

```
vsearch(1)  
vsearch(1)
```

## USER COMMANDS

## NAME

vsearch – a versatile open-source tool for microbiome analysis, including chimera detection, clustering, dereplication and rereplication, extraction, FASTA/FASTQ/SFF file processing, masking, orienting, pairwise alignment, restriction site cutting, searching, shuffling, sorting, subsampling, and taxonomic classification of amplicon sequences for metagenomics, genomics, and population genetics.

## SYNOPSIS

Chimera detection:  
vsearch (--uchime\_denovo | --uchime2\_denovo | --uchime3\_denovo) fastafilename (--) chimeras |  
--nonchimeras | --uchimealns | --uchimeout) outputfile [options]  
vsearch --uchime\_ref fastafilename (--) chimeras | --nonchimeras | --uchimealns | --uchime-  
out) outputfile --db fastafilename [options]

Clustering:  
vsearch (--) cluster\_fast | --cluster\_size | --cluster\_smallmem | --cluster\_unequal  
fastafilename (--) alnout | --biomout | --blast6out | --centroids | --clusters |  
--mothur\_shared\_out | --msaout | --otutab | --profile | --samout | --uc | --  
userout) outputfile --id real [options]

Dereplication and rereplication:  
vsearch --fastx\_uniques (fastafilename | fastqfilename) (--) fastaout | --fastqout | --  
tabbedout | --uc) outputfile [options]  
vsearch (--) derep\_fulllength | --derep\_id | --derep\_prefix) fastafilename (--) output |  
--uc) outputfile [options]  
vsearch --derep\_smallmem (fastafilename | fastqfilename) --fastaout outputfile [options]  
vsearch --rereplicate fastafilename --output outputfile [options]

Extraction of sequences:  
vsearch --fastx\_getseq fastafilename (--) fastaout | --fastqout | --notmatched | --not-  
matchedfq) outputfile --label label [options]  
vsearch --fastx\_getseqs fastafilename (--) fastaout | --fastqout | --notmatched | --  
not-matchedfq) outputfile (--) label label --labels labelfile | --label\_word label |  
--label\_words labelfile) [options]  
vsearch --fastx\_getsubseq fastafilename (--) fastaout | --fastqout | --notmatched | --  
not-matchedfq) outputfile --label label [--subseq\_start position] [--subseq\_end  
position] [options]

FASTA/FASTQ/SFF file processing:

```
vsearch --fasta2fastq fastafilename --fastqout outputfile [options]
vsearch --fastq_chars fastqfilename [options]
vsearch --fastq_convert fastqfilename --fastqout outputfile [options]
vsearch (--fastq_eestats | --fastq_eestats2) fastqfilename --output outputfile
[options]
```

fastaout\_discarded\_rev |

```
vsearch --fastq_filter fastqfilename [--reverse fastqfilename] (--fastaout | --
fastaout_discarded_rev | --fastqout | --fastqout_discarded --fastaout_rev | --
--fastqout_rev | --fastqout_discarded_rev) outputfile [options]
```

outputfile |

```
vsearch --fastq_join fastqfilename --reverse fastqfilename (--fastaout | --fastqout)
[options]
```

|

```
vsearch --fastq_mergepairs fastqfilename --reverse fastqfilename (--fastaout | --fastqout
--fastaout_notmerged_fwd | --fastaout_notmerged_rev | --fastqout_notmerged_fwd |
--fastqout_notmerged_rev | --etabbedout) outputfile [options]
```

```
vsearch --fastq_stats fastqfilename [--log logfile] [options]
```

fastaout\_discarded\_rev |

```
vsearch --fastx_filter inputfile [--reverse inputfile] (--fastaout | --
fastaout_discarded_rev | --fastqout | --fastqout_discarded --fastaout_rev | --
--fastqout_rev | --fastqout_discarded_rev) outputfile [options]
```

```
vsearch --fastx_revcomp inputfile (--fastaout | --fastqout) outputfile [options]
```

```
vsearch --sff_convert sff-file --fastqout outputfile [options]
```

Masking:

```
vsearch --fastx_mask fastxfilename (--fastaout | --fastqout) outputfile [options]
vsearch --maskfasta fastafilename --output outputfile [options]
```

Orienting:

```
vsearch --orient fastxfilename --db fastxfilename (--fastaout | --fastqout | --notmatched
--tabbedout) outputfile [options]
```

Pairwise alignment:

```
vsearch --allpairs_global fastafilename (--alnout | --blast6out | --matched | --
notmatched | --samout | --uc | --userout) outputfile (--acceptall | --id real) [options]
```

Restriction site cutting:

```
vsearch --cut fastafilename --cut_pattern pattern (--fastaout | --fastaout_rev | --
fastaout_discarded | --fastaout_discarded_rev) outputfile [options]
```

Searching:

```
vsearch --search_exact fastafilename --db fastafilename (--alnout | --biomout | --
blast6out |
```

```
--mothur_shared_out | --otutabout | --samout | --uc | --userout | --lcaout)
outputfile
[options]

biomout |
vsearch --usearch_global fastafile --db (fastafile | udbfile) (--alnout | --
--blast6out | --mothur_shared_out | --otutabout | --samout | --uc | --userout |
--lcaout) outputfile --id real [options]

Shuffling and sorting:
vsearch (--shuffle | --sortbylength | --sortbysize) fastafile --output outputfile
[op-
tions]

Subsampling:
vsearch --fastx_subsample fastafile (--fastaout | --fastqout) outputfile (-- 
sample_pct
real | --sample_size positive integer) [options]

Taxonomic classification:
vsearch --sintax fastafile --db (fastafile | udbfile) --tabbedout outputfile [--
sin-
tax_cutoff real] [options]

UDB database handling:
vsearch --makeudb_usearch fastafile --output outputfile [options]

vsearch --fdb2fasta udbfile --output outputfile [options]

vsearch (--fdbinfo | --fdbstats) udbfile [options]
```

## **DESCRIPTION**

Environmental or clinical molecular diversity studies generate large volumes of amplicons (e.g.; SSU-rRNA sequences) that need to be checked for chimeras, dereplicated, masked, sorted, searched, clustered or compared to reference sequences. The aim of vsearch is to offer a all-in-one open source tool to perform these tasks, using optimized algorithm implementations and harvesting the full potential of modern computers, thus providing fast and accurate data processing.

Comparing nucleotide sequences is at the core of vsearch. To speed up comparisons, vsearch implements an extremely fast Needleman-Wunsch algorithm, making use of the Streaming SIMD Extensions (SSE2) of post-2003 x86-64 CPUs. If SSE2 instructions are not available, vsearch exits with an error message. On Power8 CPUs it will use Altivec/VSX/VMX instructions, and on ARMv8 CPUs it will use Neon instructions. On other systems it can use the SIMD Everywhere library, if available. Memory usage increases rapidly with sequence length: for example comparing two sequences of length 1 kb requires 8 MB of memory per thread, and comparing two 10 kb sequences requires 800 MB of memory per thread. For comparisons involving sequences with a

length product greater than 25 million (for example two sequences of length 5 kb), vsearch uses a slower alignment method described by Hirschberg (1975) and Myers and Miller (1988), with much smaller memory requirements.

#### Input

vsearch accept as input fasta or fastq files containing one or several nucleotidic entries. In fasta files, each entry is made of a header and a sequence. The header is defined as the string comprised between the initial '>' symbol and the first space, tab or the end of the line, unless the --notrunclabels option is in effect, in which case the entire line is included. The header should contain printable ascii characters (33-126). The program will terminate with a fatal error if there are unprintable ascii characters. A warning will be issued if non-ascii characters (128-255) are encountered.

If the header matches the pattern '>[;]size=integer;label', the pattern '>label;size=integer;label', or the pattern '>label;size=integer[;]', vsearch will interpret integer as the number of occurrences (or abundance) of the sequence in the study. That abundance information is used or created during chimera detection, clustering, derePLICATION, sorting and searching.

The sequence is defined as a string of IUPAC symbols (ACGTURYSWKMDBHVN), starting after the end of the identifier line and ending before the next identifier line, or the file end. vsearch silently ignores ascii characters 9 to 13, and exits with an error message if ascii characters 0 to 8, 14 to 31, '.' or '-' are present. All other ascii or non-ascii characters are stripped and complained about in a warning message.

In fastq files, each entry is made of sequence header starting with a symbol '@', a nucleotidic sequence (same rules as for fasta sequences), a quality header starting with a symbol '+' and a string of ASCII characters (offset 33 or 64), each one encoding the quality value of the corresponding position in the nucleotidic sequence.

vsearch operations are case insensitive, except when soft masking is activated. DUST masking is automatically applied during chimera detection, clustering, masking, pairwise alignment and searching. Soft masking is specified with the options '--dbmask soft' (for searching and chimera detection with a reference) or '--qmask soft' (for searching, de novo chimera detection, clustering and masking). When using soft masking, lower case letters indicate masked symbols, while upper case letters indicate regular symbols. Masked symbols are never included

in the unique index words used for sequence comparisons, otherwise they are treated as normal symbols.

When comparing sequences during chimera detection, dereplication, searching and clustering, T

and U are considered identical, regardless of their case. When aligning sequences,

identical

symbols will receive a positive match score (default +2). If two symbols are not identical,

their alignment result in a negative mismatch score (default -4). Aligning a pair of symbols

where at least one of them is an ambiguous symbol (BDHKMNRSVWY) will always result in a score

of zero. Alignment of two identical ambiguous symbols (for example, R vs R) also receives a

score of zero. When computing the amount of similarity by counting matches and mismatches af-

ter alignment, ambiguous nucleotide symbols will count as matching to other symbols if they

have at least one of the nucleotides (ACGTU) they may represent in common. For example:

W will

match A and T, but also any of MRVHDN. When showing alignments (for example with the --

alnout

option) matches involving ambiguous symbols will be shown with a plus character (+)

between

them while exact matches between non-ambiguous symbols will be shown with a vertical bar char-

acter (|).

vsearch can read data from standard files and write to standard files, but it can also read

from pipes and write to pipes! For example, multiple fasta files can be piped into vsearch for

dereplication. To do so, file names can be replaced with:

output - the symbol '--', representing '/dev/stdin' for input files or '/dev/stdout' for

files (with an exception for '--db --', see \* below),

- a named pipe created with the command mkfifo,

- a process substitution '<(command)' as input or '>(command)' as output.

\* --db -- is not accepted, to prevent potential concurrent reads from stdin. A

workaround for advanced users is to call '--db /dev/stdin' directly.

vsearch can automatically read compressed gzip or bzip2 files if the appropriate libraries are

present during the compilation. vsearch can also read pipes streaming compressed gzip or bzip2

data if the options --gzip\_decompress or --bzip2\_decompress are selected. When reading from a pipe, the progress indicator is not updated.

## Options

vsearch recognizes a large number of command-line commands and options. For easier navigation,

options are grouped below by theme (chimera detection, clustering, dereplication and

rereplacement, FASTA/FASTQ file processing, masking, pairwise alignment, searching, shuffling, sorting, and subsampling). We start with the general options that apply to all themes.

**Options**

start with a double dash (--). A single dash (-) may also be used, except on NetBSD systems.

Option names may be shortened as long as they are not ambiguous (e.g. --derep\_f).

Help and version commands:

--help -h  
Display help text with brief information about all commands and options.

--version -v  
Output version information and a citation for the VSEARCH publication. Show the status of the support for gzip- and bzip2-compressed input files.

General options:

--bzip2\_decompress  
When reading from a pipe streaming bzip2-compressed data, decompress the data.  
This option is not needed when reading from a standard bzip2-compressed file.

--fasta\_width positive integer  
Fasta files produced by vsearch are wrapped (sequences are written on lines of integer nucleotides, 80 by default). Set the value to zero to eliminate the wrapping.

--gzip\_decompress  
When reading from a pipe streaming gzip-compressed data, decompress the data.  
This option is not needed when reading from a standard gzip-compressed file.

--label\_suffix string  
When writing FASTA or FASTQ files, add the suffix string to sequence headers.

--log filename  
Write messages to the specified log file. Information written includes program options, start maximum addi- version, amount of memory available, number of cores and command line and if need be, informational messages, warnings and fatal errors. The and finish times are also recorded as well as the elapsed time and the amount of memory consumed. The different vsearch commands can also write tional information to the log file.

--maxseqlength positive integer  
All vsearch operations discard sequences longer than integer (50,000 nu-

cleotides by default).

nucleotide by --minseqlength positive integer  
dereplica- All vsearch operations discard sequences shorter than integer: 1  
--use- default for sorting or shuffling, 32 nucleotides for clustering and  
arch\_global.  
header --no\_progress  
syntax com- Do not show the gradually increasing progress indicator.  
mand.  
error --quiet Suppress all messages to stdout and stderr except for warnings and fatal  
messages.  
string --sample string  
";sam- When writing FASTA or FASTQ files, add the the given sample identifier  
truncated at to sequence headers. For instance, if the given string is ABC, the text  
numerical ple=ABC" will be added to the header. Note that string will be  
the first ';' or blank character. Other characters (alphabetical,  
and punctuations) are accepted.  
should --threads positive integer  
is to Number of computation threads to use (1 to 1024). The number of threads  
following be less than or equal to the number of available CPU cores. The default  
cluster\_size, use all available resources and to launch one thread per core. The  
maskfasta, commands are multi-threaded: allpairs\_global, cluster\_fast,  
is used cluster\_smallmem, cluster\_unicore, fastq\_mergepairs, fastx\_mask,  
search\_exact, syntax, uchime\_ref, and usearch\_global. Only one thread  
for the other commands.

Chimera detection options:

(--dn, Chimera detection is based on a scoring function controlled by five options  
abun- --mindiffs, --mindiv, --minh, --xn). Sequences are first sorted by decreasing  
dance, if available, and compared on their plus strand only (case insensitive).  
options. Input sequences are masked as specified with the --qmask and --hardmask

Mask-

ing of the database for reference based chimera detection is specified with the --db- mask option.

pattern

detection, so

command

In de novo mode, input fasta file must present abundance annotations (i.e. a

[;size=integer[;] in the fasta header). Input order matters for chimera

we recommend to sort sequences by decreasing abundance (default of --

derep\_fulllength command). If your sequence set needs to be sorted, please see the --sortbysize

in the sorting section.

--abskew real

in a  
parents.

process

uchime3\_denovo the

which

their

(80 nu-

Borderline  
which are

may vary

sequences con-  
chimera-free.  
relatives,  
file or  
using the

When using --uchime\_denovo, the abundance skew is used to distinguish

three-way alignment which sequence is the chimera and which are the

The assumption is that chimeras appear later in the PCR amplification

and are therefore less abundant than their parents. For --

default value is 16.0. For the other commands, the default value is 2.0,

means that the parents should be at least 2 times more abundant than

chimera. Any positive value equal or greater than 1.0 can be used.

--alignwidth positive integer

When using --uchimealns, set the width of the three-way alignments

cleotides by default). Set to zero to eliminate wrapping.

--borderline filename

Output borderline chimeric sequences to filename, in fasta format.

chimeric sequences are sequences that have a high enough score but

not sufficiently different from their closest parent.

--chimeras filename

Output chimeric sequences to filename, in fasta format. Output order

when using multiple threads.

--db filename

When using --uchime\_ref, detect chimeras using the reference

tained in filename. Reference sequences are assumed to be

Chimeras cannot be detected if their parents, or sufficiently close

are not present in the database. The file name must refer to a FASTA

to a UDB file. If a UDB file is used, it should be created

--makeudb\_usearch command with the --dbmask dust option.

--dn strictly positive real number  
parameter n pseudo-count prior on the number of no votes, corresponding to the dn re- in the chimera scoring function (default value is 1.4). Increasing -- positives, produces the likelihood of tagging a sequence as a chimera (less false but also more false negatives).

--fasta\_score  
chimeras, Add the chimera score to the headers in the fasta output files for non-chimeras and borderline sequences, using the format  
';uchime\_den- ovo=float;'.

--lengthout  
format by Write sequence length information to the output files in FASTA adding a ";length=integer" attribute in the header.

--mindiffs positive integer  
parameter Minimum number of differences per segment (default value is 3). The is ignored with --uchime2\_denovo and --uchime3\_denovo.

--mindiv real  
parameter Minimum divergence from closest parent (default value is 0.8). The is ignored with --uchime2\_denovo and --uchime3\_denovo.

--minh real  
false Minimum score (h). Increasing this value tends to reduce the number of positives and to decrease sensitivity. Default value is 0.28, and values rang- ing from 0.0 to 1.0 included are accepted. The parameter is ignored with --uchime2\_denovo and --uchime3\_denovo.

--nonchimeras filename  
order may Output non-chimeric sequences to filename, in fasta format. Output vary when using multiple threads.

--relabel string  
to con- Relabel sequences using the prefix string and a ticker (1, 2, 3, etc.) struct the new headers. Use --sizeout to conserve the abundance annotations.

--relabel\_keep  
each se- When relabelling, keep the old identifier in the header after a space.

--relabel\_md5 converted to Relabel sequences using the MD5 message digest algorithm applied to quence. Former sequence headers are discarded. The sequence is upper case and each 'U' is replaced by a 'T' before computation of the

digest.

prob-

similar,

proba-

collision). MD5

hexadecimal

conserve

The MD5 digest is a cryptographic hash function designed to minimize the ability that two different inputs give the same output, even for very but non-identical inputs. Still, there is a very small, but non-zero, probability that two different inputs give the same digest (i.e. a generates a 128-bit (16-byte) digest that is represented by 16 numbers (using 32 symbols among 0123456789abcdef). Use --sizeout to the abundance annotations.

--relabel\_self  
Rerelabel sequences using each sequence itself as a label.

--relabel\_sha1  
Rerelabel sequences using the SHA1 message digest algorithm applied to each sequence. It is similar to the --relabel\_md5 option but uses the SHA1 instead of the MD5 algorithm. SHA1 generates a 160-bit (20-byte) is represented by 20 hexadecimal numbers (40 symbols). The collision (two non-identical sequences resulting in the same smaller for the SHA1 algorithm than it is for the MD5 algorithm.

each se-

algorithm

digest that

probability of a

digest) is

matches

rate in

present

nucleotide se-

format

always im-

external ref-

by de-

Multi-

--self When using --uchime\_ref, ignore a reference sequence when its label the label of the query sequence (useful to estimate false-positive reference sequences).

--selfid When using --uchime\_ref, ignore a reference sequence when its sequence is strictly identical to the nucleotidic sequence of the query.

--sizein In de novo mode, abundance annotations (pattern '[>]size=integer[;]') in sequence headers are taken into account by default (--sizein is implied). This option is ignored by --uchime\_ref.

--sizeout When relabelling, add abundance annotations to fasta headers (using the ';size=integer;').

--uchime\_denovo filename  
Detect chimeras present in the fasta-formatted filename, without erences (i.e. de novo). Automatically sort the sequences in filename creasing abundance beforehand (see the sorting section for details). threading is not supported.

--uchime2\_denovo filename  
Detect chimeras present in the fasta-formatted filename, using the algorithm. This algorithm is designed for denoised amplicons (see --oile). Automatically sort the sequences in filename by decreasing abundance beforehand (see the sorting section for details). Multithreading is supported.

--uchime3\_denovo filename  
Detect chimeras present in the fasta-formatted filename, using the gorithm. The only difference from --uchime2\_denovo is that the default abundance skew (--abskew) is set to 16.0 rather than 2.0.

--uchime\_ref filename  
Detect chimeras present in the fasta-formatted filename by comparing reference sequences (option --db). Multithreading is supported.

--uchimealns filename  
Write the three-way global alignments (parentA, parentB, chimera) to using a human-readable format. Use --alignwidth to modify alignment Output order may vary when using multiple threads. All sequences are to upper case before alignment. Lower case letters indicate the alignment.

