. If you are unsure how to write a Bash script, you can look at the examples in `omics_pipe/scripts` or work through this tutorial (<http://tldp.org/HOWTO/Bash-Prog-Intro-HOWTO.html>). In many cases, this will be a simple script with a one line command to call the analysis program. You should name your script something that will be easily identifiable and it should have the suffix `.sh` (e.g. `analysis_script.sh`). At the beginning of your analysis script, you should put the following lines:

```
#!/bin/bash
set -x
#Source modules for current shell
source $MODULESHOME/init/bash
#Make output directory if it doesn't exist
mkdir -p ${variable} #RESULTS_DIR
#Move tmp dir to scratch
export TMPDIR=${variable} #TEMP_DIR
#Load specified software version
module load fastqc/${variable} #VERSION
```

The `${variable}` will be changed to `${number}` (e.g. `$1`) based on the location of the variable in the input script (more on this below). These settings are assuming you are working on a cluster with a modular structure. If not, “module load” may not be appropriate to load the software, so please ask your system administrator to provide assistance with this if your cluster has a different system. After you specify the software and other configuration variables, you can write the commands for the software you would like to use. When you are finished with the commands, exit the script with ‘`exit 0.`’ An example script for running the software program FASTQC is below.

```
#Runs fastqc with $1=SAMPLE, $2=RAW_DATA_DIR, $3=QC_PATH
fastqc -o $3 $2/$1.fastq

exit 0
```

Substitute all variables that you would like to change from the parameter file with a variable notation, in the form of \$1, \$2, \$3, etc for the first, second, third, etc input parameter that will be passed to the script. Once you have appropriately parameterized the script, save the script either in your working directory (along with all the other scripts you will need, possibly copied from omics<sub>pipe</sub>/scripts) or in the omics<sub>pipe</sub>/scripts directory.

## 8.2 2. Create a Python module

Now that you have created your custom script, you can create the Python module that will handle that script and schedule a job on the compute cluster using DRMAA (<https://code.google.com/p/drmaa-python/wiki/Tutorial>). You should name the Python module the same name as your custom analysis module, but with the extension .py. In this example, your Python module would be named analysis\_script.py and the function within it would also be called analysis\_script. Save your custom Python module within the same directory as your custom pipeline script. At the top of your Python module, cut and copy the text below.

```
#!/usr/bin/env python

import drmaa
from omics_pipe.parameters.default_parameters import default_parameters
from omics_pipe.utils import *
p = Bunch(default_parameters)
```

You will then write a simple Python function that take the form of the function below. You can directly copy this function and then change the necessary names/parameters to fit your custom analysis. ::

```
def fastqc(sample, fastqc_flag):
    '''QC check of raw .fastq files using FASTQC
        input: .fastq file
        output: folder and zipped folder containing html, txt and image files
        citation: Babraham Bioinformatics
        link: http://www.bioinformatics.babraham.ac.uk/projects/fastqc/
        parameters from parameters file: RAW_DATA_DIR, QC_PATH, FASTQC_VERSION'''
    spawn_job(jobname = 'fastqc', SAMPLE = sample, LOG_PATH = p.LOG_PATH, RESULTS_EMAIL = p.RESULTS_EMAIL,
              job_status(jobname = 'fastqc', resultspath = p.QC_PATH, SAMPLE = sample, outputfilename = p.OUTPUT_FILENAME))
    return
```

Name your function the same as the names of both the Bash and Python scripts you just created for consistency. In our example, the first line would look like: “def analysis\_script(sample, analysis\_script\_flag):”. As you can see, I changed the name of the function as well as the name of the flag input file. The document string should be changed to describe what your analysis module does, what type of input file it takes, a citation and link to the tool that you are calling, as well as the parameters that are needed in the parameters file that will be passed to the Bash script that you created. After you are done documenting your function, you will change a few items within the spawn\_job and job\_status functions that are called from the omics<sub>pipe</sub>.utils module. In the spawn\_job function, you should change the job name to match the name of your function, you can customize the resources your job will request from the compute cluster, you will need to change the name of the script to match that of the Bash script that you just created, and then you will change the parameters listed in the variable “args\_list.” The variable “sample” is lower case because it is passed to this function from omics<sub>pipe</sub>, but input parameters coming from the parameters file must be prefixed with “p.” List the parameters that you need to feed into your custom analysis script in the order that you numbered them in the Bash script. In the example above, \$1 corresponds to ‘sample’ \$2 corresponds to p.RAW\_DATA\_DIR, etc. Once you have the spawn\_job function updated, you will update the job\_status function with the job name, results path and a name of an output file that will be produced from your Bash script. This can be any file that is created. This

function will check that this file exists in the specified results directory, check that its size is greater than zero, and then it will create a flag file if it exists. Once you complete this, you are finished creating your custom Python module.

### **8.3 3. Add custom Python module to your custom pipeline**

In order to use your custom analysis module, you will need to create a custom pipeline with your custom analysis module included as a step in the pipeline. For a tutorial on how to create a custom pipeline, see Section “Creating a Custom Pipeline Script.” Once you have a custom pipeline script, please make sure your custom analysis module and custom pipeline script are in the same directory.

### **8.4 4. Add new parameters to parameters file**

The final step in custom analysis module creation is to add the parameters necessary for your custom analysis module to run into the parameters file. Simply add the parameters to your parameters script, save it, and then run your custom pipeline.



---

## Omics Pipe Available Pipelines

---

### 9.1 RNA-seq (Tuxedo)

*RNA-seq Tuxedo Modules* Modules included in the Tuxedo RNA-seq pipeline.

- FASTQC
- TopHat
- Cufflinks
- Cuffmerge
- Cuffmergetocompare
- Cuffdiff
- R Summary Report - CummeRbund

### 9.2 RNA-seq(Anders 2013)

*RNA-seq Count Based Modules* Modules included in the count-based RNA-seq pipeline.

- FASTQC
- STAR
- HTSEQ
- R Summary Report - DESEQ2

### 9.3 Whole Exome Sequencing (GATK)

*Whole Genome and Whole Exome Sequencing Modules* Modules included in the whole exome sequencing pipeline.

- FASTQC
- BWA-MEM
- PICARD Mark Duplicates
- GATK Preprocessing
- GATK Variant Discovery

- GATK Variant Filtering

## 9.4 Whole Genome Sequencing (GATK)

*Whole Genome and Whole Exome Sequencing Modules* Modules included in the whole genome sequencing pipeline.

- FASTQC
- BWA-MEM
- PICARD Mark Duplicates
- GATK Preprocessing
- GATK Variant Discovery
- GATK Variant Filtering

## 9.5 Whole Genome Sequencing (MUTECT)

*Whole Genome Sequencing (MUTECT)* Modules included in the cancer (paired tumor/normal) whole genome sequencing pipeline.

- FASTQC
- BWA-MEM
- MUTECT

## 9.6 ChIP-seq (MACS)

*ChIP-seq Modules – MACS* Modules included in the ChIP-seq MACS pipeline.

- FASTQC
- Homer ChIP Trim
- Bowtie
- MACS

## 9.7 ChIP-seq (HOMER)

*ChIP-seq Modules – HOMER* Modules included in the ChIP-seq HOMER pipeline.

- FASTQC
- Homer ChIP Trim
- Bowtie
- Homer Read Density
- Homer Peaks
- Homer Peak Track



- Homer Annotate Peaks
- Homer Find Motifs

## 9.8 Breast Cancer Personalized Genomics Report- RNAseq

*Breast Cancer Personalized Genomics Report- RNAseq* Modules included in the RNAseq Cancer pipeline.

- FASTQC
- STAR
- RSEQC
- Fusion Catcher
- BWA/SNPiR
- Filter Variants
- HTseq
- Intogen
- OncoRep Cancer Report

## 9.9 TCGA Reanalysis Pipeline - RNAseq

*TCGA Reanalysis Pipeline - RNAseq* Modules included in the RNAseq Cancer pipeline.

- TCGA Download (GeneTorrent)
- FASTQC
- STAR
- RSEQC
- Fusion Catcher
- BWA/SNPiR
- Filter Variants
- HTseq
- Intogen
- OncoRep Cancer Report

## 9.10 TCGA Reanalysis Pipeline - RNAseq Counts

*RNA-seq Count Based Modules- TCGA* Modules included in the RNAseq counts pipeline for TCGA reanalysis.

- TCGA Download (GeneTorrent)
- FASTQC
- STAR
- HTSEQ

- Report

## 9.11 miRNAseq Counts (Anders 2013)

*miRNA-seq Count Based Modules* Modules included in the miRNAseq counts pipeline.

- Cutadapt
- FASTQ Length Filter
- FASTQC
- STAR
- HTSEQ
- Report

## 9.12 miRNAseq (Tuxedo)

*miRNA-seq Tuxedo Modules* Modules included in the miRNAseq Tuxedo pipeline.

- Cutadapt
- FASTQ Length Filter
- TopHat
- Cufflinks
- Cuffmerge
- Cuffmergetocompare
- Cuffdiff
- R Summary Report

## 9.13 All Available Modules

*All Available Modules*

---

## Reference Databases Needed

---

To run the pipelines, you will need to have reference databases installed on your cluster. If you are using the AWS installation, these databases are provided for you. If you need to install your references, please install the ones below. Omics Pipe is compatible with all species genome files. Examples below are for hg19, but you can substitute them for the equivalent files from other species.

### 10.1 All Pipelines

#### 10.1.1 Genome

- .fa file can be downloaded from: <http://cufflinks.cbcb.umd.edu/igenomes.html>

#### 10.1.2 Reference Annotation Files

You can use any reference annotations you would like, as long as they are GTF files.

Examples include:

- gencode.v18.annotation.gtf
- UCSC genes.gtf

### 10.2 Reference Data for Cancer Reporting Scripts (RNAseq cancer, TCGA pipelines)

For the cancer pipelines, please download the file from the link below, extract it and put the files in the respective directories. [Reporting\\_data](#)

#### 10.2.1 In your omics pipe installation directory under omics\_pipe/scripts/reporting/ref place the files.

- tega\_brca.R
- brca.txt <[http://sulab.scripps.edu/kfisch/omics\\_pipe/ref/brca.txt](http://sulab.scripps.edu/kfisch/omics_pipe/ref/brca.txt)>'\_

### 10.2.2 In your omics pipe installation directory under omics\_pipe/scripts/reporting/data place the remaining files.

- brca\_mol\_class/\*
- DoG/\*
- geneLists/\*
- SPIA
- deseq.tcga\_brca.Rdata
- loggeoameansBRCA.Rdata

## 10.3 References for Variants (RNA-seq cancer, RNA-seq cancer TCGA, WES and WGS pipelines)

For pipelines performing variant calling, please download the references below and put them in the specified directories.

### 10.3.1 You can put these files in any directory. You will point to their location in the parameters file.

- dbNSFP2.0.txt
- common\_no\_known\_medical\_impact\_00-latest.vcf

Available within the GATK resource bundle v.2.5:

- dbsnp\_137.hg19.vcf
- Mills\_and\_1000G\_gold\_standard.indels.hg19.vcf
- 1000G\_phase1.indels.hg19.vcf
- hapmap\_3.3.hg19.vcf
- 1000G\_omni2.5.hg19.vcf

### 10.3.2 In your omics pipe installation directory under omics\_pipe/scripts/reporting/ref place these files.

- cadd.tsv.gz from <http://cadd.gs.washington.edu/download>
- drugbank.tsv
- cosmic.tsv
- clinvar.txt

from PharmGKB:

- pharmgkbAllele.tsv
- pharmgkbRSID.csv

## 10.4 WES Pipeline

- truseq\_exome\_targeted\_regions.hg19.bed

## 10.5 ChIP-seq Pipelines

- hg19.chrom.sizes

## 10.6 SNPiR Pipelines (RNA-seq cancer and RNA-seq cancer TCGA pipelines)

- BWA Index
- RNA editing sites (Human\_AG\_all\_hg19.bed)
- RepeatMasker.bed
- anno\_combined\_sorted
- knowngene.bed



---

## Third Party Software Dependencies

---

Omics Pipe is dependent upon several third-party software packages. Before running Omics Pipe, please install all of the required tools for the pipeline you will be running (see below) as [Modules](#) on your local cluster. If you are running the AWS distribution, all third party software is already installed.

### 11.1 R Packages Needed

In R, you can cut and copy this to install all required packages:

```
install.packages(c("bibtex", "AnnotationDbi", "cluster", "cummeRbund", "data.table", "DBI", "DESeq2", "ggplot2", "graphite", "igraph", "KEGGREST", "knitr", "knitrBootstrap", "lattice", "locfit", "pamr", "plyr", "RColorBrewer", "Rcpp", "RcppArmadillo", "RCurl", "ReactomePA", "RefManageR", "RJSONIO", "RSQLite", "stringr", "survival", "XML", "xtable", "yaml"))
```

### 11.2 RNA-seq (Tuxedo)

FASTQC:	<a href="http://www.bioinformatics.babraham.ac.uk/projects/fastqc/">http://www.bioinformatics.babraham.ac.uk/projects/fastqc/</a>
TOPHAT:	<a href="http://tophat.cbcb.umd.edu/">http://tophat.cbcb.umd.edu/</a>
CUFFLINKS:	<a href="http://cufflinks.cbcb.umd.edu/">http://cufflinks.cbcb.umd.edu/</a>

### 11.3 RNA-seq (Anders 2013)

FASTQC:	<a href="http://www.bioinformatics.babraham.ac.uk/projects/fastqc/">http://www.bioinformatics.babraham.ac.uk/projects/fastqc/</a>
STAR:	<a href="http://code.google.com/p/rna-star/">http://code.google.com/p/rna-star/</a>
SAMTOOLS:	<a href="http://samtools.sourceforge.net/">http://samtools.sourceforge.net/</a>
HTSEQ:	<a href="http://www-huber.embl.de/users/anders/HTSeq/doc/index.html">http://www-huber.embl.de/users/anders/HTSeq/doc/index.html</a>

### 11.4 Whole Exome Sequencing (GATK)

FASTQC:	<a href="http://www.bioinformatics.babraham.ac.uk/projects/fastqc/">http://www.bioinformatics.babraham.ac.uk/projects/fastqc/</a>
BWA:	<a href="http://bio-bwa.sourceforge.net/">http://bio-bwa.sourceforge.net/</a>
PICARD:	<a href="http://picard.sourceforge.net/">http://picard.sourceforge.net/</a>
GATK:	<a href="https://www.broadinstitute.org/gatk/download">https://www.broadinstitute.org/gatk/download</a>

## 11.5 Whole Genome Sequencing (GATK)

FASTQC:	<a href="http://www.bioinformatics.babraham.ac.uk/projects/fastqc/">http://www.bioinformatics.babraham.ac.uk/projects/fastqc/</a>
BWA:	<a href="http://bio-bwa.sourceforge.net/">http://bio-bwa.sourceforge.net/</a>
PICARD:	<a href="http://picard.sourceforge.net/">http://picard.sourceforge.net/</a>
GATK:	<a href="https://www.broadinstitute.org/gatk/download">https://www.broadinstitute.org/gatk/download</a>

## 11.6 Whole Genome Sequencing (MUTECT)

FASTQC:	<a href="http://www.bioinformatics.babraham.ac.uk/projects/fastqc/">http://www.bioinformatics.babraham.ac.uk/projects/fastqc/</a>
BWA:	<a href="http://bio-bwa.sourceforge.net/">http://bio-bwa.sourceforge.net/</a>
MUTECT:	<a href="http://www.broadinstitute.org/cancer/cga/mutect">http://www.broadinstitute.org/cancer/cga/mutect</a>

## 11.7 ChIP-seq (MACS)

FASTQC:	<a href="http://www.bioinformatics.babraham.ac.uk/projects/fastqc/">http://www.bioinformatics.babraham.ac.uk/projects/fastqc/</a>
HOMER:	<a href="http://homer.salk.edu/homer/ngs/index.html">http://homer.salk.edu/homer/ngs/index.html</a>
BOWTIE:	<a href="http://bowtie-bio.sourceforge.net/index.shtml">http://bowtie-bio.sourceforge.net/index.shtml</a>
MACS:	<a href="http://liulab.dfci.harvard.edu/MACS/">http://liulab.dfci.harvard.edu/MACS/</a>
BEDTOOLS:	<a href="https://github.com/arq5x/bedtools2">https://github.com/arq5x/bedtools2</a>
SAMTOOLS:	<a href="http://samtools.sourceforge.net/">http://samtools.sourceforge.net/</a>

