

# EasyDIVER

## NAME

EasyDIVER: Easy pre-processing and Dereplication of In Vitro evolution Reads

## SYNOPSIS

easydiver -i input.directory [-o output.directory] [-p forward.primer] [-q reverse.primer] [-a] [-r] [-T threads] [-e] [-h]

## DESCRIPTION

EasyDIVER converts raw, paired-end, demultiplexed Illumina read files into processed, dereplicated data ready for analysis. The algorithm performs the following:

1. Joins the raw data using [PANDASeg](#).
2. Extracts the insert sequence based on user-supplied primer sequences.
3. Optionally translates into amino acids.
4. Generates counts files.
5. Collects sequence length distributions (histos).
6. Creates a log file.

## OPTIONS

-i input.directory	Required. Input directory path and name. If no value is provided, an error message will be printed in the terminal: <code>ERROR: No input filepath supplied</code> and no further action will be performed.
-o output.directory	Optional. Output directory path and name. If no value is provided, the default value <code>/pipeline.output</code> will be used.
-p forward.primer	Optional. Extraction forward DNA primer. If a forward primer sequence is provided, the pipeline strips out the primer from the start of the sequence. Any sequence before the provided primer will be discarded.
-q reverse.primer	Optional. Extraction reverse DNA primer. If a reverse primer sequence is provided, the pipeline strips out the primer at the start of the sequence. Any sequence after the provided primer will be discarded.
-a	Optional. Translation into amino acids is performed. DNA sequences are translated using the standard genetic code, and the resulting sequences are dereplicated. By default, translation is not performed.
-r	Optional. Files for individual lanes are retained. By default, the script will suppress outputs from individual lanes.
-T threads	Optional. Number of threads used for computation. Default value is 1.
-e	Optional. Additional internal PANDASeg flags. Values must be entered in quotation marks (e.g. <code>-e "-L 50"</code> ). Default value is <code>"-l 1 -d rbfkms"</code> . For more information
-h	If used, a help message will be printed in the terminal and no further action will be performed.

## EXAMPLE

An example of command to run the pipeline using the test data provided in the GitHub repository:

```
easydiver -i ./ -o ./output -p GGCGGAAAGCACATCTGC -T 14 -a -r
```

The above command works if it is:

1. Run from the same directory where the raw data is located.
2. The script `easydiver.sh` has been moved to `/usr/local/bin`, made executable and installed. For information on how to install EasyDIVER, see the README file available in the GitHub repository.

## INPUT

All input files must be:

3. Located in the same directory (even reads from separate lanes).
4. In FASTQ format.
5. Named using the standard Illumina naming scheme: `sample-name_S#_L00#_R#_001.fastq`
6. In either `.fastq` or `.fastq.gz` extensions.

If any of these requirements are not met, the script will not perform as intended, or more likely, outright fail.

## OUTPUT

For each sample, the pipeline combines the reads from every lane, and redirects the outputs to the following sub-directories:

**fastqs** contains the joined fastq files  
**fastas** contains the joined fasta files  
**counts** contains DNA counts files for every sample  
**counts.aa** contains peptide count files for every sample (if translation is requested using the flag `-a`)  
**histos** contain the DNA sequence length distributions and the peptide sequence length distribution (if translation is required)  
**individual.lanes** contains the files (joined fasta files joined fastq files, text counts files and text histograms) corresponding to the individual lanes (if requested using the flag `-r`)

If translation to amino acids is derided (indicated by the use of the flag `-a`) the counts files are translated using the standard genetic code, and the resulting sequences are dereplicated. Count files and length distributions are created for the amino acid sequences as well. All sequence length distributions are redirected to the directory **histos**.

By default, the script will suppress outputs from individual lanes. If you wish to retain the individual lane outputs, use the `-r` flag. If the flag `-r` is used, files corresponding to the individual lanes (joined fasta files joined fastq files, text counts files and text histograms) are retained and redirected to the subdirectory called **individual.lanes**.

For the record, and to monitor the success of the run, a single log text file with the parameters used and the number of sequences in the fastq and counts files is created at the end of the process.