--uchimeout filename  
Write chimera detection results to filename using a 18-field, tab-uchime-like format. Use --uchimeout5 to use a format compatible with v5 and earlier versions. Rows output order may vary when using threads.

1. score: higher score means a more likely chimeric alignment.
2. Q: query sequence label.
3. A: parent A sequence label.
4. B: parent B sequence label.
5. T: top parent sequence label (i.e. parent most similar query). That field is removed when using --uchimeout5.
6. idQM: percentage of similarity of query (Q) and model constructed as a part of parent A and a part of parent B.

(T).

7. idQA: percentage of similarity of query (Q) and parent A.
8. idQB: percentage of similarity of query (Q) and parent B.
9. idAB: percentage of similarity of parent A and parent B.
10. idQT: percentage of similarity of query (Q) and top parent (T).
11. LY: yes votes in the left part of the model.
12. LN: no votes in the left part of the model.
13. LA: abstain votes in the left part of the model.
14. RY: yes votes in the right part of the model.
15. RN: no votes in the right part of the model.
16. RA: abstain votes in the right part of the model.
17. div: divergence, defined as (idQM - idQT).
18. YN: query is chimeric (Y), or not (N), or is a borderline case (?).

**--uchimeout5**field,  
compati-

When using --uchimeout, write chimera detection results using a 17-field, tab-separated uchime-like format (drop the 5th field of --uchimeout), compatible with usearch version 5 and earlier versions.

**--xlength**

attribute

func-  
tagging  
negatives).

Strip header attribute ";length=integer" from input sequences. This is added to output sequences by the --lengthout option.

**--xn strictly positive real number**similar  
Important  
defini-

weight of no votes, corresponding to the parameter beta in the scoring function (default value is 8.0). Increasing --xn reduces the likelihood of a sequence as a chimera (less false positives, but also more false negatives).

file.

--xsize Strip abundance information from the headers when writing the output file.

**Clustering options:**similar  
Important  
defini-

vsearch implements a single-pass, greedy centroid-based clustering algorithm, to the algorithms implemented in usearch, DNACLUST and sumaclust for example. parameters are the global clustering threshold (--id) and the pairwise identity definition (--iddef).

Input sequences are masked as specified with the --qmask and --hardmask options.

--biomout filename  
Generate an OTU table in the biom version 1.0 JSON file format as  
specified at  
biom-1.0.html)  
biom-1.0.html>. The  
of the  
than the  
otutabout and  
least for  
Taxonomy in-  
identifiers will  
If the  
similar  
will be  
header.  
semicolons.  
of the  
allowed. OTU  
centroid se-  
somewhere,  
semicolon is  
identifier may  
label is  
and all  
identifiers can  
--rela-  
also be  
contains  
taxonomy  
mandatory  
contain

(link) ([https://biom-format.org/documentation/format\\_versions/](https://biom-format.org/documentation/format_versions/)  
[https://biom-format.org/documentation/format\\_versions/](https://biom-format.org/documentation/format_versions/))

The format describes how to store a sparse matrix containing the abundances OTUs in the different samples. This format is much more efficient than the classic and mothur OTU table formats available with the -- --mothur\_shared\_out options, respectively, and is recommended at least for large tables. The OTUs are represented by the cluster centroids. Information will be included for the OTUs if available. Sample information will be included for the OTUs if available. Sample be extracted from the headers of all sequences in the input file. header contains ';sample=abc123;' or ';barcodelabel=abc123;' or a string somewhere, then the given sample identifier (here 'abc123') used. The semicolon is not mandatory at the beginning or end of the header. The sample identifier may contain any printable character except If no such sample label is found, the identifier in the initial part header will be used, but only letters, digits and underscores are identifiers will be extracted from the headers of the cluster sequences. If the header contains ';otu=def789;' or a similar string then the given OTU identifier (here 'def789') will be used. The semicolon is not mandatory at the beginning or end of the header. The OTU contain any printable character except semicolons. If no such OTU found, the identifier in the initial part of the header will be used, characters except semicolons are allowed. Alternatively, OTU be generated using the relabelling options (--relabel, --relabel\_self, --bel\_shal, or --relabel\_md5). Taxonomy information, if present, will be extracted from the headers of the centroid sequences. If the header ';tax=Homo\_sapiens;' or a similar string somewhere, then the given taxonomy information (here 'Homo\_sapiens') will be used. The semicolon is not at the beginning or end of the header. The taxonomy information may

any printable character except semicolons. If an OTU table in the biom version 2.1 HDF5 file format is required, the biom utility may be used as described at ([link](https://biom-format.org/documentation/biom_conversion.html)) (<https://biom-format.org/documentation/> <[https://biom-format.org/documentation/biom\\_conversion.html](https://biom-format.org/documentation/biom_conversion.html)>).

--centroids filename  
Output cluster centroid sequences to filename, in fasta format. The centroid is the sequence that seeded the cluster (i.e. the first sequence of the cluster).

--clusterout\_id  
Add cluster identifier information to the output files when using the --centroids, --consout and --profile options.

--clusterout\_sort  
Sort some output files by decreasing abundance instead of input order. It applies to the --consout, --msaout, --profile, --centroids, and --uc (the C lines). For --uc, the sorting applies only to the centroid information part lines).

--cluster\_fast filename  
Clusterize the fasta sequences in filename, automatically sort by sequence length beforehand.

--cluster\_size filename  
Clusterize the fasta sequences in filename, automatically sort by sequence abundance beforehand.

--cluster\_smallmem filename  
Clusterize the fasta sequences in filename without automatically modifying their order beforehand. Sequence are expected to be sorted by sequence length, unless --usersort is used.

--cluster\_unoise filename  
Perform denoising of the fasta sequences in filename according to the UNOISE removal options specified. In sequence distance and the abundance ratio. The abundance ratio (skew) is the abundance of a new sequence divided by the abundance of the centroid sequence.

This skew together. times the mismatches in exponen- sequence to will form

and a

center- center, thresh- symbol (nu- a ma- option

description of gap ex- and a sequences

different con- (L) or Sequence sym- numer- contexts:

must not be larger than beta if the sequences should be clustered. Beta is calculated as 2 raised to the power of minus 1 minus alpha sequence distance. The sequence distance used is the number of the alignment, ignoring gaps. This means that the abundance must be tially lower as the distance increases from the centroid for a new be included in the cluster. Nearer sequences with higher abundances their own new clusters.

**--clusters string**  
Output each cluster to a separate fasta file using the prefix string ticker (0, 1, 2, etc.) to construct the path and filenames.

**--consout filename**  
Output cluster consensus sequences to filename. For each cluster, a star multiple sequence alignment is computed with the centroid as the using a fast algorithm (not accurate when using low pairwise identity olds). A consensus sequence is constructed by taking the majority cleotide or gap) from each column of the alignment. Columns containing jority of gaps are skipped, except for terminal gaps. If the --sizein is specified, sequence abundances will be taken into account.

**--cons\_truncate**  
This command is ignored. A warning is issued.

**--gapext string**  
Set penalties for a gap extension. See --gapopen for a complete the penalty declaration system. The default is to initialize the six tending penalties using a penalty of 2 for extending internal gaps penalty of 1 for extending terminal gaps, in both query and target (i.e. 2I/1E).

**--gapopen string**  
Set penalties for a gap opening. A gap opening can occur in six texts: in the query (Q) or in the target (T) sequence, at the left right (R) extremity of the sequence, or inside the sequence (I). bols (Q and T) can be combined with location symbols (L, I, and R), and ical values to declare penalties for all possible

aQL/bQL/cQR/dL/eL/fL, where abcdef are zero or positive integers, and '/' is used as a separator.

To simplify declarations, the location symbols (L, I, and R) can be combined, equally, and sequences opening right), in value is applies to be de- null- default pa- he/she symbols integer gap usearch, all maintain

the symbol (E) can be used to treat both extremities (L and R) the symbols Q and T can be omitted to treat query and target equally. For instance, the default is to declare a penalty of 20 for internal gaps and a penalty of 2 for opening terminal gaps (left or both query and target sequences (i.e. 20I/2E). If only a numerical given, without any sequence or location symbol, then the penalty all gap openings. To forbid gap-opening, an infinite penalty value can clared with the symbol '\*'. To use vsearch as a semi-global aligner, a penalty can be applied to the left (L) or right (R) gaps. vsearch always initializes the six gap opening penalties using the rameters (20I/2E). The user is then free to declare only the values wants to modify. The string is scanned from left to right, accepted are (0123456789/LIREQT\*), and later values override previous values. Please note that vsearch, in contrast to usearch, only allows penalties. Because the lowest gap penalties are 0.5 by default in default scores and gap penalties in vsearch have been doubled to equivalent penalties and to produce identical alignments.

--id real  
Do not add the target to the cluster if the pairwise identity with  
troid is lower than real (value ranging from 0.0 to 1.0 included). The  
wise identity is defined as the number of (matching columns) /  
length - terminal gaps). That definition can be modified by --iddef.

--iddef 0|1|2|3|4  
Change the pairwise identity definition used in --id. Values accepted

are:

`length).`

- 0. CD-HIT definition: (matching columns) / (shortest sequence)
  - 1. edit distance: (matching columns) / (alignment length).
  - 2. edit distance excluding terminal gaps (same as --id).
  - 3. Marine Biological Lab definition counting each gap opening  
    nal or terminal) as a single mismatch, whether or not the

(inter-

gap was sequence extended: 1.0 - [(mismatches + gap openings)/(longest length)]

global 4. BLAST definition, equivalent to --iddef 1 in a context of pairwise alignment.

format by --lengthout Write sequence length information to the output files in FASTA adding a ";length=integer" attribute in the header.

align- --match integer Score assigned to a match (i.e. identical nucleotides) in the pairwise ment. The default value is 2.

cluster\_un- --minsize positive integer Specify the minimum abundance of sequences for denoising using -- noise. The default is 8.

pairwise --mismatch integer Score assigned to a mismatch (i.e. different nucleotides) in the alignment. The default value is -4.

cluster --msaout filename Output a multiple sequence alignment and a consensus sequence for each star mul- to filename, in fasta format. Be warned that vsearch computes center tiple sequence alignments using a fast method whose accuracy can decrease sig- nificantly when using low pairwise identity thresholds. The consensus sequence is constructed by taking the majority symbol (nucleotide or gap) column of the alignment. Columns containing a majority of gaps are skipped, except for terminal gaps. If the --sizein option is specified, abundances will be taken into account when computing the consensus.

abun- format as Shared\_file) --mothur\_shared\_out filename Output an OTU table in the mothur 'shared' tab-separated plain text described at (link) ([https://www.mothur.org/wiki/Shared\\_file](https://www.mothur.org/wiki/Shared_file)). The format describes how a matrix stored. The and is for each containing the abundances of the OTUs in the different samples is first line will start with the strings 'label', 'group' and 'numOTus' followed by a list of all OTU identifiers. The following lines, one sample, starts with the string 'vsearch' followed by the sample

identifier,  
that sam-  
identifiers are  
represented by  
matrix  
first line  
separated list  
starts with  
abundances for  
OTU and  
sequences (see  
centroids. An  
information is  
'taxon-  
extracted for  
nucleotide  
FASTA-  
in a  
consensus nu-  
number  
symbols  
sizein  
not mask  
sensitive.

the total number of OTUs, and a list of abundances for each OTU in  
ple, in the order given on the first line. The OTU and sample  
extracted from the FASTA headers of the sequences. The OTUs are  
the cluster centroids. See the --biomout option for further details.

--otutabout filename  
Output an OTU table in the classic tab-separated plain text format as a  
containing the abundances of the OTUs in the different samples. The  
will start with the string '#OTU ID' and is followed by a tab-  
of all sample identifiers. The following lines, one for each OTU,  
the OTU identifier and is followed by a tab-separated list of  
that OTU in each sample, in the order given on the first line. The  
sample identifiers are extracted from the FASTA headers of the  
the --sample option). The OTUs are represented by the cluster  
extra column is added to the right of the table if taxonomy  
available for at least one of the OTUs. This column will be labelled  
omy' and each row will then contain the taxonomy information  
that OTU. See the --biomout option for further details.

--profile filename  
Output a sequence profile to a text file with the frequency of each  
in each position in the multiple alignment for each cluster. There is a  
like header line for each cluster, followed by the profile information  
tab-separated format. The eight columns are: position (0-based),  
cleotide, number of As, number of Cs, number of Gs, number of Ts or Us,  
of gap symbols, and finally the total number of ambiguous nucleotide  
(B, D, H, K, M, N, R, S, Y, V or W). All numbers are integers. If the --  
option is specified, sequence abundances will be taken into account.

--qmask none|dust|soft  
Mask regions in sequences using the dust or the soft methods, or do  
(none). Warning, when using soft masking, clustering becomes case  
The default is to mask using dust.

--qsegout filename  
Write the aligned part of each query sequence to filename in FASTA

format.

```
--relabel string
    Relabel sequence identifiers in the output files produced by --consout,
--pro-
    file and --centroids options. Please see the description of the same
option
    under Chimera detection for details.

--relabel_keep
    When relabelling, keep the old identifier in the header after a space.

--relabel_md5
    Relabel sequence identifiers in the output files produced by --consout,
--pro-
    file and --centroids options. Please see the description of the same
option
    under Chimera detection for details.

--relabel_self
    Relabel sequence identifiers in the output files produced by --consout,
--pro-
    file and --centroids options. Please see the description of the same
option
    under Chimera detection for details.

--relabel_shal
    Relabel sequence identifiers in the output files produced by --consout,
--pro-
    file and --centroids options. Please see the description of the same
option
    under Chimera detection for details.

fasta file
    --sizein Take into account the abundance annotations present in the input
        (search for the pattern '[>]size=integer[;]' in sequence headers).

distance
    --sizeorder
        When an amplicon is close to 2 or more centroids, both within the
it with
        specified with the --id option, resolve the ambiguity by clustering
closest one.
        the centroid having the highest abundance, not necessarily the
maxaccepts is
        The option only has effect when the value specified with --
referred to
        higher than one. The --sizeorder option turns on what is sometimes
default dis-
        as abundance-based greedy clustering (AGC), in contrast to the
tance-based greedy clustering (DGC).

pattern
    --sizeout
        Add abundance annotations to the output fasta files (add the
abundance an-
        ';'size=integer;' to sequence headers). If --sizein is specified,
receives a
        notations are reported to output files, and each cluster centroid
            new abundance value corresponding to the total abundance of the
```

amplicons in-  
specified, in-  
amplicons per

strand only

format

fasta se-  
(H) as-  
centroid  
uc\_allhits  
entry

set to

pairwise

Gapped

(inser-  
identical to  
sequence (S,

cluded in the cluster (--centroids option). If --sizein is not put abundances are set to 1 for amplicons, and to the number of cluster for centroids.

--strand plus|both  
When comparing sequences with the cluster seed, check the plus (default) or check both strands.

--tsegout filename  
Write the aligned part of each target sequence to filename in FASTA format.

--uc filename  
Output clustering results in filename using a tab-separated uclust-like with 10 columns and 3 different type of entries (S, H or C). Each quence in the input file can be either a cluster centroid (S) or a hit signed to a cluster. Cluster records (C) summarize information (size, label) for each cluster. In the context of clustering, the option -- has no effect on the --uc output. Column content varies with the type of (S, H or C):

1. Record type: S, H, or C.
2. Cluster number (zero-based).
3. Centroid length (S), query length (H), or cluster size (C).
4. Percentage of similarity with the centroid sequence (H), or '\*' (S, C).
5. Match orientation + or - (H), or set to '\*' (S, C).
6. Not used, always set to '\*' (S, C) or to zero (H).
7. Not used, always set to '\*' (S, C) or to zero (H).
8. set to '\*' (S, C) or, for H, compact representation of the alignment using the CIGAR format (Compact Idiosyncratic Alignment Report): M (match/mismatch), D (deletion) and I tion). The equal sign '=' indicates that the query is the centroid sequence (ignoring terminal gaps).
9. Label of the query sequence (H), or of the centroid C).

10. Label of the centroid sequence (H), or set to '\*' (S, C).

--unoise\_alpha real

Specify the alpha parameter to the --cluster\_anoise command. The default is 2.0.

--usersort

When using --cluster\_smallmem, allow any sequence input order, not just decreasing length ordering.

--xlength

Strip header attribute ";length=integer" from input sequences. This is added to output sequences by the --lengthout option.

--xsize

Strip abundance information from the headers when writing the output file.

... masking

Most searching options as well as score filtering, gap penalties and also apply to clustering (see the Searching section for definitions): --alnout, --blast6out, --fastapairs, --matched, --notmatched, --maxaccepts, --maxrejects, --samout, --userout, --userfields

#### Dereplication and rereplication options:

VSEARCH can derePLICATE sequences with the commands --derep\_fulllength, --

derep\_id, command is --derep\_smallmem, --derep\_prefix and --fastx\_uniques. The --derep\_fulllength

handle depreciated and is replaced by the new --fastx\_uniques command that can also -- derep\_smallmem, and

length, but --fastx\_uniques commands requires strictly identical sequences of the same

derep\_id command ignores upper/lower case and treats T and U as identical symbols. The --

derep\_prefix command requires both identical sequences and identical headers/labels. The --

to be command will group sequences with a common prefix and does not require them

derepli- equally long. The --derep\_smallmem uses a much smaller amount of memory when

from a cating than the other files, and may be a bit slower and cannot read the input

to the pipe. It takes both FASTA and FASTQ files as input but only writes FASTA output

FASTQ file specified with the --fastaout option. The --fastx\_uniques command can write

as well output (specified with --fastqout) or FASTA output (specified with --fastaout)

commands as a special tab-separated column text format (with --tabbedout). The other

can write FASTA output to the file specified with the --output option. All

dereplica-  
like file  
sequences in  
options  
fastq\_qmin,  
--re-  
--size-  
sequences  
(case  
sizein and  
command does  
with the  
on the  
with the  
written  
written in  
descending  
but not  
proper  
bit hash  
however  
of one  
appr. 24  
options

tion commands, except --deref\_smallmem, can write output to a special UCLUST-specified with the --uc option. The --rereplicate command can duplicate the input file according to the abundance of each input sequence. Other valid options are --fastq\_ascii, --fastq\_asciiout, --fastq\_qmax, --fastq\_qmaxout, --fastq\_qminout, --fastq\_qout\_max, --lengthout, --maxuniquesize, --minuniquesize, --relabel, --relabel\_keep, --relabel\_md5, --relabel\_self, --relabel\_shal, --sizein, --sizeout, --strand, --topn, --xlength, and --xsize.

--deref\_fulllength filename  
Merge strictly identical sequences contained in filename. Identical are defined as having the same length and the same string of nucleotides insensitive, T and U are considered the same). See the options --sizeout to take into account and compute abundance values. This not support multithreading.

--deref\_id filename  
Merge strictly identical sequences contained in filename, as --deref\_fulllength command, but the sequence labels (identifiers) header line need to be identical too.

--deref\_smallmem filename  
Merge strictly identical sequences contained in filename, as --deref\_fulllength command, but using much less memory. The output is to a FASTA file specified with the --fastaout option. The output is the order that the sequences first appear in the input, and not in abundance order as with the other dereplication commands. It can read, write FASTQ files. This command cannot read from a pipe, it must be a file, as it is read twice. Dereplication is performed with a 128 function and it is not verified that grouped sequences are identical, the probability that two different sequences are grouped in a dataset billion unique sequences is approximately 1e-21. Memory footprint is bytes times the number of unique sequence. Multithreading and the --topn, --uc, or --tabbedout are not supported.

--deref\_prefix filename  
Merge sequences with identical prefixes contained in filename. A

short sequence identical to an initial segment (prefix) of another sequence is considered a replicate of the longer sequence. If a sequence is identical prefix of two or more longer sequences, it is clustered with the them. If they are equally long, it is clustered with the most maining ties are solved using sequence headers and sequence input sequence comparisons are case insensitive, and T and U are considered identical.