## 11.8 ChIP-seq (HOMER)

FASTQC:	<a href="http://www.bioinformatics.babraham.ac.uk/projects/fastqc/">http://www.bioinformatics.babraham.ac.uk/projects/fastqc/</a>
HOMER:	<a href="http://homer.salk.edu/homer/ngs/index.html">http://homer.salk.edu/homer/ngs/index.html</a>
BOWTIE:	<a href="http://bowtie-bio.sourceforge.net/index.shtml">http://bowtie-bio.sourceforge.net/index.shtml</a>
BEDTOOLS:	<a href="https://github.com/arq5x/bedtools2">https://github.com/arq5x/bedtools2</a>
SAMTOOLS:	<a href="http://samtools.sourceforge.net/">http://samtools.sourceforge.net/</a>



## 11.9 Breast Cancer Personalized Genomics Report- RNAseq

FASTQC:	<a href="http://www.bioinformatics.babraham.ac.uk/projects/fastqc/">http://www.bioinformatics.babraham.ac.uk/projects/fastqc/</a>
STAR:	<a href="http://code.google.com/p/rna-star/">http://code.google.com/p/rna-star/</a>
SAMTOOLS:	<a href="http://samtools.sourceforge.net/">http://samtools.sourceforge.net/</a>
HTSEQ:	<a href="http://www-huber.embl.de/users/anders/HTSeq/doc/index.html">http://www-huber.embl.de/users/anders/HTSeq/doc/index.html</a>
RSEQC:	<a href="http://rseqc.sourceforge.net/">http://rseqc.sourceforge.net/</a>
PICARD:	<a href="http://picard.sourceforge.net/">http://picard.sourceforge.net/</a>
GATK:	<a href="https://www.broadinstitute.org/gatk/download">https://www.broadinstitute.org/gatk/download</a>
FusionCatcher:	<a href="https://code.google.com/p/fusioncatcher/">https://code.google.com/p/fusioncatcher/</a>
Oncofuse:	<a href="http://www.unav.es/genetica/oncofuse.html">http://www.unav.es/genetica/oncofuse.html</a>
BWA:	<a href="http://bio-bwa.sourceforge.net/">http://bio-bwa.sourceforge.net/</a>
DNANEXUS SAMTOOLS:	<a href="https://github.com/dnanexus/samtools">https://github.com/dnanexus/samtools</a>
BEDTOOLS:	<a href="https://github.com/arq5x/bedtools2">https://github.com/arq5x/bedtools2</a>
BLAT:	<a href="https://genome.ucsc.edu/FAQ/FAQblat.html#blat3">https://genome.ucsc.edu/FAQ/FAQblat.html#blat3</a>
SNPiR:	<a href="http://lilab.stanford.edu/SNPiR/">http://lilab.stanford.edu/SNPiR/</a>
SNPEFF:	<a href="http://snpeff.sourceforge.net/">http://snpeff.sourceforge.net/</a>
SNPSIFT:	<a href="http://snpeff.sourceforge.net/SnpSift.html">http://snpeff.sourceforge.net/SnpSift.html</a>
VCFTOOLS:	<a href="http://vcftools.sourceforge.net/">http://vcftools.sourceforge.net/</a>

## 11.10 TCGA Reanalysis Pipeline - RNAseq

GeneTorrent:	<a href="https://cghub.ucsc.edu/docs/user/software.html">https://cghub.ucsc.edu/docs/user/software.html</a>
FASTQC:	<a href="http://www.bioinformatics.babraham.ac.uk/projects/fastqc/">http://www.bioinformatics.babraham.ac.uk/projects/fastqc/</a>
STAR:	<a href="http://code.google.com/p/rna-star/">http://code.google.com/p/rna-star/</a>
SAMTOOLS:	<a href="http://samtools.sourceforge.net/">http://samtools.sourceforge.net/</a>
HTSEQ:	<a href="http://www-huber.embl.de/users/anders/HTSeq/doc/index.html">http://www-huber.embl.de/users/anders/HTSeq/doc/index.html</a>
RSEQC:	<a href="http://rseqc.sourceforge.net/">http://rseqc.sourceforge.net/</a>
PICARD:	<a href="http://picard.sourceforge.net/">http://picard.sourceforge.net/</a>
GATK:	<a href="https://www.broadinstitute.org/gatk/download">https://www.broadinstitute.org/gatk/download</a>
FusionCatcher:	<a href="https://code.google.com/p/fusioncatcher/">https://code.google.com/p/fusioncatcher/</a>
Oncofuse:	<a href="http://www.unav.es/genetica/oncofuse.html">http://www.unav.es/genetica/oncofuse.html</a>
BWA:	<a href="http://bio-bwa.sourceforge.net/">http://bio-bwa.sourceforge.net/</a>
DNANEXUS SAMTOOLS:	<a href="https://github.com/dnanexus/samtools">https://github.com/dnanexus/samtools</a>
BEDTOOLS:	<a href="https://github.com/arq5x/bedtools2">https://github.com/arq5x/bedtools2</a>
BLAT:	<a href="https://genome.ucsc.edu/FAQ/FAQblat.html#blat3">https://genome.ucsc.edu/FAQ/FAQblat.html#blat3</a>
SNPiR:	<a href="http://lilab.stanford.edu/SNPiR/">http://lilab.stanford.edu/SNPiR/</a>
SNPEFF:	<a href="http://snpeff.sourceforge.net/">http://snpeff.sourceforge.net/</a>
SNPSIFT:	<a href="http://snpeff.sourceforge.net/SnpSift.html">http://snpeff.sourceforge.net/SnpSift.html</a>
VCFTOOLS:	<a href="http://vcftools.sourceforge.net/">http://vcftools.sourceforge.net/</a>

## 11.11 TCGA Reanalysis Pipeline - RNAseq Counts

GeneTorrent:	<a href="https://cghub.ucsc.edu/docs/user/software.html">https://cghub.ucsc.edu/docs/user/software.html</a>
FASTQC:	<a href="http://www.bioinformatics.babraham.ac.uk/projects/fastqc/">http://www.bioinformatics.babraham.ac.uk/projects/fastqc/</a>
STAR:	<a href="http://code.google.com/p/rna-star/">http://code.google.com/p/rna-star/</a>
SAMTOOLS:	<a href="http://samtools.sourceforge.net/">http://samtools.sourceforge.net/</a>
HTSEQ:	<a href="http://www-huber.embl.de/users/anders/HTSeq/doc/index.html">http://www-huber.embl.de/users/anders/HTSeq/doc/index.html</a>

## 11.12 miRNAseq Counts (Anders 2013)

CutAdapt:	<a href="http://code.google.com/p/cutadapt/">http://code.google.com/p/cutadapt/</a>
FASTQC:	<a href="http://www.bioinformatics.babraham.ac.uk/projects/fastqc/">http://www.bioinformatics.babraham.ac.uk/projects/fastqc/</a>
STAR:	<a href="http://code.google.com/p/rna-star/">http://code.google.com/p/rna-star/</a>
SAMTOOLS:	<a href="http://samtools.sourceforge.net/">http://samtools.sourceforge.net/</a>
HTSEQ:	<a href="http://www-huber.embl.de/users/anders/HTSeq/doc/index.html">http://www-huber.embl.de/users/anders/HTSeq/doc/index.html</a>

## 11.13 miRNAseq (Tuxedo)

CutAdapt:	<a href="http://code.google.com/p/cutadapt/">http://code.google.com/p/cutadapt/</a>
FASTQC:	<a href="http://www.bioinformatics.babraham.ac.uk/projects/fastqc/">http://www.bioinformatics.babraham.ac.uk/projects/fastqc/</a>
TOPHAT:	<a href="http://tophat.cbcb.umd.edu/">http://tophat.cbcb.umd.edu/</a>
CUFFLINKS:	<a href="http://cufflinks.cbcb.umd.edu/">http://cufflinks.cbcb.umd.edu/</a>

---

## System Requirements

---

If you are running Omics Pipe on a local high performance compute cluster, please ensure that you have the following minimum resource requirements.

- A minimum of 2 processors (nodes) with at least 32GB of memory (the more nodes you have available, the more the pipeline can be parallelized)
- Scratch space that is at least 3x the size of your expected results files.
- Storage space available for ~200 GB of reference-related data
- Storage space available for raw data files
- Storage space for results files (~10x that of raw data)



---

## RNA-seq Tuxedo Modules

---

Modules available in the RNA-seq Tuxedo Pipeline.

### 13.1 FASTQC

`omics_pipe.modules.fastqc.fastqc` (*sample*, *fastqc\_flag*)  
QC check of raw .fastq files using FASTQC.

**input:** .fastq file

**output:** folder and zipped folder containing html, txt and image files

**citation:** Babraham Bioinformatics

**link:** <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

**parameters from parameters file:** RAW\_DATA\_DIR:

QC\_PATH:

FASTQC\_VERSION:

COMPRESSION:

### 13.2 TopHat

`omics_pipe.modules.tophat.tophat` (*sample*, *tophat\_flag*)  
Runs TopHat to align .fastq files.

**input:** .fastq file

**output:** accepted\_hits.bam

**citation:** Trapnell C, Pachter L, Salzberg SL. TopHat: discovering splice junctions with RNA-Seq. Bioinformatics doi:10.1093/bioinformatics/btp120

**link:** <http://tophat.cbcb.umd.edu/>

**parameters from parameters file:** RAW\_DATA\_DIR:

REF\_GENES:

TOPHAT\_RESULTS:

BOWTIE\_INDEX:

TOPHAT\_VERSION:  
TOPHAT\_OPTIONS:  
BOWTIE\_VERSION:  
SAMTOOLS\_VERSION:

## 13.3 Cufflinks

omics\_pipe.modules.cufflinks.**cufflinks** (*sample, cufflinks\_flag*)

Runs cufflinks to assemble .bam files from TopHat.

**input:** accepted\_hits.bam

**output:** transcripts.gtf

**citation:** Trapnell C, et al. Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation Nature Biotechnology doi:10.1038/nbt.1621

**link:** <http://cufflinks.cbc.umd.edu/>

**parameters from parameters file:** TOPHAT\_RESULTS:

CUFFLINKS\_RESULTS:  
REF\_GENES:  
GENOME:  
CUFFLINKS\_OPTIONS:  
CUFFLINKS\_VERSION:

## 13.4 Cuffmerge

omics\_pipe.modules.cuffmerge.**cuffmerge** (*step, cuffmerge\_flag*)

Runs cuffmerge to merge .gtf files from Cufflinks.

**input:** assembly\_GTF\_list.txt

**output:** merged.gtf

**citation:** Trapnell C, et al. Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation Nature Biotechnology doi:10.1038/nbt.1621

**link:** <http://cufflinks.cbc.umd.edu/>

**parameters from parameters file:** CUFFMERGE\_RESULTS:

REF\_GENES:  
GENOME:  
CUFFMERGE\_OPTIONS:  
CUFFLINKS\_VERSION:

## 13.5 Cuffmergetocompare

omics\_pipe.modules.cuffmergetocompare.**cuffmergetocompare** (*step*, *cuffmergetocompare\_flag*)

Runs cuffcompare to annotate merged .gtf files from Cufflinks.

**input:** assembly\_GTF\_list.txt

**output:** merged.gtf

**citation:** Trapnell C, et al. Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation Nature Biotechnology doi:10.1038/nbt.1621

**link:** <http://cufflinks.cbcb.umd.edu/>

**parameters from parameters file:** CUFFMERGE\_RESULTS:

REF\_GENES:

GENOME:

CUFFMERGETOCOMPARE\_OPTIONS:

CUFFLINKS\_VERSION:

## 13.6 Cuffdiff

omics\_pipe.modules.cuffdiff.**cuffdiff** (*step*, *cuffdiff\_flag*)

Runs Cuffdiff to perform differential expression. Runs after Cufflinks. Part of Tuxedo Suite.

**input:** .bam files

**output:** differential expression results

**citation:** Trapnell C, et al. Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation Nature Biotechnology doi:10.1038/nbt.1621

**link:** <http://cufflinks.cbcb.umd.edu/>

**parameters from parameters file:** CUFFDIFF\_RESULTS:

GENOME:

CUFFDIFF\_OPTIONS:

CUFFMERGE\_RESULTS:

CUFFDIFF\_INPUT\_LIST\_COND1:

CUFFDIFF\_INPUT\_LIST\_COND2:

CUFFLINKS\_VERSION:

## 13.7 R Summary Report

omics\_pipe.modules.RNAseq\_report\_tuxedo.**RNAseq\_report\_tuxedo** (*sample*,

*RNAseq\_report\_tuxedo\_flag*)

Runs R script with knitr to produce report from RNAseq pipeline.

**input:** results from other steps in RNAseq pipelines

**output:** html report

**citation:**

20. Meissner

**parameters from parameter file:** WORKING\_DIR:

R\_VERSION:

REPORT\_RESULTS:

DPS\_VERSION:

PARAMS\_FILE:



---

## RNA-seq Count Based Modules

---

Modules available in the count-based RNA-seq Pipeline.

### 14.1 FASTQC

`omics_pipe.modules.fastqc.fastqc` (*sample*, *fastqc\_flag*)

QC check of raw .fastq files using FASTQC.

**input:** .fastq file

**output:** folder and zipped folder containing html, txt and image files

**citation:** Babraham Bioinformatics

**link:** <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

**parameters from parameters file:** RAW\_DATA\_DIR:

QC\_PATH:

FASTQC\_VERSION:

COMPRESSION:

### 14.2 STAR Aligner

`omics_pipe.modules.star.star` (*sample*, *star\_flag*)

Runs STAR to align .fastq files.

**input:** .fastq file

**output:** Aligned.out.bam

**citation:**

1. Dobin et al, Bioinformatics 2012; doi: 10.1093/bioinformatics/bts635 “STAR: ultrafast universal RNA-seq aligner”

**link:** <https://code.google.com/p/rna-star/>

**parameters from parameters file:** ENDS:

RAW\_DATA\_DIR:

STAR\_INDEX:

STAR\_OPTIONS:  
STAR\_RESULTS:  
SAMTOOLS\_VERSION:  
STAR\_VERSION:

## 14.3 HTSEQ-count

omics\_pipe.modules.htseq.**htseq** (*sample*, *htseq\_flag*)  
Runs htseq-count to get raw count data from alignments.

**input:** Aligned.out.sort.bam

**output:** counts.txt

**citation:** Simon Anders, EMBL

**link:** <http://www-huber.embl.de/users/anders/HTSeq/doc/overview.html>

**parameters from parameters file:** STAR\_RESULTS:

HTSEQ\_OPTIONS:  
REF\_GENES:  
HTSEQ\_RESULTS:  
TEMP\_DIR:  
SAMTOOLS\_VERSION:  
BAM\_FILE\_NAME:  
PYTHON\_VERSION:

## 14.4 R Summary Report - DESEQ2

omics\_pipe.modules.RNAseq\_report\_counts.**RNAseq\_report\_counts** (*sample*,  
*RNAseq\_report\_counts\_flag*)

Runs R script with knitr to produce report from RNAseq pipeline.

**input:** results from other steps in RNAseq pipelines

**output:** html report

**citation:**

20. Meissner

**parameters from parameter file:** WORKING\_DIR:

R\_VERSION:  
REPORT\_RESULTS:  
PARAMS\_FILE:

---

## Breast Cancer Personalized Genomics Report- RNAseq

---

Modules included in the RNAseq Cancer pipeline.

### 15.1 FASTQC

`omics_pipe.modules.fastqc.fastqc` (*sample, fastqc\_flag*)

QC check of raw .fastq files using FASTQC.

**input:** .fastq file

**output:** folder and zipped folder containing html, txt and image files

**citation:** Babraham Bioinformatics

**link:** <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

**parameters from parameters file:** RAW\_DATA\_DIR:

QC\_PATH:

FASTQC\_VERSION:

COMPRESSION:

### 15.2 STAR Aligner

`omics_pipe.modules.star.star` (*sample, star\_flag*)

Runs STAR to align .fastq files.

**input:** .fastq file

**output:** Aligned.out.bam

**citation:**

1. Dobin et al, Bioinformatics 2012; doi: 10.1093/bioinformatics/bts635 “STAR: ultrafast universal RNA-seq aligner”

**link:** <https://code.google.com/p/rna-star/>

**parameters from parameters file:** ENDS:

RAW\_DATA\_DIR:

STAR\_INDEX:

STAR\_OPTIONS:  
STAR\_RESULTS:  
SAMTOOLS\_VERSION:  
STAR\_VERSION:

## 15.3 HTSEQ-count

omics<sub>pipe</sub>.modules.htseq.**htseq** (*sample*, *htseq\_flag*)

Runs htseq-count to get raw count data from alignments.