This command does not support multithreading.

sorted by first sequence (i.e. header used for valid for

--fastaout filename Write the dereplicated sequences to filename, in fasta format and decreasing abundance. Identical sequences receive the header of the sequence of their group. If --sizeout is used, the number of occurrences (abundance) of each sequence is indicated at the end of their fasta using the pattern ';size=integer;'. This option is only --fastx\_uniques and --deref\_smallmem.

sorted by first sequence (i.e. header used for valid for

--fastqout filename Write the dereplicated sequences to filename, in fastq format and decreasing abundance. Identical sequences receive the header of the sequence of their group. If --sizeout is used, the number of occurrences (abundance) of each sequence is indicated at the end of their fastq using the pattern ';size=integer;'. This option is only --fastx\_uniques.

quality FASTQ Illumina

--fastq\_ascii positive integer Define the ASCII character number used as the basis for the FASTQ score. The default is 33, which is used by the Sanger / Illumina 1.8+ format (phred+33). The value 64 is used by the Solexa, Illumina 1.3+ and mina 1.5+ formats (phred+64). Only 33 and 64 are valid arguments.

ASCII writing arguments.

--fastq\_asciiout positive integer When using --fastq\_convert, --sff\_convert or --fasta2fastq, define the character number used as the basis for the FASTQ quality score when FASTQ output files. The default is 33. Only 33 and 64 are valid

--fastq\_qmax positive integer Specify the maximum quality score accepted when reading FASTQ files.

The default formats may be used to specify the quality score range for FASTQ files. The default is 0, which is usual for recent Sanger/Illumina 1.8+ files. Older formats may use scores between -5 and 2.

The --fastq\_qminout option specifies the minimum quality score used when writing FASTQ files. The default is 0, which is usual for Sanger/Illumina 1.8+ files. Older versions of the format may use scores between -5 and 2.

The --fastq\_qout\_max option indicates that the new quality scores for dereplicating FASTQ files should be equal to the maximum (best) of the quality scores for each position (corresponding to the lowest error probability). The default is to output a quality score corresponding to the average of the error probabilities for each position.

The --fastx\_uniques option merges strictly identical sequences contained in FASTA or FASTQ file. Identical sequences are defined as having the same length and the same sequence of nucleotides (case insensitive, T and U are considered the same). Options --sizein and --sizeout take into account and compute values. This command does not support multithreading. By default, the scores in FASTQ output files will correspond to the average error of the nucleotides in the each position. If the --fastq\_qout\_max is given, the quality score will be the highest (best) quality score in each position.

The --lengthout option writes sequence length information to the output files in FASTA and FASTQ format by adding a ";length=integer" attribute in the header.

The --maxuniquesize option specifies the maximum size of unique sequences to be merged.

integers.

--minuniquesize positive integer  
Discard sequences with a post-dereplication abundance value smaller than integer.

sorted by first sequence (i.e. header used) is allowed for dereplication.

--output filename  
Write the dereplicated sequences to filename, in fasta format and decreasing abundance. Identical sequences receive the header of the sequence of their group. If --sizeout is used, the number of occurrences (abundance) of each sequence is indicated at the end of their fasta using the pattern ';size=integer;'. This option is not --fastx\_uniques or --deref\_smallmem.

for details.

--relabel string  
Please see the description of the same option under Chimera detection for details.

--relabel\_keep  
When relabelling, keep the old identifier in the header after a space.

for details.

--relabel\_md5  
Please see the description of the same option under Chimera detection for details.

for details.

--relabel\_self  
Please see the description of the same option under Chimera detection for details.

for details.

--relabel\_shal  
Please see the description of the same option under Chimera detection for details.

of each sequence. The sequence labels are identical for the same sequence, unless --relabel, --relabel\_self, --relabel\_shal or --relabel\_md5 is used to create unique output labels.

FASTA --size-out is specified, in which case an abundance of 1 is used.

fasta file --sizein Take into account the abundance annotations present in the input  
That op- (search for the pattern '[>]size=integer[;]' in sequence headers).  
tion is active by default when rereplicating.

' ;sizein- --sizeout Add abundance annotations to the output fasta file (add the pattern  
sequence teger;' to sequence headers). If --sizein is specified, each unique  
(sum of receives a new abundance value corresponding to its total abundance  
abun- the abundances of its occurrences). If --sizein is not specified, input  
value dences are set to 1, and each unique sequence receives a new abundance  
corresponding to its number of occurrences in the input file.

strand only --strand plus|both When searching for strictly identical sequences, check the plus  
(default) or check both strands.

columns --tabbedout filename Output clustering info to the specified tab-separated text file with 6  
header and a row for each input sequence. Column 1 contains the original label/  
which is of the sequence. Column 2 contains the label of the output sequence  
poten- equal to the label/header of the first sequence in each cluster, but  
tially relabelled. Column 3 contains the cluster number, starting from  
0. Col- umn 4 contains the sequence number within each cluster, starting at 0.  
Column 5 contains the number of sequences in the cluster. Column 6 contains the  
orig- inal label/header of the first sequence in the cluster before any  
potential relabelling. This option is only valid for the --fastx\_uniques command.

a tab- --topn positive integer Output only the top integer sequences (i.e. the most abundant).

entries --uc filename Output full-length or prefix-dereplication results in filename using  
cluster separated uclust-like format with 10 columns and 3 different type of  
summarize (S, H or C). Each fasta sequence in the input file can be either a  
context of centroid (S) or a hit (H) assigned to a cluster. Cluster records (C)  
output. Col- information (size, centroid label) for each cluster. In the  
dereplication, the option --uc\_allhits has no effect on the --uc  
umn content varies with the type of entry (S, H or C):

1. Record type: S, H, or C.
2. Cluster number (zero-based).
3. Sequence length (S, H), or cluster size (C).
4. Percentage of similarity with the centroid sequence (H), or '\*' (S, C).
5. Match orientation + or - (H), or set to '\*' (S, C).
6. Not used, always set to '\*' (S, C) or 0 (H).
7. Not used, always set to '\*' (S, C) or 0 (H).
8. Not used, always set to '\*'.
9. Label of the query sequence (H), or of the centroid C).
10. Label of the centroid sequence (H), or set to '\*' (S, C).

--xlength  
    Strip header attribute ";length=integer" from input sequences. This  
        added to output sequences by the --lengthout option.  
  
--xsize  
    Strip abundance information from the headers when writing the output file.

### Extraction options:

Sequences with headers matching certain criteria can be extracted from FASTA and files using the `--fastx getseq`, `--fastx getsegs` and `--fastx getsubseq` commands.

The `--fastx_getseq` command requires the header to match a label specified `--label` option. If the `--label_substr_match` option is given, the label may be string located anywhere in the header, otherwise the entire header must match bel. These matches are not case-sensitive. The headers in the input file are at the first space or tab character unless the `--notrunclabels` option is matching sequences will be written to the files specified with the `--fastqout` options, in FASTA and FASTQ format, respectively. Sequences that match are written to the files specified with the `--notmatched` and `--op-` tions, respectively.

The `--fastx_getsubseq` command is similar to the `--fastx_getseq` command, but extract a subsequence of the matching sequences. The start position is specified

with the --subseq\_start option and the end position is specified with the --subseq\_end option. The positions are 1-based, meaning that the first symbol of the sequence is at 1. If the start or end position option is not specified, the default is to start first position and end at the last position in the sequence.

The --fastx\_getseqs command is similar to the --fastx\_getseq command but flexibility in specifying the label(s) to be matched. A single label may be using the --label option as described above. Alternatively, a file containing a labels to be matched may be specified with the --labels option. The file plain text file with one label on each line. The --label\_word and --label\_words may be used to specify either a single word or a file containing a list of respectively, to be matched. Words are defined as character sequences delimited a character that is not alpha-numeric (A-Z, a-z, or 0-9) or by the beginning or the header. Word matching is case-sensitive. The --label\_field option will matching of words to a certain field in the header.

--fastaout filename

Write the extracted sequences in FASTA format to the file with the given name.

--fastqout filename

Write the extracted sequences in FASTQ format to the file with the given name.

This option is illegal if the input is in FASTA format.

--fastx\_getseq filename

Extract sequences from the given FASTA or FASTQ file. Specify a label to match using the --label option. Output files are specified with the --fastaout, --fastqout, --notmatched and --notmatchedfq options.

--fastx\_getseqs filename

Extract sequences from the given FASTA or FASTQ file. Specify the labels to match using one of the following options: --label, --labels, --label\_word, or --label\_words. Output files are specified with the --fastaout, --fastqout, --notmatched and --notmatchedfq options.

--fastx\_getsubseq filename

Extract a certain part of some of the sequences in the given FASTA or file. Specify labels to match using the --label option. Specify the sequence range to be extracted with the --subseq\_start and --subseq\_end

options.

notmatched and

Output files are specified with the --fastaout, --fastqout, --  
--notmatchedfq options.

--label string

Specify the label to match in the sequence header. Unless the --  
str\_match option is given, the label must match the entire header. The  
comparison is not case-sensitive.

--label\_field string

Specify a field name to be used when matching using the --label\_word or  
bel\_words option. The field name is a string like "abc" that must  
word to be matched with an equals sign (=) in between. The field must  
limited by semicolons or the beginning or end of the header. The  
header will match the label 123 in the field abc: "seq1;abc=123".

--label\_substr\_match

The labels specified with the --label or the --labels option may  
where in the header if this option is given. Otherwise a label needs to  
the entire header.

--label\_word string

Specify a word to match in the sequence header. Words are defined as  
delimited by either the start or end of the header or by any symbol  
not a letter (A-Z, a-z) or digit (0-9). The comparison is case-  
sensitive.

--label\_words filename

Specify a file containing words to be matched against the sequence  
headers.

defined as

symbol

sensi-

The plain text file must contain one word on each line. Words are  
strings delimited by either the start or end of the header or by any  
that is not a letter (A-Z, a-z) or digit (0-9). The comparison is case-  
tive.

--labels filename

Specify a file containing labels to be matched against the sequence  
headers.

--la-

header. The

The plain text file must contain one label on each line. Unless the  
bel\_substr\_match option is given, a label must match the entire  
comparison is not case-sensitive.

--notmatched filename

Write the sequences that were not extracted to the file with the  
given

name,  
                  in FASTA format.

name,  
                  --notmatchedfq filename  
                  Write the sequences that were not extracted to the file with the given  
                  name,  
                  in FASTQ format. This option is illegal if the input is in FASTA format.

using  
start  
option  
                  --subseq\_end positive integer  
                  Specify the end position in the sequences when extracting subsequences  
                  the --fastx\_getsubseq command. Positions are 1-based, so the sequences  
                  at position 1. The default is to end at the end of the sequence if this  
                  is not specified.

subsequences  
sequences  
sequence  
                  --subseq\_start positive integer  
                  Specify the starting position in the sequences when extracting  
                  using the --fastx\_getsubseq command. Positions are 1-based, so the  
                  start at position 1. The default is to start at the beginning of the  
                  (position 1), if this option is not specified.

FASTA/FASTQ/SFF file processing options:

FASTA,  
files to  
convert  
Statistical  
performed  
Sequences may  
commands.  
while the  
quality  
joined  
complements  
include the  
FASTA out-  
--xee  
                  Analyse, trim, filter, convert, merge, join or reverse complement sequences in  
                  FASTQ or SFF files. The --fastq\_chars command can be used to analyse FASTQ  
                  identify the quality encoding and the range of quality score values used. To  
                  between different FASTQ file variants, use the --fastq\_convert command.  
                  analysis of the quality and length of the sequences in a FASTQ file may be  
                  with the --fastq\_stats, --fastq\_eestats, and --fastq\_eestats2 commands.  
                  be trimmed, filtered and converted by the --fastq\_filter or --fastx\_filter  
                  The --sff\_convert command can be used to convert SFF files to FASTQ,  
                  --fasta2fastq command will convert a FASTA file to a FASTQ file with fake  
                  scores. Paired-end reads can be merged using the --fastq\_mergepairs command or  
                  with the --fastq\_join command. The --fastx\_revcomp command will reverse-  
                  sequences.

                  --eeout When using --fastq\_filter, --fastx\_filter or --fastq\_mergepairs,  
                  number of expected errors (ee) in the sequence header of FASTQ and  
                  put files. This option is a synonym of the --fastq\_eeout option. Use the  
                  option to remove this information from headers.

--eetabbedout filename  
When specified with the --fastq\_mergepairs command, write statistics expected errors of each merged read to the given file. The file is a tab separated file with four columns: The number of expected errors in the read, the number of expected errors in the reverse read, the number served errors in the forward read, and the number of observed errors reverse read. The observed number of errors are the number of the overlap region of the merged sequence relative to each of the reads pair.

--fasta2fastq filename  
Add a fake nucleotide quality score to the sequences in the given FASTA file and write them to the FASTQ file specified with the --fastqout option. The quality score may be adjusted using the --fastq\_qmaxout option (default 41). The --fastq\_asciiout option may be used to adjust the FASTQ output quality ASCII base character (default 33).

--fastaout filename  
When using --fastq\_filter, --fastq\_mergepairs or --fastx\_filter, write given FASTA-formatted file the sequences passing the filter, or the sequences.

--fastaout\_rev filename  
When using --fastq\_filter, or --fastx\_filter, write to the given FASTA-formatted file the reverse reads passing the filter.

--fastaout\_notmerged\_fwd filename  
When using --fastq\_mergepairs, write forward reads not merged to the specified FASTA file.

--fastaout\_notmerged\_rev filename  
When using --fastq\_mergepairs, write reverse reads not merged to the specified FASTA file.

--fastaout\_discarded filename  
Write sequences that do not pass the filter of the --fastq\_filter or --fastx\_filter command to the given FASTA-formatted file.

--fastaout\_discarded\_rev filename  
Write reverse reads that do not pass the filter of the --fastq\_filter or --fastx\_filter command to the given FASTA-formatted file.

--fastq\_allowmergestagger  
When using --fastq\_mergepairs, allow merging of staggered read pairs.  
Staggered pairs are pairs where the 3' end of the reverse read has an overhang to the left of the 5' end of the forward read. This situation can occur when a very short fragment is sequenced. The 3' overhang of the reverse read is included in the merged sequence. The opposite option is the --fastq\_nostagger option. The default is to discard staggered pairs.

--fastq\_ascii positive integer  
Define the ASCII character number used as the basis for the FASTQ quality score. The default is 33, which is used by the Sanger / Illumina 1.8+ FASTQ format (phred+33). The value 64 is used by the Solexa, Illumina 1.3+ and Illumina 1.5+ formats (phred+64). Only 33 and 64 are valid arguments.

--fastq\_asciiout positive integer  
When using --fastq\_convert, --sff\_convert or --fasta2fastq, define the character number used as the basis for the FASTQ quality score when FASTQ output files. The default is 33. Only 33 and 64 are valid arguments.

--fastq\_chars filename  
Summarize the composition of sequence and quality strings contained in the input FASTQ file. For each sequence symbol, --fastq\_chars gives the occurrences of the symbol, its relative frequency and the length of the longest run of that symbol. For each character present in the quality string, --fastq\_chars gives the ASCII value of the character, its relative frequency, and the number of times a k-mer of that character appears at the end of quality strings. The length of the k-mer can be set using --fastq\_tail (4 by default). The command --fastq\_chars tries to automatically detect the encoding (Solexa, Illumina 1.3+, Illumina 1.5+ or Illumina 1.8+/Sanger) by analyzing the range of observed quality score values. In case of --fastq\_chars suggests values for the --fastq\_ascii (33 or 64), --fastq\_qmin and --fastq\_qmax options to be used with the other commands that require a FASTQ input file.

--fastq\_convert filename  
Convert between the different variants of the FASTQ file format. The quality

encoding of the input file must be specified with the --fastq\_ascii option (either 33 or 64, the default is 33), and the output quality encoding specified with the --fastq\_asciiout option (default 33). The minimum and maximum output quality scores may be limited using the --fastq\_qminout and --fastq\_qmaxout options. The output file is specified with the --fastqout option.

--fastq\_eeout  
When using --fastq\_filter, --fastx\_filter or --fastq\_mergepairs, include the number of expected errors (ee) in the sequence header of FASTQ and FASTA files. This option is a synonym of the --eeout option. Use the --xee option to remove this information from headers.

--fastq\_eestats filename  
Analyze a FASTQ file and report statistics on the distributions of quality scores, error probabilities and expected accumulated errors. The report, a table of 21 tab-separated columns, is written to the file specified with the --output option. The first column corresponds to the position in the reads (Reads). The second and third columns correspond to the number of reads and percentage of reads (PctRecs) that include this position. The remaining columns include information about the distribution of quality scores in this position (Q), error probabilities in this position (Pe), and finally the expected number of accumulated errors from the beginning of the reads until the current position (EE). For each of the Q, Pe and EE following distributions, the following statistics are included: minimum value (Min), lower quartile (Low), median (Med), mean (Mean), upper quartile (Hi), and maximum value (Max). The quality encoding and the range of quality values may be specified with --fastq\_ascii --fastq\_qmin and --fastq\_qmax.

--fastq\_eestats2 filename  
Analyze the specified FASTQ file and report statistics on the number of sequences that would be retained at a combination of selected cutoffs for truncation and maximum expected errors, that could potentially be used as arguments to the --fastq\_trunclen and --fastq\_maxee options to the --fastq\_filter command. The result, a table of two or more columns, is written

to the  
trunc-  
truncation  
and, in  
the se-  
with the  
integers  
between  
indicates  
default  
so on  
error  
requires a  
default  
0.5, 1.0

to the  
fastx\_filter for  
between them  
with the  
reads are  
specified  
forward  
read. The  
padding  
padding se-  
IIIIIIII,  
score with

file specified with the --output option. There is a line for each length truncation cutoff. The first column on each line contains the selected length, while the following columns contain the number of sequences parenthesis, the percentage of sequences that would be retained at selected EE levels. The truncation length cutoffs may be specified --length\_cutoffs option and requires a list of three comma-separated indicating the shortest cutoff, the longest cutoff, and the increment cutoffs. The longest cutoff may be specified with a star (\*) which that the limit is equal to the longest sequence in the input file. The setting is "50,\*,50" meaning that truncation lengths of 50, 100, 150 and up to the longest sequence length should be used. The maximum expected (EE) cutoffs may be specified with the --ee\_cutoffs option which comma-separated list of floating point numbers as its argument. The setting is "0.5,1.0,2.0" that indicates that expected error levels of and 2.0 should be used.

--fastq\_filter filename  
Trim and/or filter sequences in the given FASTQ file. Similar  
--fastx\_filter command, but works only on FASTQ files. See --  
details.

--fastq\_join filename  
Join paired-end sequence reads into one sequence and add a gap  
using a padding sequence. The sequences are not merged as  
fastq\_mergepairs command, but simply joined with a gap. The forward  
specified as the argument to this option and the reverse reads are  
with the --reverse option. The resulting sequences consist of the  
read, the padding sequence and the reverse complement of the reverse  
padding sequence is specified with the --join\_padgap option and the  
quality is specified with the --join\_padgapq option. The default  
quence string is NNNNNNNN and the default padding quality string is  
corresponding to a base quality score of 40 (a very high quality  
error probability 0.0001). The joined sequences are output to the

file(s) specified with the --fastaout or --fastqout options.

--fastq\_maxdiffs positive integer  
When using --fastq\_mergepairs, specify the maximum number of non-matching nucleotides allowed in the overlap region. That option has a strong influence on the merging success rate. The default value is 10.

--fastq\_maxdiffpct real  
When using --fastq\_mergepairs, specify the maximum percentage of non-nucleotides allowed in the overlap region. The default value is 100.0%. There are other more sophisticated rules in the merging algorithm that will discard read pairs with a high fraction of mismatches.

--fastq\_maxee real  
When using --fastq\_filter, --fastq\_mergepairs or --fastx\_filter, discard sequences with an expected error greater than the specified number (value ranging from 0.0 to infinity). For a given sequence, the expected error is the sum of error probabilities for all the positions in the sequence. Since probabilities can be small but not null, the expected error is always greater than zero, and at most equal to the length of the sequence when all positions in the sequence have an error probability of 1.0.

distribution, instance,  
average to 1.0 of Using the expected error as the lambda parameter in the Poisson it is possible to compute the probability of observing k errors. For a read with an expected error of 1.0 has:

- 36.8% chance of having zero error,
- 36.8% chance of having one error,
- 18.4% chance of having two errors,
- 6.1% chance of having three errors,
- 1.5% chance of having four errors,
- 0.3% chance of having five errors,
- etc.

--fastq\_maxee\_rate real  
When using --fastq\_filter or --fastx\_filter, discard sequences with an average expected error greater than the specified number (value ranging from 0.0 to 1.0 included). For a given sequence, the average expected error is the sum of

error probabilities for all the positions in the sequence, divided by the length of the sequence.

length of

discard sequences with more than the specified number of bases.

sequence

discard sequences with more than the specified number of N's.

speci-  
with the  
reverse  
Thus,  
in both  
different  
specified with  
to the  
fastaout\_notmerged\_rev, may be trun- bases trun- bases by de- fastq\_allowmergestagger op- reads 10). The with the

--fastq\_maxlen positive integer When using --fastq\_filter, --fastq\_mergepairs or --fastx\_filter, sequences with more than the specified number of bases.

--fastq\_maxmrgelen positive integer When using --fastq\_mergepairs, specify the maximum length of the merged (default is 1,000,000).

--fastq\_maxns positive integer When using --fastq\_filter, --fastq\_mergepairs or --fastx\_filter, sequences with more than the specified number of N's.

--fastq\_mergepairs filename Merge paired-end sequence reads into one sequence. The forward reads are specified as the argument to this option and the reverse reads are specified --reverse option. Reads with the same index/position in the forward and files are considered to form a pair, even if their labels are different. forward and reverse reads must appear in the same order and total number files. A warning is emitted if the forward and reverse files contain numbers of reads. The merged sequences are written to the file(s) the --fastaout or --fastqout options. The non-merged reads can be output files specified with the --fastaout\_notmerged\_fwd, --fastqout\_notmerged\_fwd and --fastqout\_notmerged\_rev options. Statistics output to the file specified with the --eetabbedout option. Sequences are truncated as specified with the --fastq\_truncqual option to remove low-quality in the 3' end. Sequences shorter than specified with --fastq\_minlen (after truncation) are discarded (1 by default). Sequences with too many ambiguous (N's), as specified with the --fastq\_maxns are also discarded (no limit fault). Staggered reads are not merged unless the --fastq\_allowmergestagger option is specified. The minimum length of the overlap region between the may be specified with the --fastq\_minovlen option (at least 5, default overlap region may not include more mismatches than specified --fastq\_maxdiffs option (10 by default) or a higher percentage of

mismatches  
otherwise  
reads that  
length of  
fastq\_minmergelen and  
output  
options, but  
options are:  
fastq\_qmin, and

than specified with the --fastq\_maxdiffpct option (100.0% by default),  
the read pair is discarded. Additional rules will avoid merging of  
cannot be aligned reliably and unambiguously. The minimum and maximum  
the merged sequence may be specified with the --  
--fastq\_maxmergelen options, respectively. The quality value limits for  
files may be specified with the --fastq\_qminout and --fastq\_qmaxout  
they apply only to the merged region. Other relevant  
--fastq\_ascii, --fastq\_maxee, --fastq\_nostagger, --fastq\_qmax, --  
--label\_suffix.

input

--fastq\_minlen positive integer  
When using --fastq\_filter, --fastq\_mergepairs or --fastx\_filter, discard  
sequences with less than the specified number of bases (default 1).

merged se-

--fastq\_minmergelen positive integer  
When using --fastq\_mergepairs, specify the minimum length of the  
quence. The default is 1.

merged

--fastq\_minovlen positive integer  
When using --fastq\_mergepairs, specify the minimum overlap between the  
reads. The default is 10. Must be at least 5.

base with  
none.

--fastq\_minqual positive integer  
When using --fastq\_filter or --fastx\_filter, discard reads having any  
a quality score below the given value. The default is 0, which discards

pairs. This  
behaviour, see

--fastq\_nostagger  
When using --fastq\_mergepairs, forbid the merging of staggered read  
is the default behaviour of --fastq\_mergepairs. To change that  
the --fastq\_allowmergestagger option.

default

--fastq\_qmax positive integer  
Specify the maximum quality score accepted when reading FASTQ files. The  
is 41, which is usual for recent Sanger/Illumina 1.8+ files.

fasta2fastq,  
For the  
score used

--fastq\_qmaxout positive integer  
When using --fastq\_mergepairs, --fastq\_convert, --sff\_convert or --  
specify the maximum quality score used when writing FASTQ files.  
--fasta2fastq command, the value specified here is the fake quality  
for the FASTQ output file. The default is 41, which is usual for

recent Sanger/Illumina 1.8+ files. Older formats may use a maximum quality score of 40.  
of 40. The limit only applies to the merged region when using --fastq\_mergepairs.

is 0, --fastq\_qmin positive integer  
may use Specify the minimum quality score accepted for FASTQ files. The default which is usual for recent Sanger/Illumina 1.8+ files. Older formats scores between -5 and 2.

specify the --fastq\_qminout positive integer  
which is When using --fastq\_mergepairs, --fastq\_convert or --sff\_convert, minimum quality score used when writing FASTQ files. The default is 0, usual for Sanger/Illumina 1.8+ files. Older versions of the format scores between -5 and 2. The limit applies only to the merged region when using --fastq\_mergepairs.

quality en- --fastq\_stats filename  
fastq\_ascii Analyze a FASTQ file and report the number of reads it contains. The and out- coding and the range of quality values may be specified with -- length vs. --fastq\_qmin and --fastq\_qmax. That command requires the --log option puts the following detailed statistics on read length, quality score, quality distributions, and length / quality filtering:

Read length distribution:

1. L: read length.
2. N: number of reads.
3. Pct: fraction of reads with this length.
4. AccPct: fraction of reads with this length or longer.

Quality score distribution:

1. ASCII: character encoding the quality score.
2. Q: Phred quality score.
3. Pe: probability of error associated with the quality score.
4. N: number of bases with this quality score.
5. Pct: fraction of bases with this quality score.
6. AccPct: fraction of bases with this quality score or higher.

Length vs. quality distribution:

position.

1.

indi-  
fastq\_filter  
fastq\_trunclen  
or 0.1  
indicate the  
command

columns in-  
fastq\_filter  
fastq\_trun-  
lesser

bases  
null,

bases  
null,

1. L: position in reads (starting from position 2).
  2. PctRecs: fraction of reads with at least this length.
  3. AvgQ: average quality score at this position.
  4. P(AvgQ): error probability corresponding to AvgQ.
  5. AvgP: average error probability at this position.
  - 6: AvgEE: average expected error over all reads up to this position.
  - 7: Rate: growth rate of AvgEE between this position and position - 1.
  - 8: RatePct: Rate (as explained above) expressed as a percentage.
- Effect of expected error and length filtering:**  
The first column indicates read lengths (L). The next four columns indicate the number of reads that would be retained by the -- command if the reads were truncated at length L (option -- L) and filtered to have a maximum expected error of 1.0, 0.5, 0.25 (with the option --fastq\_maxee float). The last four columns fraction of reads that would be retained by the --fastq\_filter using the same length and maximum expected error parameters.
- Effect of minimum quality and length filtering:**  
The first column indicates read lengths (Len). The next four indicate the fraction of reads that would be retained by the -- command if the reads were truncated at length Len (option -- clen Len) or at the first position with a quality Q equal to or than 5, 10, 15 or 20 (option --fastq\_truncqual Q).
- fastq\_stripleft positive integer  
When using --fastq\_filter or --fastx\_filter, strip the specified number of bases from the left end of the reads. If the length of the resulting read is null, then the read is discarded.
- fastq\_stripright positive integer  
When using --fastq\_filter or --fastx\_filter, strip the specified number of bases from the right end of the reads. If the length of the resulting read is null, then the read is discarded.
- fastq\_tail positive integer  
When using --fastq\_chars, count the number of times a series of

characters of

length k appears at the end of quality strings. By default, k = 4.

--fastq\_truncee real

When using --fastq\_filter or --fastx\_filter, truncate sequences so that their total expected error is not higher than the specified value.

--fastq\_truncee\_rate real

When using --fastq\_filter or --fastx\_filter, truncate sequences so that their average expected error per base is not higher than the specified value. The error per length will happen at the first occurrence. The average expected base is calculated as the total expected number of errors divided by the length of the sequence after truncation.

--fastq\_trunclen positive integer

When using --fastq\_filter or --fastx\_filter, truncate sequences to the specified length. Shorter sequences are discarded.

--fastq\_trunclen\_keep positive integer

When using --fastq\_filter or --fastx\_filter, truncate sequences to the specified length. Shorter sequences are not discarded.

--fastq\_truncqual positive integer

When using --fastq\_filter, --fastq\_mergepairs or --fastx\_filter, truncate sequences starting from the first base with the specified base quality score or lower.

--fastqout filename

When using --fastq\_filter, --fastq\_mergepairs, --fastx\_filter or --fasta2fastq, write to the given FASTQ-formatted file the sequences passing the filter, or the merged or converted sequences.

--fastqout\_rev filename

When using --fastq\_filter or --fastx\_filter, write to the given FASTQ-formatted file the reverse reads passing the filter.

--fastqout\_discarded filename

When using --fastq\_filter or --fastx\_filter, write sequences that do not pass the filter to the given FASTQ-formatted file.

--fastqout\_discarded\_rev filename

When using --fastq\_filter or --fastx\_filter, write reverse reads that do not pass the filter to the given FASTQ-formatted file.

--fastqout\_notmerged\_fwd filename

When using --fastq\_mergepairs, write forward reads not merged to the

specified

FASTQ file.

--fastqout\_notmerged\_rev filename  
When using --fastq\_mergepairs, write reverse reads not merged to the

specified

FASTQ file.

--fastx\_filter filename  
Trim and/or filter the sequences in the given FASTA or FASTQ file and  
remaining sequences to the FASTQ file specified with the --fastqout  
and/or to the FASTA file specified with the --fastaout option.

Discarded se-

fastaout\_discarded

automatically

reverse

output will

fastaout\_rev,

can not

sequences are

may be

fastq\_stripright,

fastq\_trunclen\_keep

options

fastq\_minlen

minsize. Se-

quences,

are dis-

writ-

format.

--eeout

sequence. Af-

discarded se-

trimmed.

accepted be-

quences are written to the files specified with the --  
and  
--fastqout\_discarded options. The input format (FASTA or FASTQ) is  
detected. If the input consists of paired sequences, an input file with  
reads may be specified with the --reverse option, and corresponding  
be written to the files specified with the --fastqout\_rev, --  
--fastqout\_discarded\_rev, and --fastaout\_discarded\_rev options. Output  
be written to FASTQ files if the input is in FASTA format. The  
first trimmed and then filtered based on the remaining bases. Sequences  
trimmed using the options --fastq\_stripleft, --  
--fastq\_truncee, --fastq\_truncee\_rate, --fastq\_trunclen, --  
and --fastq\_truncqual. The sequences may be filtered using the  
--fastq\_maxee, --fastq\_maxee\_rate, --fastq maxlen, --fastq\_maxns, --  
(default 1), --fastq\_minqual, --fastq\_trunclen, --maxsize, and --  
quences not satisfying the requirements are discarded. For pairs of  
both sequences in a pair must satisfy the requirements, otherwise both  
carded. If no shortening or filtering options are given, all sequences are  
ten to the output files, possibly after conversion from FASTQ to FASTA  
The --relabel option may be used to relabel the output sequences. The  
option may be used to output the expected number of errors in each  
ter all sequences have been processed, the number of kept and  
quences will be shown, as well as how many of the kept sequences were  
When the input is in FASTA format, the following options are not  
cause quality scores are not available: --eeout, --fastq\_ascii, --

```
fastq_eeout,          --fastq_maxee, --fastq_maxee_rate, --fastq_minqual, --fastq_out, --
fastq_qmax,           --fastq_qmin,   --fastq_truncee,    --fastq_truncee_rate,   --
fastq_truncqual,     --fastqout_discarded, --fastqout_discarded_rev, --fastqout_rev.

--fastx_revcomp filename
Reverse-complement the sequences in the given FASTA or FASTQ file to
specified with the --fastaout and/or --fastqout options. If the input file
FASTA format, the output can not be written back to a FASTQ file due to
base quality scores.

--join_padgap string
When running --fastq_join, use the string as a sequence padding string.
The de-
fault is NNNNNNNN (8 N's). Option accepts an empty string, or any other
combination of ASCII characters.

--join_padgapq string
When running --fastq_join, use the string as a quality padding string.
The de-
fault is a string of I's equal in length to the sequence padding
string. The
letter I corresponds to a base quality score of 40 indicating a very high
quality base with error probability of 0.0001. Option accepts an empty
string, or
any other combination of ASCII characters.

--lengthout
Write sequence length information to the output files in FASTA or FASTQ
format
by adding a ";length=integer" attribute in the header.

--maxsize positive integer
When using --fastq_filter or --fastx_filter, discard sequences with an
abundance
higher than the specified value.

--minsize positive integer
When using --fastq_filter or --fastx_filter, discard sequences with an
abundance
lower than the specified value.

--output filename
When using --fastq_eestats or --fastq_eestats2, write tabulated results to
file-
name. See --fastq_eestats's and --fastq_eestats2's documentation for a
complete
description of the table.

--relabel_keep
When using --relabel, keep the old identifier in the header after a space.

--relabel string
```

Please see the description of the same option under Chimera detection for details.

--relabel\_md5  
Please see the description of the same option under Chimera detection for details.

--relabel\_self  
Please see the description of the same option under Chimera detection for details.

--relabel\_shal  
Please see the description of the same option under Chimera detection for details.

--reverse filename  
When using --fastq\_filter, --fastx\_filter, --fastq\_mergepairs or --fastq\_join, specify the FASTQ file containing containing the reverse reads.

--sff\_convert filename  
Convert the given SFF file to FASTQ. The FASTQ output file is specified with the --fastqout option. The sequence may be clipped as specified in the SFF file if the option --sff\_clip is specified, otherwise no clipping occurs. Bases that would have been clipped are converted to lower case, while the rest is in upper case. The output quality encoding may be specified with the --sff\_qencoding option (default 33). The minimum and maximum output quality scores may be limited using the --fastq\_qminout and --fastq\_qmaxout options.

--sff\_clip  
Specifies that the sequences converted by the --sff\_convert command should be clipped in both ends as indicated in the SFF file. By default no clipping is performed.

--xlength  
Strip header attribute ";length=integer" from input sequences. This attribute is added to output sequences by the --lengthout option.

--xsize  
Strip abundance information from the headers when writing the output file.

--xee  
Strip information about expected errors (ee) from the output file headers. This information is added by the --fastq\_eeout and --eeout options.

Masking options:

masking is An input sequence can be composed of lower- or uppercase letters. When soft specified, lower case letters are treated as symbols that should be masked. Otherwise the case of the input sequences is ignored.

detection DUST masking is performed automatically by the commands for chimera (uchime\_denovo, uchime\_ref), clustering (cluster\_fast, cluster\_smallmem, cluster\_size), masking (maskfasta, fastx\_mask), pairwise alignment (allpairs\_global) and searching (search\_exact, usearch\_global).

is used Masking is usually specified with the --qmask option, while the --dbmask option for the database sequences specified with the --db option with the --usearch\_global, --search\_exact and --uchime\_ref commands.

argu- The argument to the --qmask and --dbmask option may be none, soft or dust. If the ment is none, the no masking is performed. If the argument is soft the lower case symbols are masked. Finally, if the argument is dust, the sequence is masked using the DUST algorithm by Tatusov and Lipman to mask low-complexity regions.

other- If the --hardmask option is specified, all masked regions are converted to N's, wise masked regions are indicated by lower case letters.

letters If any sequence is masked, the masked version of the sequence (with lower case or N's) is used in all output files. Otherwise the sequence is unmodified. The excep- tion is the sequences in the output file specified with the --uchimealns option, where the input sequences are converted to upper case first and lower case letters indicate disagreement between the aligned sequences.

with the The --qmask option (or --dbmask for database sequences) may be combined qmask or --hardmask option. The results of using the none, dust or soft argument to --dbmask are presented below, assuming each input sequence contains both lower and up- percase symbols.

Results if the --hardmask option is off (default):

none:	no masking, all symbols used, no change
dust:	masked symbols lowercased, rest uppercased
soft:	lowercase symbols masked, no case changes

Results if the --hardmask option is on:

none: no masking, all symbols used, no change  
dust: masked symbols changed to Ns, rest unchanged  
soft: lowercase symbols masked and changed to Ns

indices  
treated as  
other regions.

and do  
not contribute to an alignment.

--fastaout filename  
Write the masked sequences to filename, in fasta format. Applies only  
to the  
--fastx\_mask command.

--fastqout filename  
Write the masked sequences to filename, in fastq format. Applies only  
to the  
--fastx\_mask command.

--fastx\_mask filename  
Mask regions in sequences contained in the specified fasta or fastq  
file. The  
The out-  
mini-  
with the  
default is  
file is  
please use  
unmasked  
unmasked

default is to mask using DUST (use --qmask to modify that behaviour).  
put files are specified with the --fastaout and --fastqout options. The  
imum and maximum percentage of unmasked residues may be specified  
--min\_unmasked\_pct and --max\_unmasked\_pct options, respectively.

--hardmask  
Symbols in masked regions are replaced by N's. The default is to  
replace the  
masked regions by lower case letters.

--maskfasta filename  
Mask regions in sequences contained in the fasta file filename. The  
to mask using dust (use --qmask to modify that behaviour). The output  
specified with the --output option. This command is deprecated,  
--fastx\_mask instead.

--max\_unmasked\_pct real  
Discard sequences with more than the specified maximum percentage of  
residues. Works only with --fastx\_mask.

--min\_unmasked\_pct real  
Discard sequences with less than the specified minimum percentage of  
residues. Works only with --fastx\_mask.

--output filename  
to the  
algorithm  
default  
searching.  
sequence. If  
               Write the masked sequences to filename, in fasta format. Applies only  
               --mask\_fasta command.

--qmask none|dust|soft  
               If the argument is dust, mask regions in sequences using the DUST  
               that detects simple repeats and low-complexity regions. This is the  
               for chimera detection, clustering, masking, pairwise alignment, and  
               If the argument is soft, mask the lower case letters in the input  
               the argument is none, do not mask.

#### Orienting options:

either the  
specified  
refer-  
words with  
The cor-  
fastaout,  
was also  
words is  
--not-  
specified  
direction (+,  
matching  
used for  
There has  
to be  
speci-  
relabel\_md5,  
               The --orient command can be used to orient the sequences in a given file in  
               forward or the reverse complementary direction based on a reference database  
               with the --db option. The two strands of each input sequence are compared to the  
               ence database using nucleotide words. If one of the strands shares many more  
               at least one sequence in the database than the other, that strand is chosen.  
               rectly oriented sequences may be written to a FASTA file specified with the --  
               and to a FASTQ file specified with the --fastqout option (as long as the input  
               in FASTQ format). If the result is uncertain, because the number of matching  
               too similar, the original sequence is written to the file specified with the  
               matched option. The results may also be written to a tab-delimited text file  
               with the --tabbedout option. This file will contain the query label, the  
               - or ?), the number of matching words on the forward strand, and the number of  
               words on the reverse complementary strand. By default, a word length of 12 is  
               this command. The word length may be adjusted using the --wordlength option.  
               to be at least 4 times as many matches on one strand than the other for a strand  
               selected. In addition to the common options, the following options may also be  
               fied for this command: --dbmask, --qmask, --relabel, --relabel\_keep, --  
               --relabel\_self, --relabel\_shal, --sizein, and --sizeout.