**input:** Aligned.out.sort.bam

**output:** counts.txt

**citation:** Simon Anders, EMBL

**link:** <http://www-huber.embl.de/users/anders/HTSeq/doc/overview.html>

**parameters from parameters file:** STAR\_RESULTS:

HTSEQ\_OPTIONS:  
REF\_GENES:  
HTSEQ\_RESULTS:  
TEMP\_DIR:  
SAMTOOLS\_VERSION:  
BAM\_FILE\_NAME:  
PYTHON\_VERSION:

## 15.4 RSEQC

omics<sub>pipe</sub>.modules.rseqc.**rseqc** (*sample*, *rseqc\_flag*)

Runs rseqc to determine insert size as QC for alignment.

**input:** .bam

**output:** pdf plot

**link:** <http://rseqc.sourceforge.net/>

**parameters from parameters file:** STAR\_RESULTS:

QC\_PATH:  
BAM\_FILE\_NAME:  
RSEQC\_REF:  
RSEQC\_VERSION:  
TEMP\_DIR:

## 15.5 Fusion Catcher

omics\_pipe.modules.fusion\_catcher.fusion\_catcher (sample, fusion\_catcher\_flag)

Detects fusion genes in paired-end RNAseq data.

**input:** paired end .fastq files

**output:** list of candidate fusion genes

**citation:**

19. Kangaspeska, S. Hultsch, H. Edgren, D. Nicorici, A. Murumgi, O.P. Kallioniemi, Reanalysis of RNA-sequencing data reveals several additional fusion genes with multiple isoforms, PLOS One, Oct. 2012. <http://dx.plos.org/10.1371/journal.pone.0048745>

**link:** <https://code.google.com/p/fusioncatcher>

**parameters from parameters file:** ENDS:

RAW\_DATA\_DIR:

FUSION\_RESULTS:

FUSIONCATCHERBUILD\_DIR:

TEMP\_DIR:

SAMTOOLS\_VERSION:

FUSIONCATCHER\_VERSION:

FUSIONCATCHER\_OPTIONS:

TISSUE:

PYTHON\_VERSION:

## 15.6 BWA/SNPiR

### 15.6.1 BWA

omics\_pipe.modules.bwa.bwa1 (sample, bwa1\_flag)

BWA aligner for read1 of paired\_end reads.

**input:** .fastq

**output:** .sam

**citation:** Li H. and Durbin R. (2009) Fast and accurate short read alignment with Burrows-Wheeler transform. Bioinformatics, 25, 1754-1760. [PMID: 19451168]

**link:** <http://bio-bwa.sourceforge.net/bwa.shtml>

**parameters from parameters file:** BWA\_RESULTS:

TEMP\_DIR:

SAMTOOLS\_VERSION:

BWA\_VERSION:

BWA\_INDEX:

RAW\_DATA\_DIR:

GATK\_READ\_GROUP\_INFO:

COMPRESSION:

omics<sub>pipe</sub>.modules.bwa.**bwa2** (*sample, bwa2\_flag*)

BWA aligner for read2 of paired\_end reads.

**input:** .fastq

**output:** .sam

**citation:** Li H. and Durbin R. (2009) Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics*, 25, 1754-1760. [PMID: 19451168]

**link:** <http://bio-bwa.sourceforge.net/bwa.shtml>

**parameters from parameters file:** BWA\_RESULTS:

TEMP\_DIR:

SAMTOOLS\_VERSION:

BWA\_VERSION:

BWA\_INDEX:

RAW\_DATA\_DIR:

GATK\_READ\_GROUP\_INFO:

COMPRESSION:

omics<sub>pipe</sub>.modules.bwa.**bwa\_RNA** (*sample, bwa\_flag*)

BWA aligner for single end reads.

**input:** .fastq

**output:** .sam

**citation:** Li H. and Durbin R. (2009) Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics*, 25, 1754-1760. [PMID: 19451168]

**link:** <http://bio-bwa.sourceforge.net/bwa.shtml>

**parameters from parameters file:** BWA\_RESULTS:

TEMP\_DIR:

SAMTOOLS\_VERSION:

BWA\_VERSION:

BWA\_INDEX:

RAW\_DATA\_DIR:

GATK\_READ\_GROUP\_INFO:

COMPRESSION:

## 15.6.2 SNPiR

omics<sub>pipe</sub>.modules.snpir\_variants.**snpir\_variants** (*sample, snpir\_variants\_flag*)

Calls variants using SNPiR pipeline.

**input:** Aligned.out.sort.bam or accepted\_hits.bam

**output:** final\_variants.vcf file

**citation:** Piskol, R., et al. (2013). “Reliable Identification of Genomic Variants from RNA-Seq Data.” The American Journal of Human Genetics 93(4): 641-651.

**link:** <http://lilab.stanford.edu/SNPiR/>

**parameters from parameters file:** VARIANT\_RESULTS:

TEMP\_DIR:  
 SAMTOOLS\_VERSION:  
 BWA\_VERSION:  
 PICARD\_VERSION:  
 GATK\_VERSION:  
 BEDTOOLS\_VERSION:  
 UCSC\_TOOLS\_VERSION:  
 GENOME:  
 REPEAT\_MASKER:  
 SNPIR\_ANNOTATION:  
 RNA\_EDIT:  
 DBSNP:  
 MILLS:  
 G1000:  
 WORKING\_DIR:  
 BWA\_RESULTS:  
 SNPIR\_VERSION:  
 SNPIR\_CONFIG:  
 SNPIR\_DIR:  
 ENCODING:

### 15.6.3 Filter Variants

omics\_pipe.modules.filter\_variants.**filter\_variants** (*sample, filter\_variants\_flag*)

Filters variants to remove common variants.

**input:** .bam or .sam file

**output:** .vcf file

**citation:** Piskol et al. 2013. Reliable identification of genomic variants from RNA-seq data. The American Journal of Human Genetics 93: 641-651.

**link:** <http://lilab.stanford.edu/SNPiR/>

**parameters from parameters file:** VARIANT\_RESULTS:

TEMP\_DIR:  
 SAMTOOLS\_VERSION:  
 BWA\_VERSION:

PICARD\_VERSION:  
GATK\_VERSION:  
BEDTOOLS\_VERSION:  
UCSC\_TOOLS\_VERSION:  
GENOME:  
REPEAT\_MASKER:  
SNPIR\_ANNOTATION:  
RNA\_EDIT:  
DBSNP:  
MILLS:  
G1000:  
WORKING\_DIR:  
BWA\_RESULTS:  
SNPIR\_VERSION:  
SNPIR\_CONFIG:  
SNPIR\_DIR:  
SNPEFF\_VERSION:  
dbNSFP:  
VCFTOOLS\_VERSION:  
WORKING\_DIR:  
SNP\_FILTER\_OUT\_REF:

## 15.7 Intogen

omics<sub>pipe</sub>.modules.intogen.**intogen** (*sample, intogen\_flag*)

Runs Intogen to rank mutations and implication for cancer phenotype. Follows variant calling.

**input:** .vcf

**output:** variant list

**citation:** Gonzalez-Perez et al. 2013. Intogen mutations identifies cancer drivers across tumor types. Nature Methods 10, 1081-1082.

**link:** <http://www.intogen.org/>

**parameters from parameter file:** VCF\_FILE:

INTOGEN\_OPTIONS:  
INTOGEN\_RESULTS:  
INTOGEN\_VERSION:  
USERNAME:  
WORKING\_DIR:



TEMP\_DIR:  
SCHEDULER:  
VARIANT\_RESULTS:

## 15.8 OncoRep Cancer Report

omics\_pipe.modules.BreastCancer\_RNA\_report.**BreastCancer\_RNA\_report** (*sample*,  
*Breast-*  
*Cancer\_RNA\_report\_flag*)

Runs R script with knitr to produce report from RNAseq pipeline.

**input:** results from other steps in RNAseq pipelines

**output:** html report

**citation:**

20. Meissner

**parameters from parameter file:** WORKING\_DIR:

R\_VERSION:  
REPORT\_RESULTS:  
PARAMS\_FILE:  
TABIX\_VERSION:  
TUMOR\_TYPE:  
GENELIST:  
COSMIC:  
CLINVAR:  
PHARMGKB\_rsID:  
PHARMGKB\_Allele:  
DRUGBANK:  
CADD:



---

## TCGA Reanalysis Pipeline - RNAseq

---

Modules included in the TCGA RNAseq Cancer pipeline.

### 16.1 TCGA Download

`omics_pipe.modules.TCGA_download.TCGA_download` (*sample*, *TCGA\_download\_flag*)

Downloads and unzips TCGA data from Manifest.xml downloaded from CGHub. input:

TCGA XML file

**output:** downloaded files from TCGA

**citation:** The Cancer Genome Atlas

**link:** <https://cghub.ucsc.edu/software/downloads.html>

**parameters from parameters file:** TCGA\_XML\_FILE:

TCGA\_KEY:

TCGA\_OUTPUT\_PATH:

CGATOOLS\_VERSION:

### 16.2 FASTQC

`omics_pipe.modules.fastqc.fastqc` (*sample*, *fastqc\_flag*)

QC check of raw .fastq files using FASTQC.

**input:** .fastq file

**output:** folder and zipped folder containing html, txt and image files

**citation:** Babraham Bioinformatics

**link:** <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

**parameters from parameters file:** RAW\_DATA\_DIR:

QC\_PATH:

FASTQC\_VERSION:

COMPRESSION:

## 16.3 STAR Aligner

omics\_pipe.modules.star.**star** (*sample*, *star\_flag*)

Runs STAR to align .fastq files.

**input:** .fastq file

**output:** Aligned.out.bam

**citation:**

1. Dobin et al, Bioinformatics 2012; doi: 10.1093/bioinformatics/bts635 “STAR: ultrafast universal RNA-seq aligner”

**link:** <https://code.google.com/p/rna-star/>

**parameters from parameters file:** ENDS:

RAW\_DATA\_DIR:

STAR\_INDEX:

STAR\_OPTIONS:

STAR\_RESULTS:

SAMTOOLS\_VERSION:

STAR\_VERSION:

## 16.4 HTSEQ-count

omics\_pipe.modules.htseq.**htseq** (*sample*, *htseq\_flag*)

Runs htseq-count to get raw count data from alignments.

**input:** Aligned.out.sort.bam

**output:** counts.txt

**citation:** Simon Anders, EMBL

**link:** <http://www-huber.embl.de/users/anders/HTSeq/doc/overview.html>

**parameters from parameters file:** STAR\_RESULTS:

HTSEQ\_OPTIONS:

REF\_GENES:

HTSEQ\_RESULTS:

TEMP\_DIR:

SAMTOOLS\_VERSION:

BAM\_FILE\_NAME:

PYTHON\_VERSION:

## 16.5 RSEQC

omics\_pipe.modules.rseqc.**rseqc** (*sample*, *rseqc\_flag*)

Runs rseqc to determine insert size as QC for alignment.

**input:** .bam

**output:** pdf plot

**link:** <http://rseqc.sourceforge.net/>

**parameters from parameters file:** STAR\_RESULTS:

QC\_PATH:

BAM\_FILE\_NAME:

RSEQC\_REF:

RSEQC\_VERSION:

TEMP\_DIR:

## 16.6 Fusion Catcher

omics\_pipe.modules.fusion\_catcher.**fusion\_catcher** (*sample*, *fusion\_catcher\_flag*)

Detects fusion genes in paired-end RNAseq data.

**input:** paired end .fastq files

**output:** list of candidate fusion genes

**citation:**

19. Kangaspeska, S. Hultsch, H. Edgren, D. Nicorici, A. Murumgi, O.P. Kallioniemi, Reanalysis of RNA-sequencing data reveals several additional fusion genes with multiple isoforms, PLOS One, Oct. 2012. <http://dx.plos.org/10.1371/journal.pone.0048745>

**link:** <https://code.google.com/p/fusioncatcher>

**parameters from parameters file:** ENDS:

RAW\_DATA\_DIR:

FUSION\_RESULTS:

FUSIONCATCHERBUILD\_DIR:

TEMP\_DIR:

SAMTOOLS\_VERSION:

FUSIONCATCHER\_VERSION:

FUSIONCATCHER\_OPTIONS:

TISSUE:

PYTHON\_VERSION:

## 16.7 BWA/SNPiR

### 16.7.1 BWA

omics\_pipe.modules.bwa.**bwa1** (*sample, bwa1\_flag*)

BWA aligner for read1 of paired\_end reads.

**input:** .fastq

**output:** .sam

**citation:** Li H. and Durbin R. (2009) Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics*, 25, 1754-1760. [PMID: 19451168]

**link:** <http://bio-bwa.sourceforge.net/bwa.shtml>

**parameters from parameters file:** BWA\_RESULTS:

TEMP\_DIR:

SAMTOOLS\_VERSION:

BWA\_VERSION:

BWA\_INDEX:

RAW\_DATA\_DIR:

GATK\_READ\_GROUP\_INFO:

COMPRESSION:

omics\_pipe.modules.bwa.**bwa2** (*sample, bwa2\_flag*)

BWA aligner for read2 of paired\_end reads.

**input:** .fastq

**output:** .sam

**citation:** Li H. and Durbin R. (2009) Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics*, 25, 1754-1760. [PMID: 19451168]

**link:** <http://bio-bwa.sourceforge.net/bwa.shtml>

**parameters from parameters file:** BWA\_RESULTS:

TEMP\_DIR:

SAMTOOLS\_VERSION:

BWA\_VERSION:

BWA\_INDEX:

RAW\_DATA\_DIR:

GATK\_READ\_GROUP\_INFO:

COMPRESSION:

omics\_pipe.modules.bwa.**bwa\_RNA** (*sample, bwa\_flag*)

BWA aligner for single end reads.

**input:** .fastq

**output:** .sam

**citation:** Li H. and Durbin R. (2009) Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics*, 25, 1754-1760. [PMID: 19451168]

**link:** <http://bio-bwa.sourceforge.net/bwa.shtml>

**parameters from parameters file:** BWA\_RESULTS:

TEMP\_DIR:  
 SAMTOOLS\_VERSION:  
 BWA\_VERSION:  
 BWA\_INDEX:  
 RAW\_DATA\_DIR:  
 GATK\_READ\_GROUP\_INFO:  
 COMPRESSION:

## 16.7.2 SNPIR

omics<sub>pipe</sub>.modules.snpir\_variants.**snpir\_variants** (*sample, snpir\_variants\_flag*)  
 Calls variants using SNPIR pipeline.

**input:** Aligned.out.sort.bam or accepted\_hits.bam

**output:** final\_variants.vcf file

**citation:** Piskol, R., et al. (2013). “Reliable Identification of Genomic Variants from RNA-Seq Data.” *The American Journal of Human Genetics* 93(4): 641-651.

**link:** <http://lilab.stanford.edu/SNPiR/>

**parameters from parameters file:** VARIANT\_RESULTS:

TEMP\_DIR:  
 SAMTOOLS\_VERSION:  
 BWA\_VERSION:  
 PICARD\_VERSION:  
 GATK\_VERSION:  
 BEDTOOLS\_VERSION:  
 UCSC\_TOOLS\_VERSION:  
 GENOME:  
 REPEAT\_MASKER:  
 SNPIR\_ANNOTATION:  
 RNA\_EDIT:  
 DBSNP:  
 MILLS:  
 G1000:  
 WORKING\_DIR:  
 BWA\_RESULTS:

SNPIR\_VERSION:  
SNPIR\_CONFIG:  
SNPIR\_DIR:  
ENCODING:

### 16.7.3 Filter Variants

omics<sub>pipe</sub>.modules.filter\_variants.**filter\_variants** (*sample, filter\_variants\_flag*)

Filters variants to remove common variants.

**input:** .bam or .sam file

**output:** .vcf file

**citation:** Piskol et al. 2013. Reliable identification of genomic variants from RNA-seq data. The American Journal of Human Genetics 93: 641-651.