--db filename  
FASTQ or  
with a  
               Read the reference database from the given file. It may be in FASTA,  
               UDB format. If an UDB file is used it should have been created

wordlength of 12.

--fastaout filename  
Write the correctly oriented sequences to filename, in fasta format.

--fastqout filename  
Write the correctly oriented sequences to filename, in fastq format.

--notmatched filename  
Write the sequences with undetermined direction to filename, in the original format.

--orient filename  
Orient the sequences in the given file.

--tabbedout filename  
Write the results to a tab-delimited text file with the specified filename.  
This file will contain the query label, the direction (+, - or ?), the number of matching words on the forward strand, and the number of matching words on the reverse complementary strand.

Pairwise alignment options:

The results of the  $n * (n-1) / 2$  pairwise alignments are written to the result files specified with --alnout, --blast6out, --fastapairs --matched, --notmatched, --samout, --tsegout, --uc or --userout (see Searching section below). Specify the --acceptall option to output all pairwise alignments, or specify an identity with --id to discard weak alignments. Most other accept/reject options (see options below) may also be used. Sequences are aligned on their plus strand only. Masking is performed as usual and specified with --qmask and --hardmask.

--acceptall  
Write the results of all alignments to output files. This option overrides all other accept/reject options (including --id).

--allpairs\_global filename  
Perform optimal global pairwise alignments of the fasta sequences contained in filename. Each sequence is compared to all sequences that come after it in the file, resulting in a total of  $n * (n-1) / 2$  pairwise alignments, where n is the total number of sequences. This command is multi-threaded.

--id real  
Reject the sequence match if the pairwise identity is lower than real (value ranging from 0.0 to 1.0 included).

--threads positive integer  
Number of computation threads to use (1 to 1024). The number of threads  
should be lesser or equal to the number of available CPU cores. The default is  
to use all available resources and to launch one thread per logical core.

--uc filename  
Output pairwise alignment results in filename using a tab-separated  
uclust-like format with 10 columns. Each sequence is compared to all other  
sequences, and all hits (--acceptall) or only some hits (--id float) are  
reported, with one pairwise comparison per line:

1. Record type, always set to 'H'.
2. Ordinal number of the target sequence (based on input order, starting from zero).
3. Sequence length.
4. Percentage of similarity with the target sequence.
5. Match orientation, always set to '+'.
6. Not used, always set to zero.
7. Not used, always set to zero.
8. Compact representation of the pairwise alignment using the CIGAR format (Compact Idiosyncratic Gapped Alignment (match/mismatch), D (deletion) and I (insertion). The '=' indicates that the query is identical to the centroid (ignoring terminal gaps).
9. Label of the query sequence.
10. Label of the target sequence.

#### Restriction site cutting options:

The input sequences in the file specified with the --cut command are cut into fragments at all restriction sites matching the pattern given with the --cut\_pattern option. The fragments on the forward strand are written to the file specified with the --fastaout file and the fragments on the reverse strand are written to the file specified with the --fastaout\_rev option. Input sequences that do not match are written to the file specified with the option --fastaout\_discarded, and their reverse complement are also writ-

options  
may be  
ten to the file specified with the --fastaout\_discarded\_rev option. The relabel  
(--relabel, --relabel\_self, --relabel\_keep, --relabel\_md5, and --relabel\_sh1)  
used to relabel the output sequences).

pattern is a  
that must  
characters  
posi-  
pattern  
palindromic  
output all  
be nec-  
sequences.

--cut filename  
Specify the input file with sequences in FASTA format.

--cut\_pattern string  
Specify the restriction site cutting pattern and positions. The  
string of lower- or uppercase letters specifying the nucleotides  
match, and may include ambiguous nucleotide symbols. The special  
"^" (circumflex) and "\_" (underscore) are used to indicate the cutting  
position on the forward and reverse strand, respectively. For example, the  
"G^AATT\_C" is the pattern for the EcoRI restriction site. For such  
patterns (identical to its reverse complement) the command will  
possible fragments on both strands. For non-palindromic sites, it may  
essary to run the command also on the reverse complemented input  
Exactly one cutting site on each strand must be indicated.

strand.

--fastaout filename  
Specify the output file for the resulting fragments on the forward

strand.

--fastaout\_rev filename  
Specify the output file for the resulting fragments on the reverse

complemented.

--fastaout\_discarded filename  
Specify the output file for the non-matching sequences.

--fastaout\_discarded\_rev filename  
Specify the output file for the non-matching sequences, reverse

Searching options:  
format.  
using mul-

--alnout filename  
Write pairwise global alignments to filename using a human-readable  
format.  
using mul-  
tiple threads.

format. The  
OTUs.

--biomout filename  
Write search results to an OTU table in the biom version 1.0 file  
query file contains the samples, while the database file contains the  
Sample and OTU identifiers are extracted from the header of these

sequences.

format of  
(or lack  
pair-  
common  
reported  
similar out  
query+tar-  
and de-

to ' \* '

100.0).  
columns  
other

columns).

inte-

positive

Always  
to ig-

Always  
otherwise (see

See the `--biomout` option in the Clustering section for further details.

--blast6out filename

Write search results to filename using a blast-like tab-separated twelve fields (listed below), with one line per query-target matching of matching if --output\_no\_hits is used). Warning, vsearch uses global wise alignments, not blast's seed-and-extend algorithm. Therefore, some blast output values (alignment start and end, evalue, bit score) are differently. Output order may vary when using multiple threads. A put can be obtain with --userout filename and --userfields get+id+alnlen+mism+opens+qlo+qhi+tlo+thi+evalue+bits. A complete list scription is available in the section 'Userfields' of this manual.

1. query: query label.
  2. target: target (database sequence) label. The field is set if there is no alignment.
  3. id: percentage of identity (real value ranging from 0.0 to 1.0). The percentage identity is defined as  $100 * (\text{matching length} / \text{alignment length})$ . See fields id0 to id4 for definitions.
  4. alnlen: length of the query-target alignment (number of matches). The field is set to 0 if there is no alignment.
  5. mism: number of mismatches in the alignment (zero or positive integer value).
  6. opens: number of columns containing a gap opening (zero or integer value, excluding terminal gaps).
  7. qlo: first nucleotide of the query aligned with the target. equal to 1 if there is an alignment, 0 otherwise (see qilo to ignore initial gaps).
  8. qhi: last nucleotide of the query aligned with the target. equal to the length of the pairwise alignment, 0 (see qhi to ignore terminal gaps).

Always

to ig-

Always

otherwise (see

alignments).

Always

target se-

pairwise

created us-

or the

search

file-

queries

";size=inte-

filename, in

fasta

sensitivity,

algorithm

9. tlo: first nucleotide of the target aligned with the query.  
equal to 1 if there is an alignment, 0 otherwise (see tlo  
ignore initial gaps).

10. thi: last nucleotide of the target aligned with the query.  
equal to the length of the pairwise alignment, 0  
thi to ignore terminal gaps).

11. eval: expectancy-value (not computed for nucleotide  
alignments).  
Always set to -1.

12. bits: bit score (not computed for nucleotide alignments).  
set to 0.

--db filename  
Compare query sequences (specified with --usearch\_global) to the  
sequences contained in filename in FASTA or FASTQ format, using global  
alignment. Alternatively, the name of a preformatted UDB database  
using the makeudb\_usearch command (see below) may be specified.

--dbmask none|dust|soft  
Mask regions in the target database sequences using the dust method  
soft method, or do not mask (none). Warning, when using soft masking  
commands become case sensitive. The default is to mask using dust.

--dbmatched filename  
Write database target sequences matching at least one query sequence to  
name, in fasta format. If the option --sizeout is used, the number of  
that matched each target sequence is indicated using the pattern  
ger;".

--dbnotmatched filename  
Write database target sequences not matching query sequences to  
filename, in  
fasta format.

--fastapairs filename  
Write pairwise alignments of query and target sequences to filename, in  
format.

--fulldp Dummy option for compatibility with usearch. To maximize search  
vsearch uses a 8-way 16-bit SIMD vectorized full dynamic programming  
(Needleman-Wunsch), whether or not --fulldp is specified.

--gapext string  
Set penalties for a gap extension. See --gapopen for a complete description of the penalty declaration system. The default is to initialize the six extending penalties using a penalty of 2 for extending internal gaps penalty of 1 for extending terminal gaps, in both query and target (i.e. 2I/1E).

--gapopen string  
Set penalties for a gap opening. A gap opening can occur in six texts: in the query (Q) or in the target (T) sequence, at the left right (R) extremity of the sequence, or inside the sequence (I).  
Sequence symbols (Q and T) can be combined with location symbols (L, I, and R), and numerical values to declare penalties for all possible aQL/bQI/cQR/dTL/eTI/fTR, where abcdef are zero or positive integers, is used as a separator.  
To simplify declarations, the location symbols (L, I, and R) can be the symbol (E) can be used to treat both extremities (L and R) the symbols Q and T can be omitted to treat query and target equally. For instance, the default is to declare a penalty of 20 for internal gaps and a penalty of 2 for opening terminal gaps (left or both query and target sequences (i.e. 20I/2E). If only a numerical given, without any sequence or location symbol, then the penalty all gap openings. To forbid gap-opening, an infinite penalty value can clared with the symbol '\*'. To use vsearch as a semi-global aligner, a penalty can be applied to the left (L) or right (R) gaps.  
vsearch always initializes the six gap opening penalties using the parameters (20I/2E). The user is then free to declare only the values wants to modify. The string is scanned from left to right, accepted are (0123456789/LIREQT\*), and later values override previous values. Please note that vsearch, in contrast to usearch, only allows penalties. Because the lowest gap penalties are 0.5 by default in default scores and gap penalties in vsearch have been doubled to equivalent penalties and to produce identical alignments.

--hardmask  
Mask sequence regions by replacing them with Ns instead of setting them to lower case as is the default. For more information, please see the Masking section.

--id real  
(value sequences sequence, efficient pre-weakly mers to from the minwordmatches 0.5 is identity is length - Reject the sequence match if the pairwise identity is lower than real ranging from 0.0 to 1.0 included). The search process sorts target by decreasing number of k-mers they have in common with the query using that information as a proxy for sequence similarity. That filtering also prevents pairwise alignments with very short, or with matching targets, as there needs to be by default at least 12 shared k-mers to start the pairwise alignment, and at least one out of every 16 k-mers query needs to match the target (see options --wordlength and --iddef to change that behaviour). Consequently, using values lower than --iddef is not likely to capture more weakly matching targets. The pairwise identity is by default defined as the number of (matching columns) / (alignment length - terminal gaps). That definition can be modified by --iddef.

--iddef 0|1|2|3|4  
Change the pairwise identity definition used in --id. Values accepted are:

length).  
0. CD-HIT definition: (matching columns) / (shortest sequence

1. edit distance: (matching columns) / (alignment length).

2. edit distance excluding terminal gaps (default  
--id).

3. Marine Biological Lab definition counting each gap opening (internal or terminal) as a single mismatch, whether or not the extended:  $1.0 - [(mismatches + gap openings) / (\text{longest length})]$

4. BLAST definition, equivalent to --iddef 1 for global alignments.

The option --userfields accepts the fields id0 to id4, in addition field id, to report the pairwise identity values corresponding to the

differ-

ent definitions.

--idprefix positive integer  
Reject the sequence match if the first integer nucleotides of the target do not match the query.

--idsuffix positive integer  
Reject the sequence match if the last integer nucleotides of the target do not match the query.

--lca\_cutoff real  
Adjust the fraction of matching hits required for the last common ancestor value is rank to fraction of The argument to this option must be larger than 0.5, but not larger than 1.0.

--lcaout filename  
Output last common ancestor (LCA) information about the hits of each query to query id, of the same for- "tax=k:Archaea,p:Eur- that are It is for this useful. The option --leftjust Reject the sequence match if the pairwise alignment begins with gaps.

--lengthout Write sequence length information to the output files in FASTA format by adding a ";length=integer" attribute in the header.

--match integer Score assigned to a match (i.e. identical nucleotides) in the pairwise align-

ment. The default value is 2.

**--matched filename**  
Write query sequences matching database target sequences to filename, in  
fasta  
format.

**--maxaccepts positive integer**  
Maximum number of matching target sequences to accept before  
stopping the  
in pair  
decreasing  
that in-  
alignments, if  
accepted as  
is set  
maxaccepts and  
--max-  
search  
substi-  
tutions, insertions or deletions.

**--maxdiffs positive integer**  
Reject the sequence match if the alignment contains more than integer  
inser-  
tions or deletions.

**--maxgaps positive integer**  
Reject the sequence match if the alignment contains more than integer  
query  
the plus  
targets  
However, when  
(one per  
hits per  
applies to  
output  
files.

**--maxhits non-negative integer**  
Maximum number of hits to show once the search is terminated for a given  
query  
the plus  
targets  
However, when  
(one per  
hits per  
applies to  
output  
files.  
Unlimited by default or if the argument is zero. This option  
--alnout, --blast6out, --fastapairs, --samout, --uc, or --userout

**--maxid real**

two sequences is greater than real.

--maxqsize positive integer  
Reject query sequences with an abundance greater than integer.

--maxqt real  
Reject if the query/target sequence length ratio is greater than real.

--maxrejects positive integer  
Maximum number of non-matching target sequences to consider before stopping the search for a given query. The default value is 32. This option pair with --maxaccepts. The search process sorts target sequences by increasing number of k-mers they have in common with the query sequence, information as a proxy for sequence similarity. After pairwise none of the first 32 examined target sequences pass the acceptation the search process stops for that query (no hit). If --maxrejects is higher value, more target sequences are considered. If --maxrejects are both set to 0, the complete database is searched.

stopping works in decreasing that alignments, if criteria, set to a maxaccepts and --maxsizeratio real  
Reject if the query/target abundance ratio is greater than real.

--maxsl real  
Reject if the shorter/longer sequence length ratio is greater than real.

integer  
--maxsubs positive integer  
Reject the sequence match if the pairwise alignment contains more than substitutions.

than real  
--mid real  
Reject the sequence match if the percentage of identity is lower (ignoring all gaps, internal and terminal).

integer.  
--mincols positive integer  
Reject the sequence match if the alignment length is shorter than

--minqt real  
Reject if the query/target sequence length ratio is lower than real.

--minsizeratio real  
Reject if the query/target abundance ratio is lower than real.

--minsl real  
Reject if the shorter/longer sequence length ratio is lower than real.

--mintsize positive integer  
Reject target sequences with an abundance lower than integer.

--minwordmatches non-negative integer  
Minimum number of k-mers or word matches required for a sequence to be considered further. Default value is 12 for the default word length 8.

For word lengths 3-15, the default minimum word matches are 18, 17, 16, 15, 14, 10, 9, 8, 7, 5 and 3, respectively. If the query sequence has fewer words than the number specified, all words in the query must match. If the argument is 0, no word matches are required.

--mismatch integer  
Score assigned to a mismatch (i.e. different nucleotides) in the pairwise alignment. The default value is -4.

--mothur\_shared\_out filename  
Write search results to an OTU table in the mothur 'shared' tab-separated plain text file format. The query file contains the samples, while the base file contains the OTUs. Sample and OTU identifiers are extracted from the header of these sequences. See the --otutabout option in the Clustering section for further details.

--notmatched filename  
Write query sequences not matching database target sequences to filename, in fasta format.

--otutabout filename  
Write search results to an OTU table in the classic tab-separated plain text format. The query file contains the samples, while the database file contains the OTUs. Sample and OTU identifiers are extracted from the header of these sequences (--sample option). See the --mothur\_shared\_out option in the Clustering section for further details.

--output\_no\_hits  
Write both matching and non-matching queries to --alnout, --blast6out, --samout or --userout output files. Non-matching queries are labelled 'No hits' in --alnout files.

--pattern string  
This option is ignored. It is provided for compatibility with usearch.

--qmask none|dust|soft  
Mask regions in the query sequences using the dust or the soft algorithms, or do not mask (none). Warning, when using soft masking search commands

become

format.

lower  
is com-  
terminal

Set to 0

header  
(link)  
samtools/hts-

separated text

header  
query-  
may vary  
and op-  
specs)  
the for-

hit), 4  
hit, i.e.

- case sensitive. The default is to mask using dust.
- qsegout filename  
Write the aligned part of each query sequence to filename in FASTA
- query\_cov real  
Reject if the fraction of the query aligned to the target sequence is lower than real (value ranging from 0.0 to 1.0 included). The query coverage is computed as (matches + mismatches) / query sequence length. Internal or gaps are not taken into account.
- rightjust  
Reject the sequence match if the pairwise alignment ends with gaps.
- rowlen positive integer  
Width of alignment lines in --alnout output. The default value is 64. Set to 0 to eliminate wrapping.
- samheader  
Include header lines to the SAM file when --samout is specified. The header includes lines starting with @HD, @SQ and @PG, but no @RG lines (see <https://github.com/samtools/hts-specs>) <<https://github.com/samtools/hts-specs>>). By default no header line is written.
- samout filename  
Write alignment results to filename using the SAM format (a tab-separated text file). When using the --samheader option, the SAM file starts with header lines. Each non-header line is a SAM record, which represents either a query-target alignment or the absence of match for a query (output order may vary when using multiple threads). Each record contains 11 mandatory fields and optional fields (see (link) <https://github.com/samtools/hts-specs> for a complete description of the format):
  1. query sequence label.
  2. combination of bitwise flags. Possible values are: 0 (top hit), 16 (reverse-complemented hit), 256 (secondary all hits except the top hit).
  3. target sequence label.

1 for  
of the  
query  
and al-  
the tar-

4. first position of a target aligned with the query (always global pairwise alignments, 0 if there is no match).  
5. mapping quality (ignored, always set to '\*').  
6. CIGAR string (set to '\*' if there is no match).  
7. name of the target sequence matching with the next read query (for mate reads only, ignored and always set to '\*').  
8. position of the primary alignment of the next read of the (for mate reads only, ignored and always set to 0).  
9. target sequence length (for multi-segment targets, ignored always set to 0).  
10. query sequence (complete, not only the segment aligned to get as usearch does).  
11. quality string (ignored, always set to '\*').

Optional fields for query-target matches (number and order of fields may vary):

12. AS:i:? alignment score (i.e. percentage of identity).
13. XN:i:? next best alignment score (always set to 0).
14. XM:i:? number of mismatches.
15. X0:i:? number of gap openings (excluding terminal gaps).
16. XG:i:? number of gap extensions (excluding terminal gaps).
17. NM:i:? edit distance to the target (sum of XM and XG).
18. MD:Z:? string for mismatching positions.
19. YT:Z:UU string representing the alignment type.

--search\_exact filename  
file-  
matches  
--id,  
search-  
iden-

Search for exact full-length matches to the query sequences contained in name in the fasta database of target sequences (--db). Only 100% exact are reported and this command is much faster than --usearch\_global. The --maxaccepts and --maxrejects options are ignored, but the rest of the ing options may be specified.

--self Reject the sequence match if the query and target labels are identical.  
--selfid Reject the sequence match if the query and target sequences are strictly

tical.

--sizeout  
(using the matched each (default) or aligned to computed as terminal gaps between the with the --uc, dbnot- highest per- is mea- same per- controlled by controlled by format.  
format layout is usearch\_global com- hit (N). best hit

Add abundance annotations to the output of the option --dbmatched pattern ';size=integer;'), to report the number of queries that target.

--strand plus|both  
When searching for similar sequences, check the plus strand only check both strands.

--target\_cov real  
Reject the sequence match if the fraction of the target sequence the query sequence is lower than real. The target coverage is (matches + mismatches) / target sequence length. Internal or are not taken into account.

--top\_hits\_only  
Only the top hits with an equally high percentage of identity query and database sequence sets are written to the output specified options --lcaout, --alnout, --samout, --userout, --blast6out, --fastapairs, --matched or --notmatched (but not --dbmatched and --matched). For each query, the top hit is the one presenting the centage of identity (see the --iddef option to change the way identity sured). For a given query, if several top hits present exactly the centage of identity, the number of matching targets reported is the --maxaccepts value (1 by default), and the number of hits is the --maxhits option.

--tsegout filename  
Write the aligned part of each target sequence to filename in FASTA format.

--uc filename  
Output searching results in filename using a tab-separated uclust-like with 10 columns. When using the --search\_exact command, the table the same than with the --allpairs\_global. When using the -- mand, the table present two different type of entries: hit (H) or no Each query sequence is compared to all other sequences, and the (--maxaccepts 1) or several hits (--maxaccepts > 1 and --uc\_allhits)

are re-  
content

start-

'\*' for

CIGAR

Report): M

equal sign

sequence

--uc op-

filename

file.

'Userfields' sec-

fields with  
threads. If

ported (H). Output order may vary when using multiple threads. Column varies with the type of entry (H or N):

1. Record type: H, or N ('hit' or 'no hit').
2. Ordinal number of the target sequence (based on input order, starting from zero). Set to '\*' for N.
3. Sequence length. Set to '\*' for N.
4. Percentage of similarity with the target sequence. Set to N.
5. Match orientation + or - . . Set to '.' for N.
6. Not used, always set to zero for H, or '\*' for N.
7. Not used, always set to zero for H, or '\*' for N.
8. Compact representation of the pairwise alignment using the format (Compact Idiosyncratic Gapped Alignment (match/mismatch), D (deletion) and I (insertion). The '=' indicates that the query is identical to the centroid (ignoring terminal gaps). Set to '\*' for N.
9. Label of the query sequence.
10. Label of the target centroid sequence. Set to '\*' for N.

--uc\_allhits

With the commands --search\_exact and --usearch\_global, when using the --option, show all hits, not just the top hit for each query.

--usearch\_global filename

Compare target sequences (--db) to the query sequences contained in FASTA or FASTQ format, using global pairwise alignment.

--userfields string

When using --userout, select and order the fields written to the output file. Fields are separated by '+' (e.g. query+target+id). See the section for a complete list of fields.

--userout filename

Write user-defined tab-separated output to filename. Select the option --userfields. Output order may vary when using multiple --userfields is empty or not present, filename is empty.

--weak\_id real  
terminating defined hits that high --id Logically, real must be smaller than the value indicated by --id.

--wordlength positive integer  
possible recommended. but can sensitivity increases with for very factor generally be-  
Length of words (i.e. k-mers) for database indexing. The range of values goes from 3 to 15, but values near 8 or 9 are generally Longer words may reduce the sensitivity/recall for weak similarities, increase precision. On the other hand, shorter words may increase or recall, but may reduce precision. Computation time generally shorter words and decreases with longer words, but it increases again long words. Memory requirements for a part of the index increase with a of 4 each time word length increases by one nucleotide, and this comes significant for long words (12 or more). The default value is 8.

--xlength  
attribute Strip header attribute ";length=integer" from input sequences. This is added to output sequences by the --lengthout option.

Shuffling options:  
Fasta entries in the input file are outputted in a pseudo-random order.

--lengthout  
adding Write sequence length information to the output files in FASTA format by a ";length=integer" attribute in the header.

--output filename  
Write the shuffled sequences to filename, in fasta format.

--randseed integer  
always pro- When shuffling sequence order, use integer as seed. A given seed duces the same output order (useful for replicability). Set to 0 to use a pseudo-random seed (default behaviour).

--relabel string  
to con- Relabel sequences using the prefix string and a ticker (1, 2, 3, etc.) struct the new headers. Use --sizeout to conserve the abundance

annotations.

--relabel\_keep  
When relabelling, keep the old identifier in the header after a space.

each se-  
converted to  
MD5 di-  
probability  
similar, but  
proba-  
gener-  
numbers  
abun-

--relabel\_md5  
Relabel sequences using the MD5 message digest algorithm applied to sequence. Former sequence headers are discarded. The sequence is upper case and U is replaced by T before the digest is computed. The digest is a cryptographic hash function designed to minimize the probability that two different inputs gives the same output, even for very non-identical inputs. Still, there is always a very small, but non-zero probability that two different inputs give the same result. The MD5 digest generates a 128-bit (16-byte) digest that is represented by 16 hexadecimal (using 32 symbols among 0123456789abcdef). Use --sizeout to conserve the abundance annotations.

--relabel\_self  
Relabel sequences using the sequence itself as the label.

each se-  
algorithm  
(20-byte)  
proba-  
digest) is  
sizeout

--relabel\_sha1  
Relabel sequences using the SHA1 message digest algorithm applied to sequence. It is similar to the --relabel\_md5 option but uses the SHA1 instead of the MD5 algorithm. The SHA1 digest generates a 160-bit result that is represented by 20 hexadecimal numbers (40 symbols). The probability of a collision (two non-identical sequences having the same smaller for the SHA1 algorithm than it is for the MD5 algorithm. Use -- to conserve the abundance annotations.

relabel\_shal, pre-  
(using the

--sizeout  
When using --relabel, --relabel\_self, --relabel\_md5 or -- serve and report abundance annotations to the output fasta file pattern ';size=integer;').

--shuffle filename  
Pseudo-randomly shuffle the order of sequences contained in filename.

--topn positive integer  
Output only the first integer sequences after pseudo-random reordering.

--xlength  
Strip header attribute ";length=integer" from input sequences. This attribute

is added to output sequences by the `--lengthout` option.

`--xsize` Strip abundance information from the headers when writing the output file.

Sorting options:  
length  
abun-  
or just  
assumes that  
performed  
during chimera checking (`--uchime_denovo`), dereplication (`--deref_fulllength`), and clustering (`--cluster_fast` and `--cluster_size`).

`--lengthout`  
adding  
greater  
smaller  
for de-  
for de-  
for de-  
for de-

Write sequence length information to the output files in FASTA format by adding a ";length=integer" attribute in the header.

`--maxsize positive integer`  
When using `--sortysize`, discard sequences with an abundance value greater than integer.

`--minsize positive integer`  
When using `--sortysize`, discard sequences with an abundance value smaller than integer.

`--output filename`  
Write the sorted sequences to filename, in fasta format.

`--relabel string`  
Please see the description of the same option under Chimera detection for details.

`--relabel_keep`  
When relabelling, keep the old identifier in the header after a space.

`--relabel_md5`  
Please see the description of the same option under Chimera detection for details.

`--relabel_self`  
Please see the description of the same option under Chimera detection for details.

`--relabel_shal`  
Please see the description of the same option under Chimera detection for details.

--sizeout  
When using --relabel, report abundance annotations to the output  
fasta file  
(using the pattern ';size=integer;').

--sortbylength filename  
Sort by decreasing length the sequences contained in filename. See the  
general  
long se-  
quences.  
(missing  
minsize and  
--maxsize to eliminate rare and dominant sequences.

--sortbysize filename  
Sort by decreasing abundance the sequences contained in filename  
abundance values are assumed to be ';size=1'). See the options --  
--topn positive integer  
Output only the top integer sequences (i.e. the longest or the most  
abundant).

--xlength  
Strip header attribute ";length=integer" from input sequences. This  
attribute  
is added to output sequences by the --lengthout option.

--xsize Strip abundance information from the headers when writing the output  
file.

Subsampling options:  
Subsampling randomly extracts a certain number or a certain percentage of the  
sequences  
input se-  
quences  
file.  
among the  
specified with  
fastaout and  
with the  
written to  
fastq\_discarded. The  
files specified with the options --fasta\_discarded and --  
--fastq\_ascii, --fastq\_qmin and --fastq\_qmax options are also available.

--fastaout filename  
Write the sampled sequences to filename, in fasta format.

--fastaout\_discarded filename  
Write the sequences not sampled to filename, in fasta format.

--fastq\_ascii positive integer  
Define the ASCII character number used as the basis for the FASTQ quality score. The default is 33, which is used by the Sanger / Illumina 1.8+ FASTQ format (phred+33). The value 64 is used by the Solexa, Illumina 1.3+ and Illumina 1.5+ formats (phred+64). Only 33 and 64 are valid arguments.

--fastq\_qmax positive integer  
Specify the maximum quality score accepted when reading FASTQ files. The default is 41, which is usual for recent Sanger/Illumina 1.8+ files.

--fastq\_qmin positive integer  
Specify the minimum quality score accepted for FASTQ files. The default is 0, which is usual for recent Sanger/Illumina 1.8+ files. Older formats may use scores between -5 and 2.

--fastqout filename  
Write the sampled sequences to filename, in fastq format. Requires input in fastq format.

--fastqout\_discarded filename  
Write the sequences not sampled to filename, in fastq format. Requires input in fastq format.

--fastx\_subsample filename  
Perform subsampling from the sequences in the specified input file that is in FASTA or FASTQ format.

--lengthout  
Write sequence length information to the output files in FASTA format by adding a ";length=integer" attribute in the header.

--randseed positive integer  
When subsampling, use integer as a seed for the pseudo-random generator. A given seed always produces the same output, which is useful for replicability. Set to 0 to use a pseudo-random seed (default behaviour).

--relabel string  
Relabel sequences using the prefix string and a ticker (1, 2, 3, etc.) to construct the new headers. Use --sizeout to conserve the abundance annotations.

--relabel\_keep  
When relabelling, keep the old identifier in the header after a space.

--relabel\_md5  
Relabel sequences using the MD5 message digest algorithm applied to

each sequence. Former sequence headers are discarded. The sequence is converted to upper case and U is replaced by T before the digest is computed. The MD5 digest is a cryptographic hash function designed to minimize the probability that two different inputs give the same output, even for very non-identical inputs. Still, there is always a very small, but non-zero probability that two different inputs give the same result. The MD5 generates a 128-bit (16-byte) digest that is represented by 16 hexadecimal (using 32 symbols among 0123456789abcdef). Use --sizeout to conserve the abundance annotations.

--relabel\_self  
Relabel sequences using the sequence itself as the label.

--relabel\_sha1  
Relabel sequences using the SHA1 message digest algorithm applied to each sequence. It is similar to the --relabel\_md5 option but uses the SHA1 instead of the MD5 algorithm. The SHA1 digest generates a 160-bit result that is represented by 20 hexadecimal numbers (40 symbols). The probability of a collision (two non-identical sequences having the same smaller for the SHA1 algorithm than it is for the MD5 algorithm. Use --sizeout to conserve the abundance annotations.

--sample\_pct real  
Subsample the given percentage of the input sequences. Accepted values range from 0.0 to 100.0.

--sample\_size positive integer  
Extract the given number of sequences.

--sizein Take the abundance information of the input file into account, otherwise the abundance of each sequence is considered to be 1.

--sizeout  
Write abundance information to the output file.

--xlength  
Strip header attribute ";length=integer" from input sequences. This attribute is added to output sequences by the --lengthout option.

--xsize Strip abundance information from the headers when writing the output file.

Taxonomic classification options:  
The vsearch command --syntax will classify the input sequences according to the  
Sintax algorithm as described by Robert Edgar (2016) in SINTAX: a simple non-Bayesian  
taxonomy classifier for 16S and ITS sequences, BioRxiv, 074161. Preprint. doi:  
10.1101/074161  
(link) (<https://doi.org/10.1101/074161>)

given as  
specified with  
name is  
set a  
--rand-  
number  
threads, the  
in a  
single  
randseed.  
  
option  
  
each se-  
separated  
with an  
(phylum), c  
letter is  
not al-  
names.

The name of the fasta file containing the input sequences to be classified is  
an argument to the --syntax command. The reference sequence database is  
specified with the --db option. The results are written in a tab delimited text file whose  
minimum level of bootstrap support for the taxonomic ranks to be reported. The  
seed option may be included to specify a seed for initialisation of the random  
generator used by the algorithm. Please note that when using multiple  
--randseed option may not work as intended, because sequences may be processed  
random order by different threads. To ensure the same results each time, use a  
thread --threads 1) in combination with a fixed random seed specified with --

Multithreading is supported. Databases in UDB files are supported. The strand  
may be specified.

The reference database must contain taxonomic information in the header of  
sequence in the form of a string starting with ";tax=" and followed by a comma-  
list of up to nine taxonomic identifiers. Each taxonomic identifier must start  
indication of the rank by one of the letters d (for domain) k (kingdom), p  
(class), o (order), f (family), g (genus), s (species), or t (strain). The  
followed by a colon (:) and the name of that rank. Commas and semicolons are  
allowed in the name of the rank. Non-ascii characters should be avoided in the  
names.

Example:

```
>X80725_S000004313;tax=d:Bacteria,p:Proteobacteria,c:Gammaproteobacteria,o:Enterobacte-  
rialles,f:Enterobacteriaceae,g:Escherichia/  
Shigella,s:Escherichia_coli,t:str._K-12_sub-  
str._MG1655
```

spaces in  
The option --notrunclabels is turned on by default for this command, allowing  
the taxonomic identifiers.

with the long, the option will in-

If two sequences in the reference database has equally many kmer matches query, the shortest sequence will be chosen by default. If they are equally sequence appearing first in the database will be chosen. If the recommended --syntax\_random is specified, sequences with an equal number of kmer matches stead be chosen by a random draw.

format.

--db filename  
Read the reference sequences from filename, in FASTA, FASTQ or UDB

These sequences need to be annotated with taxonomy.  
algo-  
replica-  
Does not correct

--randseed positive integer  
Use integer as seed for the random number generator used in the Syntax rithm. A given seed always produces the same output order (useful for bility). Set to 0 to use a pseudo-random seed (default behaviour). work correctly with multiple threads; please use --threads 1 to ensure behaviour.

that will corresponding to

--syntax filename  
Read the input sequences from filename, in FASTA or FASTQ format.  
--syntax\_cutoff real  
Specify a minimum level of bootstrap support for the taxonomic ranks be included in column 4 of the output file. For instance 0.9, 90%.

draw.

--syntax\_random  
Break ties between sequences with equally many kmer matches by a random draw.

This option is recommended and may be made the default in the future.  
1 con-  
the same  
paren-  
syntax\_cutoff  
while  
at or

--tabbedout filename  
Write the results to filename, in a tab-separated text format. Column tains the query label. Column 2 contains the predicted taxonomy in format as for the reference data, with bootstrap support indicated in theses after each rank. Column 3 contains the strand. If the -- option is used, the predicted taxonomy will be repeated in column 4 omitting the bootstrap values and including only the ranks with support above the threshold.

UDB options:  
Databases to be used with the --usearch\_global or --syntax commands may be prepared

from FASTA files and stored to a binary UDB formatted file in order to speed up search- sequences commands. This may be worthwhile when searching a large database repeatedly. The are indexed and stored in a way that can be quickly loaded into memory. The and options below can be used to create and inspect UDB files.

**--dbmask none|dust|soft**  
Specify the sequence masking method used with the **--makeudb\_usearch**

command, specified.

low com- letters in  
When dust is specified, the DUST algorithm will be used for masking plexity regions (short repeats and skewed composition). Lower case the input file will be masked when soft is specified (soft masking).

**--hardmask**  
Mask sequences by replacing letters with N for the **--makeudb\_usearch**

command.  
The default is to use lower case letters (soft masking).

**--makeudb\_usearch filename**  
Create an UDB database file from the FASTA-formatted sequences in the file specified with the given filename. The UDB database is written to the file with the **--output** option.

**--output filename**  
Specify the filename of a FASTA or UDB output file for the **--makeudb\_usearch** or the **--fdb2fasta** command, respectively.

**--fdb2fasta filename**  
Read the UDB database in the file with the given filename and output the sequences in FASTA format in the file specified by the **--output** option.

**--fdbinfo filename**  
Show information about the UDB database in the file with the given filename.

**--fdbstats filename**  
Report statistics about the indexed words in the UDB database in the file with the given filename.

**--wordlength positive integer**  
Specify the length of the words to be used when creating the UDB database in 15. The dex using the **--makeudb\_usearch** command. Valid numbers range from 3 to default is 8.

Userfields (fields accepted by the **--userfields** option):

aln Print a string of M (match/mismatch, i.e. not a gap), D (delete, i.e. a

gap in pairwise field	alnlen	Print the length of the query-target alignment (number of columns). The is set to 0 if there is no alignment.
(Com- (deletion) identical to is no value).	bits	Bit score (not computed for nucleotide alignments). Always set to 0.
excluding	caln	Compact representation of the pairwise alignment using the CIGAR format pact Idiosyncratic Gapped Alignment Report): M (match/mismatch), D and I (insertion). The equal sign '=' indicates that the query is the centroid sequence (ignoring terminal gaps). Empty field if there alignment.
specified by default the from 0.0 alignment defined defined value terminal) using the 100 *	evaluate	E-value (not computed for nucleotide alignments). Always set to -1.
	exts	Number of columns containing a gap extension (zero or positive integer value).
	gaps	Number of columns containing a gap (zero or positive integer value, terminal gaps).
	id	The percentage of identity, according to the identity definition the --iddef option. Equal to id0, id1, id2, id3 or id4 below. By same as id2.
	id0	CD-HIT definition of the percentage of identity (real value ranging to 100.0) using the length of the shortest sequence in the pairwise as denominator: $100 * (\text{matching columns}) / (\text{shortest sequence length})$ .
	id1	The percentage of identity (real value ranging from 0.0 to 100.0) is as the edit distance: $100 * (\text{matching columns}) / (\text{alignment length})$ .
	id2	The percentage of identity (real value ranging from 0.0 to 100.0) is as the edit distance, excluding terminal gaps.
	id3	Marine Biological Lab definition of the percentage of identity (real ranging from 0.0 to 100.0), counting each gap opening (internal or as a single mismatch, whether or not the gap was extended, and length of the longest sequence in the pairwise alignment as denominator: $(1.0 - [(\text{mismatches} + \text{gaps}) / (\text{longest sequence length})])$ .

0.0 to alignment. The value, to the positive in- alignment sequences, this 0.0 to this is (real 100.0 * gaps are alignment. is not to the terminal gaps). is no	id4 BLAST definition of the percentage of identity (real value ranging from 100.0), equivalent to --iddef 1 in a context of global pairwise alignment. The field id4 is always equal to the field id1.
	ids Number of matches in the alignment (zero or positive integer value).
	mism Number of mismatches in the alignment (zero or positive integer value).
	opens Number of columns containing a gap opening (zero or positive integer excluding terminal gaps).
	pairs Number of columns containing only nucleotides. That value corresponds length of the alignment minus the gap-containing columns (zero or integer value).
	pctgaps Number of columns containing gaps expressed as a percentage of the length (real value ranging from 0.0 to 100.0).
	pctpv Percentage of positive columns. When working with nucleotide sequences, this is equivalent to the percentage of matches (real value ranging from 100.0).
	pv Number of positive columns. When working with nucleotide sequences, equivalent to the number of matches (zero or positive integer value).
	qcov Fraction of the query sequence that is aligned with the target sequence value ranging from 0.0 to 100.0). The query coverage is computed as (matches + mismatches) / query sequence length. Internal or terminal not taken into account. The field is set to 0.0 if there is no alignment.
	qframe Query frame (-3 to +3). That field only concerns coding sequences and computed by vsearch. Always set to +0.
	qhi Last nucleotide of the query aligned with the target. Always equal length of the pairwise alignment, 0 otherwise (see qhi to ignore gaps).
	qihi Last nucleotide of the query aligned with the target (ignoring terminal Nucleotide numbering starts from 1. The field is set to 0 if there is no alignment).
	qilo First nucleotide of the query aligned with the target (ignoring initial

gaps). Nucleotide numbering starts from 1. The field is set to 0 if there is no alignment.

there ql Query sequence length (positive integer value). The field is set to 0 if 1 if there is no alignment.

alignment qlo First nucleotide of the query aligned with the target. Always equal to alignment, 0 otherwise (see qilo to ignore initial gaps).

alignment qrow Print the sequence of the query segment as seen in the pairwise alignment (i.e. with gap insertions if need be). Empty field if there is no alignment.

field if qs Query segment length. Always equal to query sequence length.

score is qstrand Query strand orientation (+ or - for nucleotide sequences). Empty exten- there is no alignment.