**link:** <http://lilab.stanford.edu/SNPiR/>

**parameters from parameters file:** VARIANT\_RESULTS:

TEMP\_DIR:  
SAMTOOLS\_VERSION:  
BWA\_VERSION:  
PICARD\_VERSION:  
GATK\_VERSION:  
BEDTOOLS\_VERSION:  
UCSC\_TOOLS\_VERSION:  
GENOME:  
REPEAT\_MASKER:  
SNPIR\_ANNOTATION:  
RNA\_EDIT:  
DBSNP:  
MILLS:  
G1000:  
WORKING\_DIR:  
BWA\_RESULTS:  
SNPIR\_VERSION:  
SNPIR\_CONFIG:  
SNPIR\_DIR:  
SNPEFF\_VERSION:  
dbNSFP:  
VCFTOOLS\_VERSION:  
WORKING\_DIR:



SNP\_FILTER\_OUT\_REF:

## 16.8 Intogen

omics\_pipe.modules.intogen.**intogen** (*sample, intogen\_flag*)

Runs Intogen to rank mutations and implication for cancer phenotype. Follows variant calling.

**input:** .vcf

**output:** variant list

**citation:** Gonzalez-Perez et al. 2013. Intogen mutations identifies cancer drivers across tumor types. Nature Methods 10, 1081-1082.

**link:** <http://www.intogen.org/>

**parameters from parameter file:** VCF\_FILE:

INTOGEN\_OPTIONS:

INTOGEN\_RESULTS:

INTOGEN\_VERSION:

USERNAME:

WORKING\_DIR:

TEMP\_DIR:

SCHEDULER:

VARIANT\_RESULTS:

## 16.9 OncoRep Cancer Report

omics\_pipe.modules.BreastCancer\_RNA\_report.**BreastCancer\_RNA\_report** (*sample, Breast-*

*Cancer\_RNA\_report\_flag*)

Runs R script with knitr to produce report from RNAseq pipeline.

**input:** results from other steps in RNAseq pipelines

**output:** html report

**citation:**

20. Meissner

**parameters from parameter file:** WORKING\_DIR:

R\_VERSION:

REPORT\_RESULTS:

PARAMS\_FILE:

TABIX\_VERSION:

TUMOR\_TYPE:

GENELIST:

COSMIC:

CLINVAR:

PHARMGKB\_rsID:

PHARMGKB\_Allele:

DRUGBANK:

CADD:

---

## RNA-seq Count Based Modules- TCGA

---

Modules available in the TCGA count-based RNA-seq Pipeline.

### 17.1 TCGA Download

`omics_pipe.modules.TCGA_download.TCGA_download` (*sample*, *TCGA\_download\_flag*)

Downloads and unzips TCGA data from Manifest.xml downloaded from CGHub. input:

TCGA XML file

**output:** downloaded files from TCGA

**citation:** The Cancer Genome Atlas

**link:** <https://cghub.ucsc.edu/software/downloads.html>

**parameters from parameters file:** TCGA\_XML\_FILE:

TCGA\_KEY:

TCGA\_OUTPUT\_PATH:

CGATOOLS\_VERSION:

### 17.2 FASTQC

`omics_pipe.modules.fastqc.fastqc` (*sample*, *fastqc\_flag*)

QC check of raw .fastq files using FASTQC.

**input:** .fastq file

**output:** folder and zipped folder containing html, txt and image files

**citation:** Babraham Bioinformatics

**link:** <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

**parameters from parameters file:** RAW\_DATA\_DIR:

QC\_PATH:

FASTQC\_VERSION:

COMPRESSION:

## 17.3 STAR Aligner

omics<sub>pipe</sub>.modules.star.**star** (*sample*, *star\_flag*)

Runs STAR to align .fastq files.

**input:** .fastq file

**output:** Aligned.out.bam

**citation:**

1. Dobin et al, Bioinformatics 2012; doi: 10.1093/bioinformatics/bts635 “STAR: ultrafast universal RNA-seq aligner”

**link:** <https://code.google.com/p/rna-star/>

**parameters from parameters file:** ENDS:

RAW\_DATA\_DIR:

STAR\_INDEX:

STAR\_OPTIONS:

STAR\_RESULTS:

SAMTOOLS\_VERSION:

STAR\_VERSION:

## 17.4 HTSEQ-count

omics<sub>pipe</sub>.modules.htseq.**htseq** (*sample*, *htseq\_flag*)

Runs htseq-count to get raw count data from alignments.

**input:** Aligned.out.sort.bam

**output:** counts.txt

**citation:** Simon Anders, EMBL

**link:** <http://www-huber.embl.de/users/anders/HTSeq/doc/overview.html>

**parameters from parameters file:** STAR\_RESULTS:

HTSEQ\_OPTIONS:

REF\_GENES:

HTSEQ\_RESULTS:

TEMP\_DIR:

SAMTOOLS\_VERSION:

BAM\_FILE\_NAME:

PYTHON\_VERSION:

## 17.5 R Summary Report - DESEQ2

omics\_pipe.modules.RNAseq\_report\_counts.**RNAseq\_report\_counts** (*sample*,  
*RNAseq\_report\_counts\_flag*)

Runs R script with knitr to produce report from RNAseq pipeline.

**input:** results from other steps in RNAseq pipelines

**output:** html report

**citation:**

20. Meissner

**parameters from parameter file:** WORKING\_DIR:

R\_VERSION:

REPORT\_RESULTS:

PARAMS\_FILE:



---

## miRNA-seq Tuxedo Modules

---

Modules available in the miRNA-seq Tuxedo Pipeline.

### 18.1 CutAdapt

`omics_pipe.modules.cutadapt_miRNA.cutadapt_miRNA` (*sample*, *cutadapt\_miRNA\_flag*)

Runs Cutadapt to trim adapters from reads.

**input:** .fastq

**output:** .fastq

**citation:** Martin 2011. Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnet.journal* 17: 10-12.

**link:** <https://code.google.com/p/cutadapt/>

**parameters from parameters file:** RAW\_DATA\_DIR:

ADAPTER:

TRIMMED\_DATA\_PATH:

PYTHON\_VERSION

### 18.2 Fastq Length Filter

`omics_pipe.modules.fastq_length_filter_miRNA.fastq_length_filter_miRNA` (*sample*, *fastq\_length\_filter\_miRNA\_flag*)

Runs custom Python script to filter miRNA reads by length.

**input:** .fastq

**output:** .fastq

**parameters from parameter file:** TRIMMED\_DATA\_PATH:

### 18.3 FASTQC

`omics_pipe.modules.fastqc.fastqc` (*sample*, *fastqc\_flag*)

QC check of raw .fastq files using FASTQC.

**input:** .fastq file

**output:** folder and zipped folder containing html, txt and image files

**citation:** Babraham Bioinformatics

**link:** <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

**parameters from parameters file:** RAW\_DATA\_DIR:

QC\_PATH:

FASTQC\_VERSION:

COMPRESSION:

## 18.4 TopHat

omics<sub>pipe</sub>.modules.tophat.**tophat** (*sample, tophat\_flag*)

Runs TopHat to align .fastq files.

**input:** .fastq file

**output:** accepted\_hits.bam

**citation:** Trapnell C, Pachter L, Salzberg SL. TopHat: discovering splice junctions with RNA-Seq. Bioinformatics doi:10.1093/bioinformatics/btp120

**link:** <http://tophat.cbcb.umd.edu/>

**parameters from parameters file:** RAW\_DATA\_DIR:

REF\_GENES:

TOPHAT\_RESULTS:

BOWTIE\_INDEX:

TOPHAT\_VERSION:

TOPHAT\_OPTIONS:

BOWTIE\_VERSION:

SAMTOOLS\_VERSION:

## 18.5 Cufflinks

omics<sub>pipe</sub>.modules.cufflinks.**cufflinks** (*sample, cufflinks\_flag*)

Runs cufflinks to assemble .bam files from TopHat.

**input:** accepted\_hits.bam

**output:** transcripts.gtf

**citation:** Trapnell C, et al. Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation Nature Biotechnology doi:10.1038/nbt.1621

**link:** <http://cufflinks.cbcb.umd.edu/>



**parameters from parameters file:** TOPHAT\_RESULTS:

CUFFLINKS\_RESULTS:

REF\_GENES:

GENOME:

CUFFLINKS\_OPTIONS:

CUFFLINKS\_VERSION:

## 18.6 Cuffmerge

omics<sub>pipe</sub>.modules.cuffmerge.**cuffmerge** (*step, cuffmerge\_flag*)

Runs cuffmerge to merge .gtf files from Cufflinks.

**input:** assembly\_GTF\_list.txt

**output:** merged.gtf

**citation:** Trapnell C, et al. Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation Nature Biotechnology doi:10.1038/nbt.1621

**link:** <http://cufflinks.cbcb.umd.edu/>

**parameters from parameters file:** CUFFMERGE\_RESULTS:

REF\_GENES:

GENOME:

CUFFMERGE\_OPTIONS:

CUFFLINKS\_VERSION:

## 18.7 Cuffmergetocompare

omics<sub>pipe</sub>.modules.cuffmergetocompare.**cuffmergetocompare** (*step, cuffmergetocompare\_flag*)

Runs cuffcompare to annotate merged .gtf files from Cufflinks.

**input:** assembly\_GTF\_list.txt

**output:** merged.gtf

**citation:** Trapnell C, et al. Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation Nature Biotechnology doi:10.1038/nbt.1621

**link:** <http://cufflinks.cbcb.umd.edu/>

**parameters from parameters file:** CUFFMERGE\_RESULTS:

REF\_GENES:

GENOME:

CUFFMERGETOCOMPARE\_OPTIONS:

CUFFLINKS\_VERSION:

## 18.8 Cuffdiff

omics<sub>pipe</sub>.modules.cuffdiff.**cuffdiff** (*step, cuffdiff\_flag*)

Runs Cuffdiff to perform differential expression. Runs after Cufflinks. Part of Tuxedo Suite.

**input:** .bam files

**output:** differential expression results

**citation:** Trapnell C, et al. Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation Nature Biotechnology doi:10.1038/nbt.1621

**link:** <http://cufflinks.cbc.umd.edu/>

**parameters from parameters file:** CUFFDIFF\_RESULTS:

GENOME:

CUFFDIFF\_OPTIONS:

CUFFMERGE\_RESULTS:

CUFFDIFF\_INPUT\_LIST\_COND1:

CUFFDIFF\_INPUT\_LIST\_COND2:

CUFFLINKS\_VERSION:

## 18.9 R Summary Report

omics<sub>pipe</sub>.modules.RNAseq\_report\_tuxedo.**RNAseq\_report\_tuxedo** (*sample,*

*RNAseq\_report\_tuxedo\_flag*)

Runs R script with knitr to produce report from RNAseq pipeline.

**input:** results from other steps in RNAseq pipelines

**output:** html report

**citation:**

20. Meissner

**parameters from parameter file:** WORKING\_DIR:

R\_VERSION:

REPORT\_RESULTS:

DPS\_VERSION:

PARAMS\_FILE:

---

## miRNA-seq Count Based Modules

---

Modules available in the count-based miRNA-seq Pipeline.

### 19.1 CutAdapt

`omics_pipe.modules.cutadapt_miRNA.cutadapt_miRNA` (*sample, cutadapt\_miRNA\_flag*)

Runs Cutadapt to trim adapters from reads.

**input:** .fastq

**output:** .fastq

**citation:** Martin 2011. Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnet.journal* 17: 10-12.

**link:** <https://code.google.com/p/cutadapt/>

**parameters from parameters file:** RAW\_DATA\_DIR:

ADAPTER:

TRIMMED\_DATA\_PATH:

PYTHON\_VERSION

### 19.2 Fastq Length Filter

`omics_pipe.modules.fastq_length_filter_miRNA.fastq_length_filter_miRNA` (*sample, fastq\_length\_filter\_miRNA\_f*)

Runs custom Python script to filter miRNA reads by length.

**input:** .fastq

**output:** .fastq

**parameters from parameter file:** TRIMMED\_DATA\_PATH:

### 19.3 FASTQC

`omics_pipe.modules.fastqc.fastqc` (*sample, fastqc\_flag*)

QC check of raw .fastq files using FASTQC.

**input:** .fastq file

**output:** folder and zipped folder containing html, txt and image files

**citation:** Babraham Bioinformatics

**link:** <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

**parameters from parameters file:** RAW\_DATA\_DIR:

QC\_PATH:

FASTQC\_VERSION:

COMPRESSION:

## 19.4 STAR Aligner

omics<sub>pipe</sub>.modules.star.star (sample, star\_flag)

Runs STAR to align .fastq files.

**input:** .fastq file

**output:** Aligned.out.bam

**citation:**

1. Dobin et al, Bioinformatics 2012; doi: 10.1093/bioinformatics/bts635 “STAR: ultrafast universal RNA-seq aligner”

**link:** <https://code.google.com/p/rna-star/>

**parameters from parameters file:** ENDS:

RAW\_DATA\_DIR:

STAR\_INDEX:

STAR\_OPTIONS:

STAR\_RESULTS:

SAMTOOLS\_VERSION:

STAR\_VERSION:

## 19.5 HTSEQ

omics<sub>pipe</sub>.modules.htseq.htseq (sample, htseq\_flag)

Runs htseq-count to get raw count data from alignments.

**input:** Aligned.out.sort.bam

**output:** counts.txt

**citation:** Simon Anders, EMBL

**link:** <http://www-huber.embl.de/users/anders/HTSeq/doc/overview.html>

**parameters from parameters file:** STAR\_RESULTS:

HTSEQ\_OPTIONS:

REF\_GENES:

HTSEQ\_RESULTS:  
TEMP\_DIR:  
SAMTOOLS\_VERSION:  
BAM\_FILE\_NAME:  
PYTHON\_VERSION:

## 19.6 R Summary Report - DESEQ2

omics\_pipe.modules.RNAseq\_report\_counts.**RNAseq\_report\_counts** (*sample*,  
*RNAseq\_report\_counts\_flag*)

Runs R script with knitr to produce report from RNAseq pipeline.

**input:** results from other steps in RNAseq pipelines

**output:** html report

**citation:**

20. Meissner

**parameters from parameter file:** WORKING\_DIR:

R\_VERSION:

REPORT\_RESULTS:

PARAMS\_FILE:



---

## ChIP-seq Modules – HOMER

---

### 20.1 FASTQC

`omics_pipe.modules.fastqc.fastqc` (*sample, fastqc\_flag*)

QC check of raw .fastq files using FASTQC.

**input:** .fastq file

**output:** folder and zipped folder containing html, txt and image files

**citation:** Babraham Bioinformatics

**link:** <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

**parameters from parameters file:** RAW\_DATA\_DIR:

QC\_PATH:

FASTQC\_VERSION:

COMPRESSION:

### 20.2 ChIP trim

`omics_pipe.modules.ChIP_trim.ChIP_trim` (*sample, ChIP\_trim\_flag*)

Runs Homer Tools to trim adapters from .fastq files.

**input:** .fastq file

**output:** .fastq file

**citation:** Heinz S, Benner C, Spann N, Bertolino E et al. Simple Combinations of Lineage-Determining Transcription Factors Prime cis-Regulatory Elements Required for Macrophage and B Cell Identities. Mol Cell 2010 May 28;38(4):576-589. PMID: 20513432

**link:** <http://homer.salk.edu/homer/>

**parameters from parameters file:** ENDS:

RAW\_DATA\_DIR:

HOMER\_TRIM\_OPTIONS:

TRIMMED\_DATA\_PATH:

HOMER\_VERSION:

## 20.3 Bowtie

omics<sub>pipe</sub>.modules.bowtie.**bowtie** (*sample, bowtie\_flag*)

Runs Bowtie to align .fastq files.

**input:** .fastq file

**output:** sample.bam

**citation:** Langmead B, Trapnell C, Pop M, Salzberg SL. Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. Genome Biology 10:R25

**link:** <http://bowtie-bio.sourceforge.net/index.shtml>

**parameters from parameters file:** ENDS:

TRIMMED\_DATA\_PATH:

BOWTIE\_OPTIONS:

BOWTIE\_INDEX:

BOWTIE\_RESULTS:

BOWTIE\_VERSION:

SAMTOOLS\_VERSION:

BEDTOOLS\_VERSION:

## 20.4 Read Density -HOMER

omics<sub>pipe</sub>.modules.read\_density.**read\_density** (*sample, read\_density\_flag*)

Runs HOMER to visualize read density from ChIPseq data.

**input:** .bam file

**output:** .txt file

**citation:** Heinz S, Benner C, Spann N, Bertolino E et al. Simple Combinations of Lineage-Determining Transcription Factors Prime cis-Regulatory Elements Required for Macrophage and B Cell Identities. Mol Cell 2010 May 28;38(4):576-589. PMID: 20513432

**link:** <http://homer.salk.edu/homer/>

**parameters from parameters file:** BOWTIE\_RESULTS:

CHROM\_SIZES:

HOMER\_RESULTS:

HOMER\_VERSION:

TEMP\_DIR:

## 20.5 Peak Detection - HOMER

omics<sub>pipe</sub>.modules.homer\_peaks.**homer\_peaks** (*step, homer\_peaks\_flag*)

Runs HOMER to call peaks from ChIPseq data.