(real 100.0 \* terminal gaps alignment. query Query label.

target raw Raw alignment score (negative, null or positive integer value). The the sum of match rewards minus mismatch penalties, gap openings and gap sions. The field is set to 0 if there is no alignment.

target tcov Fraction of the target sequence that is aligned with the query sequence value ranging from 0.0 to 100.0). The target coverage is computed as (matches + mismatches) / target sequence length. Internal or are not taken into account. The field is set to 0.0 if there is no alignment.

is not tframe Target frame (-3 to +3). That field only concerns coding sequences and computed by vsearch. Always set to +0.

to the thi Last nucleotide of the target aligned with the query. Always equal terminal length of the pairwise alignment, 0 otherwise (see tihi to ignore gaps).

gaps). tihi Last nucleotide of the target aligned with the query (ignoring terminal Nucleotide numbering starts from 1. The field is set to 0 if there alignment.

is no tilo First nucleotide of the target aligned with the query (ignoring initial

gaps). Nucleotide numbering starts from 1. The field is set to 0 if there is no alignment.

0 if tl Target sequence length (positive integer value). The field is set to there is no alignment.

to 1 if tlo First nucleotide of the target aligned with the query. Always equal there is an alignment, 0 otherwise (see tilo to ignore initial gaps).

alignment trow Print the sequence of the target segment as seen in the pairwise alignment (i.e. with gap insertions if need be). Empty field if there is no alignment.

field is ts Target segment length. Always equal to target sequence length. The set to 0 if there is no alignment.

set to tstrand Target strand orientation (+ or - for nucleotide sequences). Always '+' , so reverse strand matches have tstrand '+' and qstrand '-'. Empty field if there is no alignment.

#### DELIBERATE CHANGES

If you are a usearch user, our objective is to make you feel at home. That's why vsearch was designed to behave like usearch, to some extent. Like any complex software, usearch is not free from quirks and inconsistencies. We decided not to reproduce some of them, and for complete transparency, to document here the deliberate changes we made.

During a search with usearch, when using the options --blast6out and --output\_no\_hits, for queries with no match the number of fields reported is 13, where it should be 12. This is corrected in vsearch.

The field raw of the --userfields option is not informative in usearch. This is corrected in vsearch.

The fields qlo, qhi, tlo, thi now have counterparts (qilo, qihi, tilo, tihi) reporting alignment coordinates ignoring terminal gaps.

In usearch, when using the option --output\_no\_hits, queries that receive no match are reported in --blast6out file, but not in the alignment output file. This is corrected in vsearch.

vsearch introduces a new --cluster\_size command that sorts sequences by decreasing abundance before clustering.

vsearch reintroduces --iddef alternative pairwise identity definitions that were

removed from  
usearch.

vsearch extends the --topn option to sorting commands.

vsearch extends the --sizein option to dereplication (--deref\_fulllength) and clustering (--cluster\_fast).

vsearch treats T and U as identical nucleotides during dereplication.

vsearch sorting is stabilized by using sequence abundances or sequences labels as secondary or tertiary keys.

vsearch by default uses the DUST algorithm for masking low-complexity regions. Masking behaviour is also slightly changed to be more consistent.

## NOVELTIES

vsearch introduces new commands and new options not present in usearch 7. They are described

in the 'Options' section of this manual. Here is a short list:

check-  
ing)

(clustering)

- uchime2\_denovo, uchime3\_denovo, alignwidth, borderline, fasta\_score (chimera checking)
- cluster\_size, cluster\_unequal, clusterout\_id, clusterout\_sort, profile (clustering)
- fasta\_width, gzip\_decompress, bzip2\_decompress (general option)
- iddef (clustering, pairwise alignment, searching)
- maxuniquesize (dereplication)
- relabel\_md5, relabel\_self and relabel\_shal (chimera detection, dereplication, processing, shuffling, sorting)
- shuffle (shuffling)
- fastq\_eestats, fastq\_eestats2, fastq\_maxlen, fastq\_trunce (FASTQ processing)
- fastqout\_discarded, fastqout\_discarded (subsampling)
- rereplicate (dereplication/rereplication)

## FASTQ

## EXAMPLES

Align all sequences in a database with each other and output all pairwise alignments:

```
vsearch --allpairs_global database.fas --alnout results.aln --acceptall
```

abun-  
Check for the presence of chimeras (de novo); parents should be at least 1.5 times more abundant than chimeras. Output non-chimeric sequences in fasta format (no wrapping):

```
vsearch --uchime_denovo queries.fas --abskew 1.5 --nonchimeras  
results.fas
```

```
--fasta_width 0
```

Cluster with a 97% similarity threshold, collect cluster centroids, and write cluster descriptions using a uclust-like format:

```
vsearch --cluster_fast queries.fas --id 0.97 --centroids centroids.fas --uc clusters.uc
```

Dereplicate the sequences contained in queries.fas, take into account the abundance information already present, write unwrapped fasta sequences to queries\_unique.fas with the new abundance information, discard all sequences with an abundance of 1:

```
vsearch --derep_fulllength queries.fas --sizein --fasta_width 0 --sizeout --output queries_unique.fas --minuniquesize 2
```

Mask simple repeats and low complexity regions in the input fasta file with the DUST algorithm

(masked regions are lowercased), and write the results to the output file:

```
vsearch --maskfasta queries.fas --qmask dust --output queries_masked.fas
```

Search queries in a reference database, with a 80%-similarity threshold, take terminal gaps

into account when calculating pairwise similarities, output pairwise alignments:

```
vsearch --usearch_global queries.fas --db references.fas --id 0.8 --iddef 1 --alnout results.aln
```

Search a sequence dataset against itself (ignore self hits), get all matches with at least 60%

similarity, and collect results in a blast-like tab-separated format. Accept an unlimited number of hits (--maxaccepts 0), and compare each query to all other sequences, including unlikely candidates (--maxrejects 0):

```
vsearch --usearch_global queries.fas --db queries.fas --self --id 0.6 --blast6out results.blast6 --maxaccepts 0 --maxrejects 0
```

Shuffle the input fasta file (change the order of sequences) in a repeatable fashion (fixed seed), and write unwrapped fasta sequences to the output file:

```
vsearch --shuffle queries.fas --output queries_shuffled.fas --randseed 13 --fasta_width 0
```

Sort by decreasing abundance the sequences contained in queries.fas (using the 'size=integer' information), relabel the sequences while preserving the abundance information (with --size-out), keep only sequences with an abundance equal to or greater than 2:

```
vsearch --sortbysize queries.fas --output queries_sorted.fas --relabel sampleA_
```

--size-

out --minsize 2

#### AUTHORS

Implementation and documentation by Torbjørn Rognes, Frédéric Mahé and Tomás Flouri.

#### CITATION

Rognes T, Flouri T, Nichols B, Quince C, Mahé F. (2016) VSEARCH: a versatile open source tool for metagenomics. PeerJ 4:e2584 doi: 10.7717/peerj.2584 (link) (<https://doi.org/10.7717/peerj.2584>)

#### REPORTING BUGS

Submit suggestions and bug-reports at (link) (<https://github.com/torognes/vsearch/issues>) <<https://github.com/torognes/vsearch/issues>>, send a pull request on (link) (<https://github.com/torognes/vsearch>) <<https://github.com/torognes/vsearch>>, or compose a friendly or curmudgeon e-mail to Torbjørn Rognes (link) ([torognes@ifi.uio.no](mailto:torognes@ifi.uio.no) <[torognes@ifi.uio.no](mailto:torognes@ifi.uio.no)>).

#### AVAILABILITY

Source code and binaries are available at <<https://github.com/torognes/vsearch>>.

#### COPYRIGHT

Copyright (C) 2014-2025, Torbjørn Rognes, Frédéric Mahé and Tomás Flouri

All rights reserved.

Contact: Torbjørn Rognes <[torognes@ifi.uio.no](mailto:torognes@ifi.uio.no)>, Department of Informatics, University of Oslo, PO Box 1080 Blindern, NO-0316 Oslo, Norway

This software is dual-licensed and available under a choice of one of two licenses, either under the terms of the GNU General Public License version 3 or the BSD 2-Clause License.

GNU General Public License version 3

This program is free software: you can redistribute it and/or modify it under the terms of the GNU General Public License as published by the Free Software Foundation, either version 3 of the License, or (at your option) any later version.

This program is distributed in the hope that it will be useful, but WITHOUT ANY WARRANTY;

without even the implied warranty of MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the GNU General Public License for more details.

You should have received a copy of the GNU General Public License along with this program. If not, see (link) (<https://www.gnu.org/licenses/>) <<https://www.gnu.org/licenses/>>.

The BSD 2-Clause License

Redistribution and use in source and binary forms, with or without modification, are permitted

provided that the following conditions are met:

1. Redistributions of source code must retain the above copyright notice, this list of conditions and the following disclaimer.
2. Redistributions in binary form must reproduce the above copyright notice, this list of conditions and the following disclaimer in the documentation and/or other materials provided with the distribution.

THIS SOFTWARE IS PROVIDED BY THE COPYRIGHT HOLDERS AND CONTRIBUTORS "AS IS" AND ANY EXPRESS OR

IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY

AND FITNESS FOR A PARTICULAR PURPOSE ARE DISCLAIMED. IN NO EVENT SHALL THE COPYRIGHT HOLDER OR

CONTRIBUTORS BE LIABLE FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL, EXEMPLARY, OR CONSEQUEN-

TIAL DAMAGES (INCLUDING, BUT NOT LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES; LOSS

OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION) HOWEVER CAUSED AND ON ANY THEORY OF LIA-

BILITY, WHETHER IN CONTRACT, STRICT LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE)

ARISING IN ANY WAY OUT OF THE USE OF THIS SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE.

We would like to thank the authors of the following projects for making their source code available:

- vsearch includes code from Google's CityHash project by Geoff Pike and Jyrki Alakui-jala, providing some excellent hash functions available under a MIT license.
- vsearch includes code derived from Tatusov and Lipman's DUST program that is in the public domain.
- vsearch includes public domain code written by Alexander Peslyak for the MD5 message digest algorithm.
- vsearch includes public domain code written by Steve Reid and others for the SHA1 message digest algorithm.
- vsearch binaries may include code from the zlib library, copyright Jean-Loup Gailly and Mark Adler.
- vsearch binaries may include code from the bzip2 library, copyright Julian R. Seward.

#### SEE ALSO

swipe, an extremely fast pairwise local (Smith-Waterman) database search tool by Torbjørn

Rognes, available at (link) (<https://github.com/torognes/swipe>)  
<<https://github.com/torognes/swipe>>.

Rognes, swarm, a fast and accurate amplicon clustering method by Frédéric Mahé and Torbjørn  
available at (link) (<https://github.com/torognes/swarm>) <<https://github.com/torognes/swarm>>.

## VERSION HISTORY

New features and important modifications of vsearch (short lived or minor bug releases  
may not be mentioned):

v1.0.0 released November 28th, 2014  
First public release.

v1.0.1 released December 1st, 2014  
changes Bug fixes (sortbysize, semicolon after size annotation in headers) and minor  
derepli- (labels as secondary sort key for most sorts, treat T and U as identical for  
cation, only output size in --dbmatched file if --sizeout specified).

v1.0.2 released December 6th, 2014  
Bug fixes (ssse3/sse4.1 requirement, memory leak).

v1.0.3 released December 6th, 2014  
Bug fix (now writes help to stdout instead of stderr).

v1.0.4 released December 8th, 2014  
eliminate mem- Added --allpairs\_global option. Reduce memory requirements slightly and  
ory leaks.

v1.0.5 released December 9th, 2014  
Fixes a minor bug with --allpairs\_global and --acceptall options.

v1.0.6 released December 14th, 2014  
Fixes a memory allocation bug in chimera detection (--uchime\_ref option).

v1.0.7 released December 19th, 2014  
Fixes a bug in the output from chimera detection with the --uchimeout option.

v1.0.8 released January 22nd, 2015  
Introduces several changes and bug fixes:

nucleotides, - a new linear memory aligner for alignment of sequences longer than 5,000  
before - a new --cluster\_size command that sorts sequences by decreasing abundance  
clustering,  
usearch, - meaning of userfields qlo, qhi, tlo, thi changed for compatibility with  
terminal - new userfields qilo, qihi, tilo, tihi give alignment coordinates ignoring  
gaps,

- in --uc output files, a perfect alignment is indicated with a '=' sign,
- the option --cluster\_fast now sorts sequences by decreasing length, then by decreasing abundance and finally by sequence identifier,
- default --maxseqlength value set to 50,000 nucleotides,
- fix for bug in alignment in rare cases,
- fix for lack of detection of under- or overflow in SIMD aligner.

v1.0.9 released January 22nd, 2015

decreas- Fixes a bug in the function sorting sequences by decreasing abundance (-- sortbysize).

v1.0.10 released January 23rd, 2015

affecting Fixes a bug where the --sizein option was ignored and always treated as on, clustering and derePLICATION commands.

v1.0.11 released February 5th, 2015

pairwise Introduces the possibility to output results in SAM format (for clustering, alignment and searching).

v1.0.12 released February 6th, 2015

Temporarily fixes a problem with long headers in FASTA files.

v1.0.13 released February 17th, 2015

with the Fix a memory allocation problem when computing multiple sequence alignments buffer --msaout and --consout options, as well as a memory leak. Also increased line for reading FASTA files to 4MB.

v1.0.14 released February 17th, 2015

clustering Fix a bug where the multiple alignment and consensus sequence computed after reading ignored the strand of the sequences. Also decreased size of line buffer for FASTA files to 1MB again due to excessive stack memory usage.

v1.0.15 released February 18th, 2015

defini- Fix bug in calculation of identity metric between sequences when using the MBL tion (--iddef 3).

v1.0.16 released February 19th, 2015

Integrated patches from Debian for increased compatibility with various architectures.

v1.1.0 released February 20th, 2015

warn- Added the --quiet option to suppress all output to stdout and stderr except for ings and fatal errors. Added the --log option to write messages to a log file.

v1.1.1 released February 20th, 2015

Added info about --log and --quiet options to help text.

v1.1.2 released March 18th, 2015  
Fix bug with large datasets. Fix format of help info.

v1.1.3 released March 18th, 2015  
Fix more bugs with large datasets.

v1.2.0-1.2.19 released July 6th to September 8th, 2015  
Several new commands and options added. Bugs fixed. Documentation updated.

v1.3.0 released September 9th, 2015  
Changed to autotools build system.

v1.3.1 released September 14th, 2015  
Several new commands and options. Bug fixes.

v1.3.2 released September 15th, 2015  
Fixed memory leaks. Added '-h' shortcut for help. Removed extra 'v' in version number.

v1.3.3 released September 15th, 2015  
Fixed bug in hexadecimal digits of MD5 and SHA1 digests. Added --samheader option.

v1.3.4 released September 16th, 2015  
Fixed compilation problems with zlib and bzip2lib.

v1.3.5 released September 17th, 2015  
Minor configuration/makefile changes to compile to native CPU and simplify makefile.

v1.4.0 released September 25th, 2015  
Added --sizeorder option.

v1.4.1 released September 29th, 2015  
Inserted public domain MD5 and SHA1 code to eliminate dependency on crypto and openssl libraries and their licensing issues.

v1.4.2 released October 2nd, 2015  
Dynamic loading of libraries for reading gzip and bzip2 compressed files if available.  
Circumvention of missing gzoffset function in zlib 1.2.3 and earlier.

v1.4.3 released October 3rd, 2015  
Fix a bug with determining amount of memory on some versions of Apple OS X.

v1.4.4 released October 3rd, 2015  
Remove debug message.

v1.4.5 released October 6th, 2015  
Fix memory allocation bug when reading long FASTA sequences.

v1.4.6 released October 6th, 2015  
Fix subtle bug in SIMD alignment code that reduced accuracy.

v1.4.7 released October 7th, 2015  
Fixes a problem with searching for or clustering sequences with repeats. In this new version, vsearch looks at all words occurring at least once in the sequences

in the initial step. Previously only words occurring exactly once were considered. In addition, vsearch now requires at least 10 words to be shared by the sequences, previously only 6 were required. If the query contains less than 10 words, all words must be present for a match. This change seems to lead to slightly reduced recall, but somewhat increased precision, ending up with slightly improved overall accuracy.

v1.5.0 released October 7th, 2015  
This version introduces the new option --minwordmatches that allows the user to specify the minimum number of matching unique words before a sequence is considered further. New default values for different word lengths are also set. The minimum word length is increased to 7.

v1.6.0 released October 9th, 2015  
This version adds the relabeling options (--relabel, --relabel\_md5 and --relabel\_shal) to the shuffle command. It also adds the --xsize option to the clustering, dereplication, shuffling and sorting commands.

v1.6.1 released October 14th, 2015  
Fix bugs and update manual and help text regarding relabelling. Add all relabelling options to the subsampling command. Add the --xsize option to chimera detection, dereplication and fastq filtering commands. Refactoring of code.

v1.7.0 released October 14th, 2015  
Add --relabel\_keep option.

v1.8.0 released October 19th, 2015  
Added --search\_exact, --fastx\_mask and --fastq\_convert commands. Changed most commands to read FASTQ input files as well as FASTA files. Modified --fastx\_revcomp and --fastx\_subsample to write FASTQ files.

v1.8.1 released November 2nd, 2015  
Fixes for compatibility with QIIME and older OS X versions.

v1.9.0 released November 12th, 2015  
not been Added the --fastq\_mergepairs command and associated options. This command has compilation tested well yet. Included additional files to avoid dependency of autoconf for tabs, compilation. Fixed an error where identifiers in fasta headers were not truncated at gzip compression just spaces. Fixed a bug in detection of the file format (FASTA/FASTQ) of a compressed input file.

v1.9.1 released November 13th, 2015  
Fixed memory leak and a bug in score computation in --fastq\_mergepairs, and

improved

speed.

v1.9.2 released November 17th, 2015

Fixed a bug in the computation of some values with --fastq\_stats.

v1.9.3 released November 19th, 2015

Workaround for missing x86intrin.h with old compilers.

v1.9.4 released December 3rd, 2015

Fixed incrementation of counter when relabeling dereplicated sequences.

v1.9.5 released December 3rd, 2015

Fixed bug resulting in inferior chimera detection performance.

v1.9.6 released January 8th, 2016

Fixed bug in aligned sequences produced with --fastapairs and --userout (qrow, trow)

options.

v1.9.7 released January 12th, 2016

Masking behaviour is changed somewhat to keep the letter case of the input sequences unchanged when no masking is performed. Masking is now performed also during chimera detection. Documentation updated.

v1.9.8 released January 22nd, 2016

Fixed bug causing segfault when chimera detection is performed on extremely short sequences.

v1.9.9 released January 22nd, 2016

Adjusted default minimum number of word matches during searches for improved performance.

v1.9.10 released January 25th, 2016

Fixed bug related to masking and lower case database sequences.