**input:** .tag input file



**output:** .txt file

**citation:** Heinz S, Benner C, Spann N, Bertolino E et al. Simple Combinations of Lineage-Determining Transcription Factors Prime cis-Regulatory Elements Required for Macrophage and B Cell Identities. Mol Cell 2010 May 28;38(4):576-589. PMID: 20513432

**link:** <http://homer.salk.edu/homer/>

**parameters from parameters file:** PAIR\_LIST:

HOMER\_RESULTS:

HOMER\_PEAKE\_OPTIONS:

HOMER\_VERSION:

TEMP\_DIR:

## 20.6 Peak Annotation & Visualization - HOMER

omics<sub>pipe</sub>.modules.peak\_track.**peak\_track** (*step, peak\_track\_flag*)

Runs HOMER to create peak track from ChIPseq data.

**input:** .tag input file

**output:** .txt file

**citation:** Heinz S, Benner C, Spann N, Bertolino E et al. Simple Combinations of Lineage-Determining Transcription Factors Prime cis-Regulatory Elements Required for Macrophage and B Cell Identities. Mol Cell 2010 May 28;38(4):576-589. PMID: 20513432

**link:** <http://homer.salk.edu/homer/>

**parameters from parameters file:** PAIR\_LIST:

HOMER\_RESULTS:

HOMER\_VERSION:

TEMP\_DIR:

omics<sub>pipe</sub>.modules.annotate\_peaks.**annotate\_peaks** (*step, annotate\_peaks\_flag*)

Runs HOMER to annotate peak track from ChIPseq data.

**input:** .tag input file

**output:** .txt file

**citation:** Heinz S, Benner C, Spann N, Bertolino E et al. Simple Combinations of Lineage-Determining Transcription Factors Prime cis-Regulatory Elements Required for Macrophage and B Cell Identities. Mol Cell 2010 May 28;38(4):576-589. PMID: 20513432

**link:** <http://homer.salk.edu/homer/>

**parameters from parameters file:** PAIR\_LIST:

HOMER\_RESULTS:

HOMER\_VERSION:

TEMP\_DIR:

HOMER\_GENOME:

HOMER\_ANNOTATE\_OPTIONS:

## 20.7 Find Motifs - HOMER

omics<sub>pipe</sub>.modules.find\_motifs.**find\_motifs** (*step, find\_motifs\_flag*)

Runs HOMER to find motifs from ChIPseq data.

**input:** .txt peak file from Homer

**output:** .txt file

**citation:** Heinz S, Benner C, Spann N, Bertolino E et al. Simple Combinations of Lineage-Determining Transcription Factors Prime cis-Regulatory Elements Required for Macrophage and B Cell Identities. Mol Cell 2010 May 28;38(4):576-589. PMID: 20513432

**link:** <http://homer.salk.edu/homer/>

**parameters from parameters file:** PAIR\_LIST:

HOMER\_RESULTS:

HOMER\_VERSION:

TEMP\_DIR:

HOMER\_GENOME:

HOMER\_MOTIFS\_OPTIONS:

---

## ChIP-seq Modules – MACS

---

### 21.1 FASTQC

`omics_pipe.modules.fastqc.fastqc` (*sample, fastqc\_flag*)

QC check of raw .fastq files using FASTQC.

**input:** .fastq file

**output:** folder and zipped folder containing html, txt and image files

**citation:** Babraham Bioinformatics

**link:** <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

**parameters from parameters file:** RAW\_DATA\_DIR:

QC\_PATH:

FASTQC\_VERSION:

COMPRESSION:

### 21.2 ChIP trim

`omics_pipe.modules.ChIP_trim.ChIP_trim` (*sample, ChIP\_trim\_flag*)

Runs Homer Tools to trim adapters from .fastq files.

**input:** .fastq file

**output:** .fastq file

**citation:** Heinz S, Benner C, Spann N, Bertolino E et al. Simple Combinations of Lineage-Determining Transcription Factors Prime cis-Regulatory Elements Required for Macrophage and B Cell Identities. Mol Cell 2010 May 28;38(4):576-589. PMID: 20513432

**link:** <http://homer.salk.edu/homer/>

**parameters from parameters file:** ENDS:

RAW\_DATA\_DIR:

HOMER\_TRIM\_OPTIONS:

TRIMMED\_DATA\_PATH:

HOMER\_VERSION:

## 21.3 Bowtie

omics<sub>pipe</sub>.modules.bowtie.**bowtie** (*sample, bowtie\_flag*)

Runs Bowtie to align .fastq files.

**input:** .fastq file

**output:** sample.bam

**citation:** Langmead B, Trapnell C, Pop M, Salzberg SL. Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. Genome Biology 10:R25

**link:** <http://bowtie-bio.sourceforge.net/index.shtml>

**parameters from parameters file:** ENDS:

TRIMMED\_DATA\_PATH:

BOWTIE\_OPTIONS:

BOWTIE\_INDEX:

BOWTIE\_RESULTS:

BOWTIE\_VERSION:

SAMTOOLS\_VERSION:

BEDTOOLS\_VERSION:

## 21.4 MACS

omics<sub>pipe</sub>.modules.macs.**macs** (*step, macs\_flag*)

Runs MACS to call peaks from ChIPseq data. input:

.fastq file

**output:** peaks and .bed file

**citation:** Zhang et al. Model-based Analysis of ChIP-Seq (MACS). Genome Biol (2008) vol. 9 (9) pp. R137

**link:** <http://liulab.dfci.harvard.edu/MACS/>

**parameters from parameters file:** PAIR\_LIST:

BOWTIE\_RESULTS:

CHROM\_SIZES:

MACS\_RESULTS:

MACS\_VERSION:

TEMP\_DIR:

BEDTOOLS\_VERSION:

PYTHON\_VERSION:

---

## Whole Genome and Whole Exome Sequencing Modules

---

### 22.1 FASTQC

`omics_pipe.modules.fastqc.fastqc` (*sample*, *fastqc\_flag*)

QC check of raw .fastq files using FASTQC.

**input:** .fastq file

**output:** folder and zipped folder containing html, txt and image files

**citation:** Babraham Bioinformatics

**link:** <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

**parameters from parameters file:** RAW\_DATA\_DIR:

QC\_PATH:

FASTQC\_VERSION:

COMPRESSION:

### 22.2 BWA-MEM

`omics_pipe.modules.bwa.bwa_mem` (*sample*, *bwa\_mem\_flag*)

BWA aligner with BWA-MEM algorithm.

**input:** .fastq

**output:** .sam

**citation:** Li H. and Durbin R. (2009) Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics*, 25, 1754-1760. [PMID: 19451168]

**link:** <http://bio-bwa.sourceforge.net/bwa.shtml>

**parameters from parameters file:** BWA\_RESULTS:

TEMP\_DIR:

SAMTOOLS\_VERSION:

BWA\_VERSION:

GENOME:

RAW\_DATA\_DIR:

BWA\_OPTIONS:

COMPRESSION:

## 22.3 PICARD Mark Duplicates

omics\_pipe.modules.picard\_mark\_duplicates.**picard\_mark\_duplicates** (*sample*, *picard\_mark\_duplicates\_flag*)

Picard tools Mark Duplicates.

**input:** sorted.bam

**output:** \_sorted.rg.md.bam

**citation:** <http://picard.sourceforge.net/>

**link:** <http://picard.sourceforge.net/>

**parameters from parameters file:** BWA\_RESULTS:

TEMP\_DIR:

PICARD\_VERSION:

SAMTOOLS\_VERSION:

## 22.4 GATK Preprocessing

### 22.4.1 WES

omics\_pipe.modules.GATK\_preprocessing\_WES.**GATK\_preprocessing\_WES** (*sample*, *GATK\_preprocessing\_WES\_flag*)

GATK preprocessing steps for whole exome sequencing.

**input:** sorted.rg.md.bam

**output:** .ready.bam

**citation:** McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernyttsky A, Garimella K, Altshuler D, Gabriel S, Daly M, DePristo MA (2010). The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res.* 20:1297-303.

**link:** <http://www.broadinstitute.org/gatk/>

**parameters from parameters file:** BWA\_RESULTS:

TEMP\_DIR:

GATK\_VERSION:

GENOME:

DBSNP:

MILLS:

G1000:

CAPTURE\_KIT\_BED:

SAMTOOLS\_VERSION:

## 22.4.2 WGS

omics\_pipe.modules.GATK\_preprocessing\_WGS.**GATK\_preprocessing\_WGS** (*sample*,  
*GATK\_preprocessing\_WGS\_flag*)

GATK preprocessing steps for whole genome sequencing.

**input:** sorted.rg.md.bam

**output:** .ready.bam

**citation:** McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytzky A, Garimella K, Altshuler D, Gabriel S, Daly M, DePristo MA (2010). The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res.* 20:1297-303.

**link:** <http://www.broadinstitute.org/gatk/>

**parameters from parameters file:** BWA\_RESULTS:

TEMP\_DIR:

GATK\_VERSION:

GENOME:

DBSNP:

MILLS:

G1000:

SAMTOOLS\_VERSION:

## 22.5 GATK Variant Discovery

omics\_pipe.modules.GATK\_variant\_discovery.**GATK\_variant\_discovery** (*sample*,  
*GATK\_variant\_discovery\_flag*)

GATK\_variant\_discovery.

**input:** sorted.rg.md.bam

**output:** .ready.bam

**citation:** McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytzky A, Garimella K, Altshuler D, Gabriel S, Daly M, DePristo MA (2010). The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res.* 20:1297-303.

**link:** **GATK\_variant\_discovery** <http://www.broadinstitute.org/gatk/>

**parameters from parameters file:** BWA\_RESULTS:

TEMP\_DIR:

GATK\_VERSION:

GENOME:

DBSNP:

VARIANT\_RESULTS:

## 22.6 GATK Variant Filtering

omics\_pipe.modules.GATK\_variant\_filtering.**GATK\_variant\_filtering** (*sample*,  
GATK\_variant\_filtering\_flag)

GATK\_variant\_filtering.

**input:** sorted.rg.md.bam

**output:** .ready.bam

**citation:** McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, Garimella K, Altshuler D, Gabriel S, Daly M, DePristo MA (2010). The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res.* 20:1297-303.

**link:** **GATK\_variant\_filtering** <http://www.broadinstitute.org/gatk/>

**parameters from parameters file:** VARIANT\_RESULTS:

TEMP\_DIR:

GATK\_VERSION:

GENOME:

DBSNP:

MILLS:

OMNI:

HAPMAP:

R\_VERSION:

G1000\_SNPs:

G1000\_Indels:

omics\_pipe.modules.GATK\_variant\_filtering.**GATK\_variant\_filtering\_group** (*sample*,  
GATK\_variant\_filtering\_group)

GATK\_variant\_filtering.

**input:** sorted.rg.md.bam

**output:** .ready.bam

**citation:** McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, Garimella K, Altshuler D, Gabriel S, Daly M, DePristo MA (2010). The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res.* 20:1297-303.

**link:** **GATK\_variant\_filtering** <http://www.broadinstitute.org/gatk/>

parameters from parameters file:

VARIANT\_RESULTS:

TEMP\_DIR:

GATK\_VERSION:

GENOME:

DBSNP:

MILLS\_G1000:

OMNI:



HAPMAP:  
R\_VERSION:  
G1000:



---

## Whole Genome Sequencing (MUTECT)

---

### 23.1 FASTQC

`omics_pipe.modules.fastqc.fastqc` (*sample*, *fastqc\_flag*)

QC check of raw .fastq files using FASTQC.

**input:** .fastq file

**output:** folder and zipped folder containing html, txt and image files

**citation:** Babraham Bioinformatics

**link:** <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

**parameters from parameters file:** RAW\_DATA\_DIR:

QC\_PATH:

FASTQC\_VERSION:

COMPRESSION:

### 23.2 BWA-MEM

`omics_pipe.modules.bwa.bwa_mem` (*sample*, *bwa\_mem\_flag*)

BWA aligner with BWA-MEM algorithm.

**input:** .fastq

**output:** .sam

**citation:** Li H. and Durbin R. (2009) Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics*, 25, 1754-1760. [PMID: 19451168]

**link:** <http://bio-bwa.sourceforge.net/bwa.shtml>

**parameters from parameters file:** BWA\_RESULTS:

TEMP\_DIR:

SAMTOOLS\_VERSION:

BWA\_VERSION:

GENOME:

RAW\_DATA\_DIR:

BWA\_OPTIONS:

COMPRESSION:

## 23.3 MUTECT

omics\_pipe.modules.mutect.**mutect** (*sample, mutect\_flag*)

Runs MuTect on paired tumor/normal samples to detect somatic point mutations in cancer genomes.

**input:** .bam

**output:** call\_stats.txt

**citation:** Cibulskis, K. et al. Sensitive detection of somatic point mutations in impure and heterogeneous cancer samples. Nat Biotechnology (2013).doi:10.1038/nbt.2514

**link:** <http://www.broadinstitute.org/cancer/cga/mutect>

**parameters from parameters file:** BWA\_RESULTS:

TEMP\_DIR:

GATK\_VERSION:

GENOME:

DBSNP:

MILLS:

G1000:

CAPTURE\_KIT\_BED:

---

## All Available Modules

---

Below are all available modules in the current release of Omics Pipe in alphabetical order. When creating a custom pipeline, you can choose from these modules or create your own.

`omics_pipe.modules.annotate_peaks.annotate_peaks` (*step, annotate\_peaks\_flag*)

Runs HOMER to annotate peak track from CHIPseq data.

**input:** .tag input file

**output:** .txt file

**citation:** Heinz S, Benner C, Spann N, Bertolino E et al. Simple Combinations of Lineage-Determining Transcription Factors Prime cis-Regulatory Elements Required for Macrophage and B Cell Identities. Mol Cell 2010 May 28;38(4):576-589. PMID: 20513432

**link:** <http://homer.salk.edu/homer/>

**parameters from parameters file:** PAIR\_LIST:

HOMER\_RESULTS:

HOMER\_VERSION:

TEMP\_DIR:

HOMER\_GENOME:

HOMER\_ANNOTATE\_OPTIONS:

`omics_pipe.modules.annotate_variants.annotate_variants` (*sample, annotate\_variants\_flag*)

Annotates variants with ANNOVAR variant annotator. Follows VarCall. input:

.vcf

**output:** .vcf

**citation:** Wang K, Li M, Hakonarson H. ANNOVAR: Functional annotation of genetic variants from next-generation sequencing data Nucleic Acids Research, 38:e164, 2010

**link:** <http://www.openbioinformatics.org/annovar/>

**parameters from parameters file:** VARIANT\_RESULTS:

ANNOVARDB:

ANNOVAR\_OPTIONS:

ANNOVAR\_OPTIONS2:

TEMP\_DIR:

ANNOVAR\_VERSION:

VCFTOOLS\_VERSION:

omics\_pipe.modules.bowtie.**bowtie** (*sample, bowtie\_flag*)

Runs Bowtie to align .fastq files.

**input:** .fastq file

**output:** sample.bam

**citation:** Langmead B, Trapnell C, Pop M, Salzberg SL. Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. *Genome Biology* 10:R25

**link:** <http://bowtie-bio.sourceforge.net/index.shtml>

**parameters from parameters file:** ENDS:

TRIMMED\_DATA\_PATH:

BOWTIE\_OPTIONS:

BOWTIE\_INDEX:

BOWTIE\_RESULTS:

BOWTIE\_VERSION:

SAMTOOLS\_VERSION:

BEDTOOLS\_VERSION:

omics\_pipe.modules.BreastCancer\_RNA\_report.**BreastCancer\_RNA\_report** (*sample, Breast-Cancer\_RNA\_report\_flag*)

Runs R script with knitr to produce report from RNAseq pipeline.

**input:** results from other steps in RNAseq pipelines

**output:** html report

**citation:**

20. Meissner

**parameters from parameter file:** WORKING\_DIR:

R\_VERSION:

REPORT\_RESULTS:

PARAMS\_FILE:

TABIX\_VERSION:

TUMOR\_TYPE:

GENELIST:

COSMIC:

CLINVAR:

PHARMGKB\_rsID:

PHARMGKB\_Allele:

DRUGBANK:

CADD:

omics\_pipe.modules.bwa.**bwa1** (*sample, bwa1\_flag*)  
BWA aligner for read1 of paired\_end reads.

**input:** .fastq

**output:** .sam

**citation:** Li H. and Durbin R. (2009) Fast and accurate short read alignment with Burrows-Wheeler transform. Bioinformatics, 25, 1754-1760. [PMID: 19451168]

**link:** <http://bio-bwa.sourceforge.net/bwa.shtml>

**parameters from parameters file:** BWA\_RESULTS:

TEMP\_DIR:

SAMTOOLS\_VERSION:

BWA\_VERSION:

BWA\_INDEX:

RAW\_DATA\_DIR:

GATK\_READ\_GROUP\_INFO:

COMPRESSION:

omics\_pipe.modules.bwa.**bwa2** (*sample, bwa2\_flag*)  
BWA aligner for read2 of paired\_end reads.

**input:** .fastq

**output:** .sam

**citation:** Li H. and Durbin R. (2009) Fast and accurate short read alignment with Burrows-Wheeler transform. Bioinformatics, 25, 1754-1760. [PMID: 19451168]

**link:** <http://bio-bwa.sourceforge.net/bwa.shtml>

**parameters from parameters file:** BWA\_RESULTS:

TEMP\_DIR:

SAMTOOLS\_VERSION:

BWA\_VERSION:

BWA\_INDEX:

RAW\_DATA\_DIR:

GATK\_READ\_GROUP\_INFO:

COMPRESSION:

omics\_pipe.modules.bwa.**bwa\_RNA** (*sample, bwa\_flag*)  
BWA aligner for single end reads.

**input:** .fastq

**output:** .sam

**citation:** Li H. and Durbin R. (2009) Fast and accurate short read alignment with Burrows-Wheeler transform. Bioinformatics, 25, 1754-1760. [PMID: 19451168]

**link:** <http://bio-bwa.sourceforge.net/bwa.shtml>

**parameters from parameters file:** BWA\_RESULTS:

TEMP\_DIR:  
SAMTOOLS\_VERSION:  
BWA\_VERSION:  
BWA\_INDEX:  
RAW\_DATA\_DIR:  
GATK\_READ\_GROUP\_INFO:  
COMPRESSION:

omics<sub>pipe</sub>.modules.bwa.**bwa\_mem** (*sample, bwa\_mem\_flag*)

BWA aligner with BWA-MEM algorithm.

**input:** .fastq

**output:** .sam

**citation:** Li H. and Durbin R. (2009) Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics*, 25, 1754-1760. [PMID: 19451168]

**link:** <http://bio-bwa.sourceforge.net/bwa.shtml>

**parameters from parameters file:** BWA\_RESULTS:

TEMP\_DIR:  
SAMTOOLS\_VERSION:  
BWA\_VERSION:  
GENOME:  
RAW\_DATA\_DIR:  
BWA\_OPTIONS:  
COMPRESSION:

omics<sub>pipe</sub>.modules.bwa.**bwa\_mem\_pipe** (*sample, bwa\_mem\_pipe\_flag*)

BWA aligner with BWA-MEM algorithm.

**input:** .fastq

**output:** .sam

**citation:** Li H. and Durbin R. (2009) Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics*, 25, 1754-1760. [PMID: 19451168]

**link:** <http://bio-bwa.sourceforge.net/bwa.shtml>

**parameters from parameters file:** BWA\_RESULTS:

TEMP\_DIR:  
SAMTOOLS\_VERSION:  
BWA\_VERSION:  
GENOME:  
RAW\_DATA\_DIR:  
BWA\_OPTIONS:



COMPRESSION:

SAMBAMBA\_VERSION:

SAMBLASTER\_VERSION:

SAMBAMBA\_OPTIONS:

omics\_pipe.modules.call\_variants.**call\_variants** (*sample, call\_variants\_flag*)

Calls variants from alignment .bam files using Varcall.

**input:** Aligned.out.sort.bam or accepted\_hits.bam

**output:** .vcf file

**citation:** Erik Aronesty (2011). ea-utils : “Command-line tools for processing biological sequencing data”;

**link:** <https://code.google.com/p/ea-utils/wiki/Varcall>

**parameters from parameters file:** STAR\_RESULTS:

GENOME:

VARSCAN\_PATH:

VARSCAN\_OPTIONS:

VARIANT\_RESULTS:

TEMP\_DIR:

SAMTOOLS\_VERSION:

ANNOVAR\_VERSION:

VCFTOOLS\_VERSION:

VARSCAN\_VERSION:

SAMTOOLS\_OPTIONS:

omics\_pipe.modules.ChIP\_trim.**ChIP\_trim** (*sample, ChIP\_trim\_flag*)

Runs Homer Tools to trim adapters from .fastq files.

**input:** .fastq file

**output:** .fastq file

**citation:** Heinz S, Benner C, Spann N, Bertolino E et al. Simple Combinations of Lineage-Determining Transcription Factors Prime cis-Regulatory Elements Required for Macrophage and B Cell Identities. Mol Cell 2010 May 28;38(4):576-589. PMID: 20513432

**link:** <http://homer.salk.edu/homer/>

**parameters from parameters file:** ENDS:

RAW\_DATA\_DIR:

HOMER\_TRIM\_OPTIONS:

TRIMMED\_DATA\_PATH:

HOMER\_VERSION:

omics\_pipe.modules.cuffdiff.**cuffdiff** (*step, cuffdiff\_flag*)

Runs Cuffdiff to perform differential expression. Runs after Cufflinks. Part of Tuxedo Suite.

**input:** .bam files

**output:** differential expression results

**citation:** Trapnell C, et al. Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation Nature Biotechnology doi:10.1038/nbt.1621

**link:** <http://cufflinks.cbc.umd.edu/>

**parameters from parameters file:** CUFFDIFF\_RESULTS:

GENOME:

CUFFDIFF\_OPTIONS:

CUFFMERGE\_RESULTS:

CUFFDIFF\_INPUT\_LIST\_COND1:

CUFFDIFF\_INPUT\_LIST\_COND2:

CUFFLINKS\_VERSION:

`omics_pipe.modules.cuffdiff_miRNA.cuffdiff_miRNA` (*step, cuffdiff\_miRNA\_flag*)

Runs Cuffdiff to perform differential expression. Runs after Cufflinks. Part of Tuxedo Suite.

**input:** .bam files

**output:** differential expression results

**citation:** Trapnell C, et al. Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation Nature Biotechnology doi:10.1038/nbt.1621

**link:** <http://cufflinks.cbc.umd.edu/>

**parameters from parameters file:** CUFFDIFF\_RESULTS:

GENOME:

CUFFDIFF\_OPTIONS:

CUFFMERGE\_RESULTS:

CUFFDIFF\_INPUT\_LIST\_COND1:

CUFFDIFF\_INPUT\_LIST\_COND2:

CUFFLINKS\_VERSION:

`omics_pipe.modules.cufflinks.cufflinks` (*sample, cufflinks\_flag*)

Runs cufflinks to assemble .bam files from TopHat.

**input:** accepted\_hits.bam

**output:** transcripts.gtf

**citation:** Trapnell C, et al. Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation Nature Biotechnology doi:10.1038/nbt.1621

**link:** <http://cufflinks.cbc.umd.edu/>

**parameters from parameters file:** TOPHAT\_RESULTS:

CUFFLINKS\_RESULTS:

REF\_GENES:

GENOME:

CUFFLINKS\_OPTIONS:

CUFFLINKS\_VERSION:

omics\_pipe.modules.cufflinks\_miRNA.**cufflinks\_miRNA** (*sample, cufflinks\_miRNA\_flag*)  
Runs cufflinks to assemble .bam files from TopHat. Takes parameter MIRNA\_GTF.

**input:** accepted\_hits.bam

**output:** transcripts.gtf

**citation:** Trapnell C, et al. Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation Nature Biotechnology doi:10.1038/nbt.1621

**link:** <http://cufflinks.cbcb.umd.edu/>

**parameters from parameters file:** TOPHAT\_RESULTS:

CUFFLINKS\_RESULTS:

miRNA\_GTF:

GENOME:

CUFFLINKS\_OPTIONS:

CUFFLINKS\_VERSION:

omics\_pipe.modules.cufflinks\_ncRNA.**cufflinks\_ncRNA** (*sample, cufflinks\_ncRNA\_flag*)  
Runs cufflinks to assemble .bam files from TopHat. Takes parameters LNCRNA\_GTF and NONCODE\_FASTA.

**input:** accepted\_hits.bam

**output:** transcripts.gtf

**citation:** Trapnell C, et al. Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation Nature Biotechnology doi:10.1038/nbt.1621

**link:** <http://cufflinks.cbcb.umd.edu/>

**parameters from parameters file:** TOPHAT\_RESULTS:

CUFFLINKS\_RESULTS:

LNCRNA\_GTF:

NONCODE\_FASTA:

CUFFLINKS\_OPTIONS:

CUFFLINKS\_VERSION:

omics\_pipe.modules.cuffmerge.**cuffmerge** (*step, cuffmerge\_flag*)  
Runs cuffmerge to merge .gtf files from Cufflinks.

**input:** assembly\_GTF\_list.txt

**output:** merged.gtf

**citation:** Trapnell C, et al. Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation Nature Biotechnology doi:10.1038/nbt.1621

**link:** <http://cufflinks.cbcb.umd.edu/>

**parameters from parameters file:** CUFFMERGE\_RESULTS:

REF\_GENES:

GENOME:

CUFFMERGE\_OPTIONS:

CUFFLINKS\_VERSION:

omics\_pipe.modules.cuffmerge\_miRNA.**cuffmerge\_miRNA** (*step, cuffmerge\_miRNA\_flag*)

Runs cuffmerge to merge .gtf files from Cufflinks.

**input:** assembly\_GTF\_list.txt

**output:** merged.gtf

**citation:** Trapnell C, et al. Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation Nature Biotechnology doi:10.1038/nbt.1621

**link:** <http://cufflinks.cbcb.umd.edu/>

**parameters from parameters file:** CUFFMERGE\_RESULTS:

miRNA\_GTF:

GENOME:

CUFFMERGE\_OPTIONS:

CUFFLINKS\_VERSION:

omics\_pipe.modules.cuffmergetocompare.**cuffmergetocompare** (*step, cuffmergetocompare\_flag*)

Runs cuffcompare to annotate merged .gtf files from Cufflinks.

**input:** assembly\_GTF\_list.txt

**output:** merged.gtf

**citation:** Trapnell C, et al. Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation Nature Biotechnology doi:10.1038/nbt.1621

**link:** <http://cufflinks.cbcb.umd.edu/>

**parameters from parameters file:** CUFFMERGE\_RESULTS:

REF\_GENES:

GENOME:

CUFFMERGETOCOMPARE\_OPTIONS:

CUFFLINKS\_VERSION:

omics\_pipe.modules.cuffmergetocompare\_miRNA.**cuffmergetocompare\_miRNA** (*step, cuffmergetocompare\_miRNA\_flag*)

Runs cuffcompare to annotate merged .gtf files from Cufflinks.

**input:** assembly\_GTF\_list.txt

**output:** merged.gtf

**citation:** Trapnell C, et al. Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation Nature Biotechnology doi:10.1038/nbt.1621

**link:** <http://cufflinks.cbcb.umd.edu/>

**parameters from parameters file:** CUFFMERGE\_RESULTS:

miRNA\_GTF:

GENOME:

CUFFMERGETOCOMPARE\_OPTIONS:

CUFFLINKS\_VERSION:

omics\_pipe.modules.custom\_R\_report.**custom\_R\_report** (*sample, custom\_R\_report\_flag*)  
Runs R script with knitr to produce report from omics pipeline.

**input:** results from other steps in RNAseq pipelines

**output:** html report

**citation:**

20. Meissner

**parameters from parameter file:** REPORT\_SCRIPT:

R\_VERSION:

REPORT\_RESULTS:

R\_MARKUP\_FILE:

DPS\_VERSION:

PARAMS\_FILE:

omics\_pipe.modules.cutadapt\_miRNA.**cutadapt\_miRNA** (*sample, cutadapt\_miRNA\_flag*)  
Runs Cutadapt to trim adapters from reads.

**input:** .fastq

**output:** .fastq

**citation:** Martin 2011. Cutadapt removes adapter sequences from high-throughput sequencing reads. EMBnet.journal 17: 10-12.

**link:** <https://code.google.com/p/cutadapt/>

**parameters from parameters file:** RAW\_DATA\_DIR:

ADAPTER:

TRIMMED\_DATA\_PATH:

PYTHON\_VERSION

omics\_pipe.modules.fastq\_length\_filter\_miRNA.**fastq\_length\_filter\_miRNA** (*sample, fastq\_length\_filter\_miRNA\_f*)

Runs custom Python script to filter miRNA reads by length.

**input:** .fastq

**output:** .fastq

**parameters from parameter file:** TRIMMED\_DATA\_PATH:

omics\_pipe.modules.fastqc.**fastqc** (*sample, fastqc\_flag*)  
QC check of raw .fastq files using FASTQC.

**input:** .fastq file

**output:** folder and zipped folder containing html, txt and image files

**citation:** Babraham Bioinformatics

**link:** <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

**parameters from parameters file:** RAW\_DATA\_DIR:

QC\_PATH:

FASTQC\_VERSION:

COMPRESSION:

omics<sub>pipe</sub>.modules.fastqc\_miRNA.**fastqc\_miRNA** (*sample, fastqc\_miRNA\_flag*)  
QC check of raw .fastq files using FASTQC.

**input:** .fastq file

**output:** folder and zipped folder containing html, txt and image files

**citation:** Babraham Bioinformatics

**link:** <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

**parameters from parameters file:** RAW\_DATA\_DIR:

QC\_PATH:

FASTQC\_VERSION:

omics<sub>pipe</sub>.modules.filter\_variants.**filter\_variants** (*sample, filter\_variants\_flag*)  
Filters variants to remove common variants.

**input:** .bam or .sam file

**output:** .vcf file

**citation:** Piskol et al. 2013. Reliable identification of genomic variants from RNA-seq data. The American Journal of Human Genetics 93: 641-651.

**link:** <http://lilab.stanford.edu/SNPiR/>

**parameters from parameters file:** VARIANT\_RESULTS:

TEMP\_DIR:

SAMTOOLS\_VERSION:

BWA\_VERSION:

PICARD\_VERSION:

GATK\_VERSION:

BEDTOOLS\_VERSION:

UCSC\_TOOLS\_VERSION:

GENOME:

REPEAT\_MASKER:

SNPIR\_ANNOTATION:

RNA\_EDIT:

DBSNP:

MILLS:

G1000:

WORKING\_DIR:

BWA\_RESULTS:

SNPIR\_VERSION:

SNPIR\_CONFIG:

SNPIR\_DIR:

SNPEFF\_VERSION:

dbNSFP:

VCFTOOLS\_VERSION:

WORKING\_DIR:

SNP\_FILTER\_OUT\_REF:

`omics_pipe.modules.find_motifs.find_motifs` (*step, find\_motifs\_flag*)

Runs HOMER to find motifs from ChIPseq data.

**input:** .txt peak file from Homer

**output:** .txt file

**citation:** Heinz S, Benner C, Spann N, Bertolino E et al. Simple Combinations of Lineage-Determining Transcription Factors Prime cis-Regulatory Elements Required for Macrophage and B Cell Identities. Mol Cell 2010 May 28;38(4):576-589. PMID: 20513432

**link:** <http://homer.salk.edu/homer/>

**parameters from parameters file:** PAIR\_LIST:

HOMER\_RESULTS:

HOMER\_VERSION:

TEMP\_DIR:

HOMER\_GENOME:

HOMER\_MOTIFS\_OPTIONS:

`omics_pipe.modules.fusion_catcher.fusion_catcher` (*sample, fusion\_catcher\_flag*)

Detects fusion genes in paired-end RNAseq data.

**input:** paired end .fastq files

**output:** list of candidate fusion genes

**citation:**

19. Kangaspeska, S. Hultsch, H. Edgren, D. Nicorici, A. Murumgi, O.P. Kallioniemi, Reanalysis of RNA-sequencing data reveals several additional fusion genes with multiple isoforms, PLOS One, Oct. 2012. <http://dx.plos.org/10.1371/journal.pone.0048745>

**link:** <https://code.google.com/p/fusioncatcher>

**parameters from parameters file:** ENDS:

RAW\_DATA\_DIR:

FUSION\_RESULTS:

FUSIONCATCHERBUILD\_DIR:

TEMP\_DIR:

SAMTOOLS\_VERSION:

FUSIONCATCHER\_VERSION:

FUSIONCATCHER\_OPTIONS:

TISSUE:

PYTHON\_VERSION:

omics\_pipe.modules.GATK\_preprocessing\_WES.**GATK\_preprocessing\_WES** (*sample,*  
*GATK\_preprocessing\_WES\_flag*)

GATK preprocessing steps for whole exome sequencing.

**input:** sorted.rg.md.bam

**output:** .ready.bam

**citation:** McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, Garimella K, Altshuler D, Gabriel S, Daly M, DePristo MA (2010). The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res.* 20:1297-303.

**link:** <http://www.broadinstitute.org/gatk/>

**parameters from parameters file:** BWA\_RESULTS:

TEMP\_DIR:

GATK\_VERSION:

GENOME:

DBSNP:

MILLS:

G1000:

CAPTURE\_KIT\_BED:

SAMTOOLS\_VERSION:

omics\_pipe.modules.GATK\_preprocessing\_WGS.**GATK\_preprocessing\_WGS** (*sample,*  
*GATK\_preprocessing\_WGS\_flag*)

GATK preprocessing steps for whole genome sequencing.

**input:** sorted.rg.md.bam

**output:** .ready.bam

**citation:** McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, Garimella K, Altshuler D, Gabriel S, Daly M, DePristo MA (2010). The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res.* 20:1297-303.

**link:** <http://www.broadinstitute.org/gatk/>

**parameters from parameters file:** BWA\_RESULTS:

TEMP\_DIR:

GATK\_VERSION:

GENOME:

DBSNP:

MILLS:

G1000:

SAMTOOLS\_VERSION:

omics\_pipe.modules.GATK\_variant\_discovery.**GATK\_variant\_discovery** (*sample,*  
*GATK\_variant\_discovery\_flag*)

GATK\_variant\_discovery.

**input:** sorted.rg.md.bam

**output:** .ready.bam



**citation:** McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, Garimella K, Altshuler D, Gabriel S, Daly M, DePristo MA (2010). The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res.* 20:1297-303.

**link:** [GATK\\_variant\\_discovery](http://www.broadinstitute.org/gatk/) <http://www.broadinstitute.org/gatk/>

**parameters from parameters file:** BWA\_RESULTS:

TEMP\_DIR:

GATK\_VERSION:

GENOME:

DBSNP:

VARIANT\_RESULTS:

omics\_pipe.modules.GATK\_variant\_filtering.**GATK\_variant\_filtering** (*sample,*  
*GATK\_variant\_filtering\_flag*)

GATK\_variant\_filtering.

**input:** sorted.rg.md.bam

**output:** .ready.bam

**citation:** McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, Garimella K, Altshuler D, Gabriel S, Daly M, DePristo MA (2010). The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res.* 20:1297-303.

**link:** [GATK\\_variant\\_filtering](http://www.broadinstitute.org/gatk/) <http://www.broadinstitute.org/gatk/>

**parameters from parameters file:** VARIANT\_RESULTS:

TEMP\_DIR:

GATK\_VERSION:

GENOME:

DBSNP:

MILLS:

OMNI:

HAPMAP:

R\_VERSION:

G1000\_SNPs:

G1000\_Indels:

omics\_pipe.modules.GATK\_variant\_filtering.**GATK\_variant\_filtering\_group** (*sample,*  
*GATK\_variant\_filtering\_group*)

GATK\_variant\_filtering.

**input:** sorted.rg.md.bam

**output:** .ready.bam

**citation:** McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, Garimella K, Altshuler D, Gabriel S, Daly M, DePristo MA (2010). The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res.* 20:1297-303.

**link:** [GATK\\_variant\\_filtering](http://www.broadinstitute.org/gatk/) <http://www.broadinstitute.org/gatk/>

parameters from parameters file:

VARIANT\_RESULTS:

TEMP\_DIR:

GATK\_VERSION:

GENOME:

DBSNP:

MILLS\_G1000:

OMNI:

HAPMAP:

R\_VERSION:

G1000:

omics<sub>pipe</sub>.modules.homer\_peaks.**homer\_peaks** (*step, homer\_peaks\_flag*)

Runs HOMER to call peaks from ChIPseq data.

**input:** .tag input file

**output:** .txt file

**citation:** Heinz S, Benner C, Spann N, Bertolino E et al. Simple Combinations of Lineage-Determining Transcription Factors Prime cis-Regulatory Elements Required for Macrophage and B Cell Identities. Mol Cell 2010 May 28;38(4):576-589. PMID: 20513432

**link:** <http://homer.salk.edu/homer/>

**parameters from parameters file:** PAIR\_LIST:

HOMER\_RESULTS:

HOMER\_PEAKE\_OPTIONS:

HOMER\_VERSION:

TEMP\_DIR:

omics<sub>pipe</sub>.modules.htseq.**htseq** (*sample, htseq\_flag*)

Runs htseq-count to get raw count data from alignments.

**input:** Aligned.out.sort.bam

**output:** counts.txt

**citation:** Simon Anders, EMBL

**link:** <http://www-huber.embl.de/users/anders/HTSeq/doc/overview.html>

**parameters from parameters file:** STAR\_RESULTS:

HTSEQ\_OPTIONS:

REF\_GENES:

HTSEQ\_RESULTS:

TEMP\_DIR:

SAMTOOLS\_VERSION:

BAM\_FILE\_NAME:

PYTHON\_VERSION:

omics\_pipe.modules.htseq\_gencode.**htseq\_gencode** (*sample*, *htseq\_flag*)

Runs htseq-count to get raw count data from alignments.

**input:** Aligned.out.sort.bam

**output:** counts.txt

**citation:** Simon Anders, EMBL

**link:** <http://www-huber.embl.de/users/anders/HTSeq/doc/overview.html>

**parameters from parameters file:** STAR\_RESULTS:

HTSEQ\_OPTIONS:

REF\_GENES\_GENCODE:

HTSEQ\_GENCODE\_RESULTS:

TEMP\_DIR:

SAMTOOLS\_VERSION:

BAM\_FILE\_NAME:

omics\_pipe.modules.htseq\_miRNA.**htseq\_miRNA** (*sample*, *htseq\_miRNA\_flag*)

Runs htseq-count to get raw count data from alignments.

**input:** Aligned.out.sort.bam

**output:** counts.txt

**citation:** Simon Anders, EMBL

**link:** <http://www-huber.embl.de/users/anders/HTSeq/doc/overview.html>

**parameters from parameters file:** TOPHAT\_RESULTS:

HTSEQ\_OPTIONS:

miRNA\_GFF:

HTSEQ\_RESULTS:

TEMP\_DIR:

SAMTOOLS\_VERSION:

BAM\_FILE\_NAME:

omics\_pipe.modules.intogen.**intogen** (*sample*, *intogen\_flag*)

Runs Intogen to rank mutations and implication for cancer phenotype. Follows variant calling.

**input:** .vcf

**output:** variant list

**citation:** Gonzalez-Perez et al. 2013. Intogen mutations identifies cancer drivers across tumor types. Nature Methods 10, 1081-1082.

**link:** <http://www.intogen.org/>

**parameters from parameter file:** VCF\_FILE:

INTOGEN\_OPTIONS:

INTOGEN\_RESULTS:

INTOGEN\_VERSION:

USERNAME:  
WORKING\_DIR:  
TEMP\_DIR:  
SCHEDULER:  
VARIANT\_RESULTS:

omics\_pipe.modules.macs.**macs** (*step, macs\_flag*)

Runs MACS to call peaks from ChIPseq data. input:

.fastq file

**output:** peaks and .bed file

**citation:** Zhang et al. Model-based Analysis of ChIP-Seq (MACS). Genome Biol (2008) vol. 9 (9) pp. R137

**link:** <http://liulab.dfci.harvard.edu/MACS/>

**parameters from parameters file:** PAIR\_LIST:

BOWTIE\_RESULTS:  
CHROM\_SIZES:  
MACS\_RESULTS:  
MACS\_VERSION:  
TEMP\_DIR:  
BEDTOOLS\_VERSION:  
PYTHON\_VERSION:

omics\_pipe.modules.mutect.**mutect** (*sample, mutect\_flag*)

Runs MuTect on paired tumor/normal samples to detect somatic point mutations in cancer genomes.

**input:** .bam

**output:** call\_stats.txt

**citation:** Cibulskis, K. et al. Sensitive detection of somatic point mutations in impure and heterogeneous cancer samples. Nat Biotechnology (2013).doi:10.1038/nbt.2514

**link:** <http://www.broadinstitute.org/cancer/cga/mutect>

**parameters from parameters file:** BWA\_RESULTS:

TEMP\_DIR:  
GATK\_VERSION:  
GENOME:  
DBSNP:  
MILLS:  
G1000:  
CAPTURE\_KIT\_BED:

omics\_pipe.modules.peak\_track.**peak\_track** (*step, peak\_track\_flag*)

Runs HOMER to create peak track from ChIPseq data.

**input:** .tag input file

**output:** .txt file

**citation:** Heinz S, Benner C, Spann N, Bertolino E et al. Simple Combinations of Lineage-Determining Transcription Factors Prime cis-Regulatory Elements Required for Macrophage and B Cell Identities. Mol Cell 2010 May 28;38(4):576-589. PMID: 20513432

**link:** <http://homer.salk.edu/homer/>

**parameters from parameters file:** PAIR\_LIST:

HOMER\_RESULTS:

HOMER\_VERSION:

TEMP\_DIR:

omics\_pipe.modules.picard\_mark\_duplicates.**picard\_mark\_duplicates** (*sample, picard\_mark\_duplicates\_flag*)

Picard tools Mark Duplicates.

**input:** sorted.bam

**output:** \_sorted.rg.md.bam

**citation:** <http://picard.sourceforge.net/>

**link:** <http://picard.sourceforge.net/>

**parameters from parameters file:** BWA\_RESULTS:

TEMP\_DIR:

PICARD\_VERSION:

SAMTOOLS\_VERSION:

omics\_pipe.modules.read\_density.**read\_density** (*sample, read\_density\_flag*)

Runs HOMER to visualize read density from ChIPseq data.

**input:** .bam file

**output:** .txt file

**citation:** Heinz S, Benner C, Spann N, Bertolino E et al. Simple Combinations of Lineage-Determining Transcription Factors Prime cis-Regulatory Elements Required for Macrophage and B Cell Identities. Mol Cell 2010 May 28;38(4):576-589. PMID: 20513432

**link:** <http://homer.salk.edu/homer/>

**parameters from parameters file:** BOWTIE\_RESULTS:

CHROM\_SIZES:

HOMER\_RESULTS:

HOMER\_VERSION:

TEMP\_DIR:

omics\_pipe.modules.RNAseq\_QC.**RNAseq\_QC** (*sample, RNAseq\_QC\_flag*)

Runs rseqc to determine insert size as QC for alignment.

**input:** .bam

**output:** pdf plot

**link:** <http://rseqc.sourceforge.net/>

**parameters from parameters file:** STAR\_RESULTS:

QC\_PATH:  
BAM\_FILE\_NAME:  
RSEQC\_REF:  
TEMP\_DIR:  
PICARD\_VERSION:  
R\_VERSION:

omics<sub>pipe</sub>.modules.RNAseq\_report.**RNAseq\_report** (*sample*, *RNAseq\_report\_flag*)  
Runs R script with knitr to produce report from RNAseq pipeline.

**input:** results from other steps in RNAseq pipelines

**output:** html report

**citation:**

20. Meissner

**parameters from parameter file:** REPORT\_SCRIPT:

R\_VERSION:  
REPORT\_RESULTS:  
R\_MARKUP\_FILE:  
DPS\_VERSION:  
PARAMS\_FILE:

omics<sub>pipe</sub>.modules.RNAseq\_report\_counts.**RNAseq\_report\_counts** (*sample*,  
*RNAseq\_report\_counts\_flag*)

Runs R script with knitr to produce report from RNAseq pipeline.

**input:** results from other steps in RNAseq pipelines

**output:** html report

**citation:**

20. Meissner

**parameters from parameter file:** WORKING\_DIR:

R\_VERSION:  
REPORT\_RESULTS:  
PARAMS\_FILE:

omics<sub>pipe</sub>.modules.RNAseq\_report\_tuxedo.**RNAseq\_report\_tuxedo** (*sample*,  
*RNAseq\_report\_tuxedo\_flag*)

Runs R script with knitr to produce report from RNAseq pipeline.

**input:** results from other steps in RNAseq pipelines

**output:** html report

**citation:**

20. Meissner

**parameters from parameter file:** WORKING\_DIR:

R\_VERSION:  
REPORT\_RESULTS:  
DPS\_VERSION:  
PARAMS\_FILE:

omics\_pipe.modules.rseqc.**rseqc** (*sample, rseqc\_flag*)  
Runs rseqc to determine insert size as QC for alignment.

**input:** .bam

**output:** pdf plot

**link:** <http://rseqc.sourceforge.net/>

**parameters from parameters file:** STAR\_RESULTS:

QC\_PATH:  
BAM\_FILE\_NAME:  
RSEQC\_REF:  
RSEQC\_VERSION:  
TEMP\_DIR:

omics\_pipe.modules.snpir\_variants.**snpir\_variants** (*sample, snpir\_variants\_flag*)  
Calls variants using SNPIR pipeline.

**input:** Aligned.out.sort.bam or accepted\_hits.bam

**output:** final\_variants.vcf file

**citation:** Piskol, R., et al. (2013). "Reliable Identification of Genomic Variants from RNA-Seq Data." The American Journal of Human Genetics 93(4): 641-651.

**link:** <http://lilab.stanford.edu/SNPiR/>

**parameters from parameters file:** VARIANT\_RESULTS:

TEMP\_DIR:  
SAMTOOLS\_VERSION:  
BWA\_VERSION:  
PICARD\_VERSION:  
GATK\_VERSION:  
BEDTOOLS\_VERSION:  
UCSC\_TOOLS\_VERSION:  
GENOME:  
REPEAT\_MASKER:  
SNPIR\_ANNOTATION:  
RNA\_EDIT:  
DBSNP:  
MILLS:

G1000:  
WORKING\_DIR:  
BWA\_RESULTS:  
SNPIR\_VERSION:  
SNPIR\_CONFIG:  
SNPIR\_DIR:  
ENCODING:

omics<sub>pipe</sub>.modules.star.star (sample, star\_flag)

Runs STAR to align .fastq files.

**input:** .fastq file

**output:** Aligned.out.bam

**citation:**

1. Dobin et al, Bioinformatics 2012; doi: 10.1093/bioinformatics/bts635 “STAR: ultrafast universal RNA-seq aligner”

**link:** <https://code.google.com/p/rna-star/>

**parameters from parameters file:** ENDS:

RAW\_DATA\_DIR:  
STAR\_INDEX:  
STAR\_OPTIONS:  
STAR\_RESULTS:  
SAMTOOLS\_VERSION:  
STAR\_VERSION:

omics<sub>pipe</sub>.modules.star\_piRNA.star\_piRNA (sample, star\_flag)

Runs STAR to align .fastq files.

**input:** .fastq file

**output:** Aligned.out.bam

**citation:**

1. Dobin et al, Bioinformatics 2012; doi: 10.1093/bioinformatics/bts635 “STAR: ultrafast universal RNA-seq aligner”

**link:** <https://code.google.com/p/rna-star/>

**parameters from parameters file:** ENDS:

RAW\_DATA\_DIR:  
STAR\_INDEX:  
STAR\_OPTIONS:  
STAR\_RESULTS:  
SAMTOOLS\_VERSION:  
STAR\_VERSION:



omics\_pipe.modules.TCGA\_download.**TCGA\_download** (*sample*, *TCGA\_download\_flag*)

Downloads and unzips TCGA data from Manifest.xml downloaded from CGHub. input:

TCGA XML file

**output:** downloaded files from TCGA

**citation:** The Cancer Genome Atlas

**link:** <https://cghub.ucsc.edu/software/downloads.html>

**parameters from parameters file:** TCGA\_XML\_FILE:

TCGA\_KEY:

TCGA\_OUTPUT\_PATH:

CGATOOLS\_VERSION:

omics\_pipe.modules.tophat.**tophat** (*sample*, *tophat\_flag*)

Runs TopHat to align .fastq files.

**input:** .fastq file

**output:** accepted\_hits.bam

**citation:** Trapnell C, Pachter L, Salzberg SL. TopHat: discovering splice junctions with RNA-Seq. Bioinformatics doi:10.1093/bioinformatics/btp120

**link:** <http://tophat.cbcb.umd.edu/>

**parameters from parameters file:** RAW\_DATA\_DIR:

REF\_GENES:

TOPHAT\_RESULTS:

BOWTIE\_INDEX:

TOPHAT\_VERSION:

TOPHAT\_OPTIONS:

BOWTIE\_VERSION:

SAMTOOLS\_VERSION:

omics\_pipe.modules.tophat\_miRNA.**tophat\_miRNA** (*sample*, *tophat\_miRNA\_flag*)

Runs TopHat to align .fastq files.

**input:** .fastq file

**output:** accepted\_hits.bam

**citation:** Trapnell C, Pachter L, Salzberg SL. TopHat: discovering splice junctions with RNA-Seq. Bioinformatics doi:10.1093/bioinformatics/btp120

**link:** <http://tophat.cbcb.umd.edu/>

**parameters from parameters file:** RAW\_DATA\_DIR:

miRNA\_GTF:

TOPHAT\_RESULTS:

miRNA\_BOWTIE\_INDEX:

TOPHAT\_VERSION:

TOPHAT\_OPTIONS:

BOWTIE\_VERSION:

SAMTOOLS\_VERSION:

omics<sub>pipe</sub>.modules.tophat\_ncRNA.tophat\_ncRNA(*sample, tophat\_ncRNA\_flag*)

Runs TopHat to align .fastq files.

**input:** .fastq file

**output:** accepted\_hits.bam

**citation:** Trapnell C, Pachter L, Salzberg SL. TopHat: discovering splice junctions with RNA-Seq. Bioinformatics doi:10.1093/bioinformatics/btp120

**link:** <http://tophat.cbcb.umd.edu/>

**parameters from parameters file:** RAW\_DATA\_DIR:

REF\_GENES:

TOPHAT\_RESULTS:

NONCODE\_BOWTIE\_INDEX:

TOPHAT\_VERSION:

TOPHAT\_OPTIONS:

BOWTIE\_VERSION:

SAMTOOLS\_VERSION:

---

## Version History

---

### 25.1 1.1.3 (2014/08/22)

\*New release for PyPi with bug fixes

### 25.2 1.1.2b (2014/08/05)

#### 25.2.1 New Features

- Added support for latest GATK version
- Added GATK Group Variant Calling pipeline
- Added noncoding RNA HTseq module

#### 25.2.2 Bug Fixes

- AMI memory handling
- Fixed Sumatra parameters file handling
- RNAseq count based pipeline produces single report for all samples

### 25.3 1.1.0 (2014/07/09)

First public release!

### 25.4 1.0.16 (2014/07/01)

#### 25.4.1 New Features

- Added Sumatra for logging of each run to Sumatra database
- Added reading samples from a text file
- Added AWS distribution and docs

## 25.4.2 Bug Fixes

- Import modules for utils.py sumatra and hgapi

## 25.5 1.0.15 (2014/06/20)

### 25.5.1 New Features

- Added TCGA download support

---

## Copyright & License

---

### 26.1 Omics Pipe

MIT License (MIT)

Copyright (c) 2013 Kathleen Marie Fisch

Permission is hereby granted, free of charge, to any person obtaining a copy of this software and associated documentation files (the “Software”), to deal in the Software without restriction, including without limitation the rights to use, copy, modify, merge, publish, distribute, sublicense, and/or sell copies of the Software, and to permit persons to whom the Software is furnished to do so, subject to the following conditions:

The above copyright notice and this permission notice shall be included in all copies or substantial portions of the Software.

THE SOFTWARE IS PROVIDED “AS IS”, WITHOUT WARRANTY OF ANY KIND, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT. IN NO EVENT SHALL THE AUTHORS OR COPYRIGHT HOLDERS BE LIABLE FOR ANY CLAIM, DAMAGES OR OTHER LIABILITY, WHETHER IN AN ACTION OF CONTRACT, TORT OR OTHERWISE, ARISING FROM, OUT OF OR IN CONNECTION WITH THE SOFTWARE OR THE USE OR OTHER DEALINGS IN THE SOFTWARE.



**O**

omics\_pipe.modules.annotate\_peaks, 95  
omics\_pipe.modules.annotate\_variants,  
95  
omics\_pipe.modules.bowtie, 96  
omics\_pipe.modules.BreastCancer\_RNA\_report,  
96  
omics\_pipe.modules.bwa, 93  
omics\_pipe.modules.call\_variants, 99  
omics\_pipe.modules.ChIP\_trim, 99  
omics\_pipe.modules.cuffdiff, 76  
omics\_pipe.modules.cuffdiff\_miRNA, 100  
omics\_pipe.modules.cufflinks, 74  
omics\_pipe.modules.cufflinks\_miRNA, 100  
omics\_pipe.modules.cufflinks\_ncRNA, 101  
omics\_pipe.modules.cuffmerge, 75  
omics\_pipe.modules.cuffmerge\_miRNA, 101  
omics\_pipe.modules.cuffmergetocompare,  
75  
omics\_pipe.modules.cuffmergetocompare\_miRNA,  
102  
omics\_pipe.modules.custom\_R\_report, 102  
omics\_pipe.modules.cutadapt\_miRNA, 77  
omics\_pipe.modules.fastq\_length\_filter\_miRNA,  
77  
omics\_pipe.modules.fastqc, 93  
omics\_pipe.modules.fastqc\_miRNA, 104  
omics\_pipe.modules.filter\_variants, 104  
omics\_pipe.modules.find\_motifs, 105  
omics\_pipe.modules.fusion\_catcher, 105  
omics\_pipe.modules.GATK\_preprocessing\_WES,  
105  
omics\_pipe.modules.GATK\_preprocessing\_WGS,  
106  
omics\_pipe.modules.GATK\_variant\_discovery,  
106  
omics\_pipe.modules.GATK\_variant\_filtering,  
107  
omics\_pipe.modules.homer\_peaks, 108  
omics\_pipe.modules.htseq, 78  
omics\_pipe.modules.htseq\_gencode, 108  
omics\_pipe.modules.htseq\_miRNA, 109  
omics\_pipe.modules.intogen, 109  
omics\_pipe.modules.macs, 110  
omics\_pipe.modules.mutect, 94  
omics\_pipe.modules.peak\_track, 110  
omics\_pipe.modules.picard\_mark\_duplicates,  
111  
omics\_pipe.modules.read\_density, 111  
omics\_pipe.modules.RNAseq\_QC, 111  
omics\_pipe.modules.RNAseq\_report, 112  
omics\_pipe.modules.RNAseq\_report\_counts,  
79  
omics\_pipe.modules.RNAseq\_report\_tuxedo,  
76  
omics\_pipe.modules.rseqc, 113  
omics\_pipe.modules.snpir\_variants, 113  
omics\_pipe.modules.star, 78  
omics\_pipe.modules.star\_piRNA, 114  
omics\_pipe.modules.TCGA\_download, 114  
omics\_pipe.modules.tophat, 74  
omics\_pipe.modules.tophat\_miRNA, 115  
omics\_pipe.modules.tophat\_ncRNA, 116





**A**

annotate\_peaks() (in module omics\_pipe.modules.annotate\_peaks), 83, 95

annotate\_variants() (in module omics\_pipe.modules.annotate\_variants), 95

**B**

bowtie() (in module omics\_pipe.modules.bowtie), 82, 86, 96

BreastCancer\_RNA\_report() (in module omics\_pipe.modules.BreastCancer\_RNA\_report), 59, 67, 96

bwa1() (in module omics\_pipe.modules.bwa), 55, 64, 97

bwa2() (in module omics\_pipe.modules.bwa), 56, 64, 97

bwa\_mem() (in module omics\_pipe.modules.bwa), 87, 93, 98

bwa\_mem\_pipe() (in module omics\_pipe.modules.bwa), 98

bwa\_RNA() (in module omics\_pipe.modules.bwa), 56, 64, 97

**C**

call\_variants() (in module omics\_pipe.modules.call\_variants), 99

ChIP\_trim() (in module omics\_pipe.modules.ChIP\_trim), 81, 85, 99

cuffdiff() (in module omics\_pipe.modules.cuffdiff), 49, 76, 99

cuffdiff\_miRNA() (in module omics\_pipe.modules.cuffdiff\_miRNA), 100

cufflinks() (in module omics\_pipe.modules.cufflinks), 48, 74, 100

cufflinks\_miRNA() (in module omics\_pipe.modules.cufflinks\_miRNA), 100

cufflinks\_ncRNA() (in module omics\_pipe.modules.cufflinks\_ncRNA), 101

cuffmerge() (in module omics\_pipe.modules.cuffmerge), 48, 75, 101

cuffmerge\_miRNA() (in module omics\_pipe.modules.cuffmerge\_miRNA), 101

cuffmergetocompare() (in module omics\_pipe.modules.cuffmergetocompare), 49, 75, 102

cuffmergetocompare\_miRNA() (in module omics\_pipe.modules.cuffmergetocompare\_miRNA), 102

custom\_R\_report() (in module omics\_pipe.modules.custom\_R\_report), 102

cutadapt\_miRNA() (in module omics\_pipe.modules.cutadapt\_miRNA), 73, 77, 103

**F**

fastq\_length\_filter\_miRNA() (in module omics\_pipe.modules.fastq\_length\_filter\_miRNA), 73, 77, 103

fastqc() (in module omics\_pipe.modules.fastqc), 47, 51, 53, 61, 69, 73, 77, 81, 85, 87, 93, 103

fastqc\_miRNA() (in module omics\_pipe.modules.fastqc\_miRNA), 104

filter\_variants() (in module omics\_pipe.modules.filter\_variants), 57, 66, 104

find\_motifs() (in module omics\_pipe.modules.find\_motifs), 84, 105

fusion\_catcher() (in module omics\_pipe.modules.fusion\_catcher), 55, 63, 105

**G**

GATK\_preprocessing\_WES() (in module omics\_pipe.modules.GATK\_preprocessing\_WES), 88, 105

GATK\_preprocessing\_WGS() (in module omics\_pipe.modules.GATK\_preprocessing\_WGS), 89, 106

GATK\_variant\_discovery() (in module omics\_pipe.modules.GATK\_variant\_discovery),

89, 106  
 GATK\_variant\_filtering() (in module omics\_pipe.modules.GATK\_variant\_filtering), 90, 107  
 GATK\_variant\_filtering\_group() (in module omics\_pipe.modules.GATK\_variant\_filtering), 90, 107

## H

homer\_peaks() (in module omics\_pipe.modules.homer\_peaks), 82, 108  
 htseq() (in module omics\_pipe.modules.htseq), 52, 54, 62, 70, 78, 108  
 htseq\_gencode() (in module omics\_pipe.modules.htseq\_gencode), 108  
 htseq\_miRNA() (in module omics\_pipe.modules.htseq\_miRNA), 109

## I

intogen() (in module omics\_pipe.modules.intogen), 58, 67, 109

## M

macs() (in module omics\_pipe.modules.macs), 86, 110  
 mutect() (in module omics\_pipe.modules.mutect), 94, 110

## O

omics\_pipe  
   about, 2  
   Adding Modules, 28  
   AWS\_installation, 8  
   Custom Pipeline, 23  
   history, 116  
   installation, 5  
   license, 118  
   parameters, 18  
   pipelines, 31  
   tutorial, 16  
 omics\_pipe.modules.annotate\_peaks (module), 83, 95  
 omics\_pipe.modules.annotate\_variants (module), 95  
 omics\_pipe.modules.bowtie (module), 82, 86, 96  
 omics\_pipe.modules.BreastCancer\_RNA\_report (module), 59, 67, 96  
 omics\_pipe.modules.bwa (module), 55, 64, 87, 93, 97  
 omics\_pipe.modules.call\_variants (module), 99  
 omics\_pipe.modules.ChIP\_trim (module), 81, 85, 99  
 omics\_pipe.modules.cuffdiff (module), 49, 76, 99  
 omics\_pipe.modules.cuffdiff\_miRNA (module), 100  
 omics\_pipe.modules.cufflinks (module), 48, 74, 100  
 omics\_pipe.modules.cufflinks\_miRNA (module), 100  
 omics\_pipe.modules.cufflinks\_ncRNA (module), 101  
 omics\_pipe.modules.cuffmerge (module), 48, 75, 101

omics\_pipe.modules.cuffmerge\_miRNA (module), 101  
 omics\_pipe.modules.cuffmergetocompare (module), 49, 75, 102  
 omics\_pipe.modules.cuffmergetocompare\_miRNA (module), 102  
 omics\_pipe.modules.custom\_R\_report (module), 102  
 omics\_pipe.modules.cutadapt\_miRNA (module), 73, 77, 103  
 omics\_pipe.modules.fastq\_length\_filter\_miRNA (module), 73, 77, 103  
 omics\_pipe.modules.fastqc (module), 47, 51, 53, 61, 69, 73, 77, 81, 85, 87, 93, 103  
 omics\_pipe.modules.fastqc\_miRNA (module), 104  
 omics\_pipe.modules.filter\_variants (module), 57, 66, 104  
 omics\_pipe.modules.find\_motifs (module), 84, 105  
 omics\_pipe.modules.fusion\_catcher (module), 55, 63, 105  
 omics\_pipe.modules.GATK\_preprocessing\_WES (module), 88, 105  
 omics\_pipe.modules.GATK\_preprocessing\_WGS (module), 89, 106  
 omics\_pipe.modules.GATK\_variant\_discovery (module), 89, 106  
 omics\_pipe.modules.GATK\_variant\_filtering (module), 90, 107  
 omics\_pipe.modules.homer\_peaks (module), 82, 108  
 omics\_pipe.modules.htseq (module), 52, 54, 62, 70, 78, 108  
 omics\_pipe.modules.htseq\_gencode (module), 108  
 omics\_pipe.modules.htseq\_miRNA (module), 109  
 omics\_pipe.modules.intogen (module), 58, 67, 109  
 omics\_pipe.modules.macs (module), 86, 110  
 omics\_pipe.modules.mutect (module), 94, 110  
 omics\_pipe.modules.peak\_track (module), 83, 110  
 omics\_pipe.modules.picard\_mark\_duplicates (module), 88, 111  
 omics\_pipe.modules.read\_density (module), 82, 111  
 omics\_pipe.modules.RNAseq\_QC (module), 111  
 omics\_pipe.modules.RNAseq\_report (module), 112  
 omics\_pipe.modules.RNAseq\_report\_counts (module), 52, 71, 79, 112  
 omics\_pipe.modules.RNAseq\_report\_tuxedo (module), 49, 76, 112  
 omics\_pipe.modules.rseqc (module), 54, 63, 113  
 omics\_pipe.modules.snpir\_variants (module), 56, 65, 113  
 omics\_pipe.modules.star (module), 51, 53, 62, 70, 78, 114  
 omics\_pipe.modules.star\_piRNA (module), 114  
 omics\_pipe.modules.TCGA\_download (module), 61, 69, 114  
 omics\_pipe.modules.tophat (module), 47, 74, 115  
 omics\_pipe.modules.tophat\_miRNA (module), 115  
 omics\_pipe.modules.tophat\_ncRNA (module), 116

## P

peak\_track() (in module omics\_pipe.modules.peak\_track), 83, 110  
 picard\_mark\_duplicates() (in module omics\_pipe.modules.picard\_mark\_duplicates), 88, 111

## R

read\_density() (in module omics\_pipe.modules.read\_density), 82, 111  
 RNAseq\_QC() (in module omics\_pipe.modules.RNAseq\_QC), 111  
 RNAseq\_report() (in module omics\_pipe.modules.RNAseq\_report), 112  
 RNAseq\_report\_counts() (in module omics\_pipe.modules.RNAseq\_report\_counts), 52, 71, 79, 112  
 RNAseq\_report\_tuxedo() (in module omics\_pipe.modules.RNAseq\_report\_tuxedo), 49, 76, 112  
 rseqc() (in module omics\_pipe.modules.rseqc), 54, 63, 113

## S

snpir\_variants() (in module omics\_pipe.modules.snpir\_variants), 56, 65, 113  
 star() (in module omics\_pipe.modules.star), 51, 53, 62, 70, 78, 114  
 star\_piRNA() (in module omics\_pipe.modules.star\_piRNA), 114

## T

TCGA\_download() (in module omics\_pipe.modules.TCGA\_download), 61, 69, 114  
 tophat() (in module omics\_pipe.modules.tophat), 47, 74, 115  
 tophat\_miRNA() (in module omics\_pipe.modules.tophat\_miRNA), 115  
 tophat\_ncRNA() (in module omics\_pipe.modules.tophat\_ncRNA), 116