v1.10.0 released February 11th, 2016

Parallelized and improved merging of paired-end reads and adjusted some defaults. Re-option to commands. Fixed typos. Added relabelling to files produced with --consout and --profile options.

v1.10.1 released February 23rd, 2016

Fixed a bug affecting the --fastq\_mergepairs command causing FASTQ headers to be truncated at first space (despite the bug fix release 1.9.0 of November 12th, 2015). Full headers are now included in the output (no matter if --notrunclabels is in effect or not).

v1.10.2 released March 18th, 2016

empty Fixed a bug causing a segmentation fault when running --usearch\_global with an  
with an query sequence. Also fixed a bug causing imperfect alignments to be reported  
fasta/fastq alignment string of '=' in uc output files. Fixed typos in man file. Fixed  
files. processing code regarding presence or absence of compression library header

v1.11.1 released April 13th, 2016  
Added strand information in UC file for --derep\_fulllength and --derep\_prefix.  
Added --fas- expected errors (ee) to header of FASTA files specified with --fastaout and  
fastq\_filter taout\_discarded when --eeout or --fastq\_eeout option is in effect for  
and fastq\_mergepairs. The options --eeout and --fastq\_eeout are now equivalent.

v1.11.2 released June 21st, 2016  
used a Two bugs were fixed. The first issue was related to the --query\_cov option that  
defined as different coverage definition than the qcov userfield. The coverage is now  
or mis- the fraction of the whole query sequence length that is aligned with matching  
related to matching residues in the target. All gaps are ignored. The other issue was  
in some the consensus sequences produced during clustering when only N's were present  
behaviour positions. Previously these would be converted to A's in the consensus. The  
compat- is changed so that N's are produced in the consensus, and it should now be more  
ible with usearch.

v2.0.0 released June 24th, 2016  
added: This major new version supports reading from pipes. Two new options are  
specified if --gzip\_decompress and --bzip2\_decompress. One of these options must be  
ordinary reading compressed input from a pipe, but are not required when reading from  
written to files. The vsearch header that was previously written to stdout is now  
'-' now stderr. This enables piping of results for further processing. The file name  
reading or represent standard input (/dev/stdin) or standard output (/dev/stdout) when  
refac- writing files, respectively. Code for reading FASTA and FASTQ files has been  
tored.

v2.0.1 released June 30th, 2016  
Avoid segmentation fault when masking very long sequences.

v2.0.2 released July 5th, 2016  
Avoid warnings when compiling with GCC 6.

v2.0.3 released August 2nd, 2016  
Fixed bad compiler options resulting in Illegal instruction errors when running

precom- piled binaries.

v2.0.4 released September 1st, 2016  
Improved error message for bad FASTQ quality values. Improved manual.

v2.0.5 released September 9th, 2016  
Add options --fastaout\_discarded and --fastqout\_discarded to output discarded sequences from subsampling to separate files. Updated manual.

v2.1.0 released September 16th, 2016  
New command: --fastx\_filter. New options: --fastq\_maxlen, --fastq\_truncee. Allow --min- wordmatches down to 3.

v2.1.1 released September 23rd, 2016  
Fixed bugs in output to UC-files. Improved help text and manual.

v2.1.2 released September 28th, 2016  
Fixed incorrect abundance output from fastx\_filter and fastq\_filter when relabelling.

v2.2.0 released October 7th, 2016  
Added OTU table generation options --biomout, --mothur\_shared\_out and --otutabout to the clustering and searching commands.

v2.3.0 released October 10th, 2016  
Allowed zero-length sequences in FASTA and FASTQ files. Added --fastq\_trunclen\_keep option. Fixed bug with output of OTU tables to pipes.

v2.3.1 released November 16th, 2016  
Fixed bug where --minwordmatches 0 was interpreted as the default minimum word matches for the given word length instead of zero. When used in combination with --maxaccepts 0 and --maxrejects 0 it will allow complete bypass of kmer-based heuristics.

v2.3.2 released November 18th, 2016  
Fixed bug where vsearch reported the ordinal number of the target sequence instead of the cluster number in column 2 on H-lines in the uc output file after clustering. For search and alignment commands both usearch and vsearch reports the target sequence number here.

v2.3.3 released December 5th, 2016  
A minor speed improvement.

v2.3.4 released December 9th, 2016  
Fixed bug in output of sequence profiles and updated documentation.

v2.4.0 released February 8th, 2017  
Added support for Linux on Power8 systems (ppc64le) and Windows on x86\_64. Improved detection of pipes when reading FASTA and FASTQ files. Corrected option for specifying

output from fastq\_eestats command in help text.

v2.4.1 released March 1st, 2017

Fixed an overflow bug in fastq\_stats and fastq\_eestats affecting analysis of very large FASTQ files. Fixed maximum memory usage reporting on Windows.

v2.4.2 released March 10th, 2017

Default value for fastq\_minovlen increased to 16 in accordance with help text and for compatibility with usearch. Minor changes for improved accuracy of paired-end read merging.

v2.4.3 released April 6th, 2017

Fixed bug with progress bar for shuffling. Fixed missing N-lines in UC files with use-arch\_global, search\_exact and allpairs\_global when the output\_no\_hits option was not specified.

v2.4.4 released August 28th, 2017

Fixed a few minor bugs, improved error messages and updated documentation.

v2.5.0 released October 5th, 2017

Support for UDB database files. New commands: fastq\_stripright, fastq\_eestats2, makeudb\_usearch, udb2fasta, udbinfo, and udbstats. New general option: no\_progress. New options minsize and maxsize to fastx\_filter. Minor bug fixes, error message improvements and documentation updates.

v2.5.1 released October 25th, 2017

Fixed bug with bad default value of 1 instead of 32 for minseqlength when using the makeudb\_usearch command.

v2.5.2 released October 30th, 2017

Fixed bug where '-' as an argument to the fastq\_eestats2 option was treated literally instead of equivalent to stdin.

v2.6.0 released November 10th, 2017

Rewritten paired-end reads merger with improved accuracy. Decreased default value for fastq\_minovlen option from 16 to 10. The default value for the fastq\_maxdiffs option is increased from 5 to 10. There are now other more important restrictions that will avoid merging reads that cannot be reliably aligned.

v2.6.1 released December 8th, 2017

Improved parallelisation of paired end reads merging.

v2.6.2 released December 18th, 2017

Fixed option xsize that was partially inactive for commands uchime\_denovo, uchime\_ref, and fastx\_filter.

v2.7.0 released February 13th, 2018  
Added commands cluster\_unoise, uchime2\_denovo and uchime3\_denovo contributed by Davide Albanese based on Robert Edgar's papers. Refactored fasta and fastq print functions as well as code for extraction of abundance and other attributes from the headers.

v2.7.1 released February 16th, 2018  
Fix several bugs on Windows related to large files, use of "-" as a file name to mean stdin or stdout, alignment errors, missed kmers and corrupted UDB files. Added documentation of UDB-related commands.

v2.7.2 released April 20th, 2018  
Added the syntax command for taxonomic classification. Fixed a bug with incorrect FASTA headers of consensus sequences after clustering.

v2.8.0 released April 24th, 2018  
Added the fastq\_maxdiffpct option to the fastq\_mergepairs command.

v2.8.1 released June 22nd, 2018  
Fixes for compilation warnings with GCC 8.

v2.8.2 released August 21st, 2018  
Fix for wrong placement of semicolons in header lines in some cases when using the derePLICATION in sizeout or xsize options. Reduced memory requirements for full-length report. Updated manual regarding use of sizein and sizeout with derePLICATION. Changed a compiler option.

v2.8.3 released August 31st, 2018  
Fix for segmentation fault for --derep\_fulllength with --uc.

v2.8.4 released September 3rd, 2018  
Further reduce memory requirements for derePLICATION when not using the uc option. Fix output during subsampling when quiet or log options are in effect.

v2.8.5 released September 26th, 2018  
Fixed a bug in fastq\_eestats2 that caused the values for large lengths to be much too high when the input sequences had varying lengths.

v2.8.6 released October 9th, 2018  
Fixed a bug introduced in version 2.8.2 that caused derep\_fulllength to include the full FASTA header in its output instead of stopping at the first space (unless the notrunclabels option is in effect).

v2.9.0 released October 10th, 2018  
Added the fastq\_join command.

v2.9.1 released October 29th, 2018

software to illegal instruction error on some architectures. Update documentation of rereplicate command.

v2.10.0 released December 6th, 2018  
Added the sff\_convert command to convert SFF files to FASTQ. Added some additional option argument checks. Fixed segmentation fault bug after some fatal errors when a log file was specified.

v2.10.1 released December 7th, 2018  
Improved sff\_convert command. It will now read several variants of the SFF format. It is also able to read from a pipe. Warnings are given if there are minor problems. Error messages have been improved. Minor speed and memory usage improvements.

v2.10.2 released December 10th, 2018  
Fixed bug in syntax with reversed order of domain and kingdom.

v2.10.3 released December 19th, 2018  
8.1.0 Ported to Linux on ARMv8 (aarch64). Fixed compilation warning with gcc version 8.2.0.

v2.10.4 released January 4th, 2019  
Added link to BioConda in README. Fixed bug in fastq\_stats with sequence length 1. Fixed use of equals symbol in UC files for identical sequences with cluster\_fast.

v2.11.0 released February 13th, 2019  
with the attributes from Added ability to trim and filter paired-end reads using the reverse option fastx\_filter and fastq\_filter commands. Added --xee option to remove ee FASTA headers. Minor invisible improvement to the progress indicator.

v2.11.1 released February 28th, 2019  
cluster\_unico Minor change to the handling of the weak\_id and id options when using cluster\_unico.

v2.12.0 released March 19th, 2019  
profiles after offset option Take sequence abundance into account when computing consensus sequences or ter clustering. Warn when rereplicating sequences without abundance info. Guess 33 in more cases with fastq\_chars. Stricter checking of option arguments and combinations.

v2.13.0 released April 11th, 2019  
extract se- Added the --fastx\_getseq, --fastx\_getseqs and --fastx\_getsubseq commands to quences from a FASTA or FASTQ file based on their labels. Improved handling of

ambiguous nucleotide symbols. Corrected behaviour of --uchime\_ref command with and options --self and --selfid. Strict detection of illegal options for each command.

v2.13.1 released April 26th, 2019  
the log, Minor changes to the allowed options for each command. All commands now allow quiet and threads options. If more than 1 thread is specified for commands that are not multi-threaded, a warning will be issued. Minor changes to the manual.

v2.13.2 released April 30th, 2019  
strand Fixed bug related to improper handling of newlines on Windows. Allowed option plus to uchime\_ref for compatibility.

v2.13.3 released April 30th, 2019  
the out- Fixed bug in FASTQ parsing introduced in version 2.13.2.

v2.13.4 released May 10th, 2019  
NetBSD Added information about support for gzip- and bzip2-compressed input files to output of the version command. Adapted source code for compilation on FreeBSD and systems.

v2.13.5 released July 2nd, 2019  
from the Added cut command to fragment sequences at restriction sites. Silenced output fastq\_stats command if quiet option was given. Updated manual.

v2.13.6 released July 2nd, 2019  
Added info about cut command to output of help command.

v2.13.7 released September 2nd, 2019  
Fixed bug in consensus sequence introduced in version 2.13.0.

v2.14.0 released September 11th, 2019  
options Added relabel\_self option. Made fasta\_width, sizein, sizeout and relabelling valid for certain commands.

v2.14.1 released September 18th, 2019  
commands Fixed bug with sequences written to file specified with fastaout\_rev for fastx\_filter and fastq\_filter.

v2.14.2 released January 28th, 2020  
fastq\_mergelpairs, and Fixed some issues with the cut, fastx\_revcomp, fastq\_convert, makeudb\_usearch commands. Updated manual.

v2.15.0 released June 19th, 2020  
by de- Update manual and documentation. Turn on notrunclabels option for syntax command char- fault. Change maxhits 0 to mean unlimited hits, like the default. Allow non-ascii spec- characters in headers, with a warning. Sort centroids and uc too when clusterout\_sort

ified. Add cluster id to centroids output when clusterout\_id specified. Improve error messages when parsing FASTQ files. Add missing fastq\_qminout option and fix fix option for fastq\_mergepairs. Add derep\_id command that dereplicates based on both label and sequence. Remove compilation warnings.

v2.15.1 released October 28th, 2020  
Fix for dereplication when including reverse complement sequences and headers. Make some extra checks when loading compression libraries and add more diagnostic output about them to the output of the version command. Report an error when fastx\_filter is used with FASTA input and options that require FASTQ input. Update manual.

v2.15.2 released January 26th, 2021  
No real functional changes, but some code and compilation changes. Compiles fully on macOS running on Apple Silicon (ARMv8). Binaries available. Code C++11. Minor adaptations for Windows compatibility, including the use of the C++ standard library for regular expressions. Minor changes for compatibility with Power8.  
Switch to C++ header files.

v2.16.0 released March 22nd, 2021  
This version adds the orient command. It also handles empty input files properly. Documentation has been updated.

v2.17.0 released March 29nd, 2021  
The fastq\_mergepairs command has been changed. It now allows merging of sequences with overlaps as short as 5 bp if the --fastq\_minovlen option has been adjusted down from the default 10. In addition, much fewer pairs of reads should now be rejected with the reason 'multiple potential alignments' as the algorithm for detecting those have been changed.

v2.17.1 released June 14th, 2021  
Modernized code. Minor changes to help info.

v2.18.0 released August 27th, 2021  
Added the fasta2fastq command. Fixed search bug on ppc64le. Fixed bug with removal of size and ee info in uc files. Fixed compilation errors in some cases. Made some general code improvements. Updated manual.

v2.19.0 released December 21st, 2021  
Added the lcaout and lca\_cutoff options to enable the output of last common ancestor valid (LCA) information about hits when searching. The randseed option was added as a option to the syntax command. Code improvements.

v2.20.0 released January 10th, 2022  
Added the fastx\_uniques command and the fastq\_qout\_max option for dereplication of FASTQ files. Some code cleaning.

v2.20.1 released January 11th, 2022  
Fixes a bug in fastq\_mergelpair that caused an occasional hang at the end when using multiple threads.

v2.21.0 released January 12th, 2022  
This version adds the sample, qsegout and tsegout options. It enables the use of UDB databases with uchime\_ref.

v2.21.1 released January 18th, 2022  
Fix a problem with dereplication of empty input files. Update Altivec code on ppc64le for improved compiler compatibility (vector->\_\_vector).

v2.21.2 released September 12th, 2022  
Fix problems with the lcaout option when using maxaccepts above 1 and either lca\_cutoff below 1 or with top\_hits\_only enabled. Update documentation. Update code to avoid compiler warnings.

v2.22.0 released September 19th, 2022  
Add the derep\_smallmem command for dereplication using little memory.

v2.22.1 released September 19th, 2022  
Fix compiler warning.

v2.23.0 released July 7th, 2023  
Update documentation. Add citation file. Modernize and improve code. Fix several minor bugs. Fix compilation with GCC 13. Print stats after fastq\_mergelpairs to log instead of stderr. Handle sizein option correctly with dbmatched option for use-arch\_global. Allow maxseqlength option for makefdb\_usearch. Fix memory allocation prob-lem with chimera detection. Add lengthout and xlength options. Increase precision for eeout option. Add warning about syntax algorithm, random seed and multiple threads.

minor file in-for use-prob-precision for long\_chimeras\_denovo  
Refactor chimera detection code. Add undocumented experimental command. Fix segfault with clustering. Add more references.

v2.24.0 released October 26th, 2023  
Update documentation. Improve code. Allow up to 20 parents for the undocumented and experimental chimeras\_denovo command. Fix compilation warnings for shal.c. Compile for release (not debug) by default.

v2.25.0 released November 10th, 2023  
Allow a given percentage of mismatches between chimeras and parents for the

experimental chimeras\_denovo command.

v2.26.0 released November 24th, 2023  
Enable the maxseqlength and minseqlength options for the chimera detection commands.

include samples and OTUs with no matches.

v2.26.1 released November 25th, 2023  
No real changes, but the previous version was released without proper updates to the source code.

v2.27.0 released January 19th, 2024  
The usearch\_global and search\_exact commands now support FASTQ files as well as FASTA files as input. This version of vsearch includes clarifications and updates to the manual. Some code has been refactored. Generic Dockerfiles for major Linux distributions have been included. Some warnings from compilers and other tools have been eliminated. The release for Windows will also include DLL's for the two compression libraries.

v2.27.1 released April 6th, 2024  
This version fixes the weak\_id option and makes searches report weak hits in some cases. It also updates the names of the compression libraries to libz.so.1 and libbz2.so.1 on Linux to make them work on common Linux distributions without installing additional packages. README.md has been updated with information about compression libraries on Windows.

v2.28.0 released April 26th, 2024  
The syntax command has been improved in several ways in this version of vsearch. Please note that several details of this algorithm is not clearly described in the preprint, and the implementation in vsearch differs from that in usearch. The former vsearch version did not always choose the most common taxonomic entity over the 100 bootstraps among the database sequences with the highest amount of word similarity to the query. Instead, if several sequences had an equal similarity with the query, the sequence entity calculated bootstraps. This could lead to a suboptimal choice with a low confidence. In the new version, the most common of the sequences with the highest amount of word similarity across the 100 bootstraps will be selected, and ties will be broken randomly. Another problem

with the  
similar-  
they were  
called  
the se-  
considering  
sequences. This  
Furthermore, a  
syntax  
vsearch with  
in the  
old implementation was that if several sequences had the same amount of word  
length, the shortest one in the reference database would be chosen, and if  
equally long, the earliest in the database file would be chosen. A new option  
syntax\_random has now been introduced. This option will randomly select one of  
sequences with the highest number of shared words with the query, without  
their length or position. This avoids a bias towards shorter reference  
option is strongly recommended and will probably soon be the default.  
ninth taxonomic rank, strain (letter t), is now recognized. The speed of the  
command has also been significantly improved at least in some cases. Run  
the randseed option and 1 thread to ensure reproducibility of the random choices  
algorithm.

v2.28.1 released April 26th, 2024  
Fix a segmentation fault that could occur with the blast6out and  
output\_no\_hits options.

v2.29.0 released September 26th, 2024  
This version fixes seven bugs (see changelog below), adds initial support for  
RISC-V architectures, and improves code quality and code testing (1,210 new tests):  
- add: experimental support for RISCV64 and other 64-bit little-endian  
architectures, thanks to Michael R. Crusoe and his fellow Debian developers (issue #566),  
- add: official support for clang-19 and gcc 14,  
- add: beta support for clang-20,  
- remove: unused --output option for command --fastq\_stats (issue #572),  
- fix: bug in --syntax when selecting the best lineage (only low confidence  
values below 0.5 were affected) (issue #573),  
- fix: out-of-bounds error in --fastq\_stats when processing empty reads (issue  
#571),  
- fix: bug in --cut, patterns with multiple cutting sites were not detected  
(commit 4c4f9fa70f14b28d50185dbf322cf5727087e86a),  
- fix: memory error (segmentation fault) when using --derep\_id and --strand  
(issue #565),  
- fix: --fastq\_join now obeys to --quiet and --log options  
(commit

87f968b09f17c17ebf8db00aebe86e89b13a3948),  
is 64  
- fix: --fastq\_join quality padding is now also set to Q40 when quality offset  
(commit  
(commit  
that they  
(commit  
(commit  
- improve: additional safeguards to validate input values and to make sure  
are within acceptable limits. Changes concern options --abskew  
a530dd8990f8a05cb25fc0b6a5da5a14d28fbbedd) and --fastq\_maxdiffs  
4b254db7f120bfd49e86185ef3cd9070c236f940),  
- improve: code quality (1.3k+ commits, 6k+ clang-tidy warnings eliminated),  
- improve: documentation and help messages (issue #568),  
- improve: complete refactoring and modernization of a subset of commands  
(--sort-bylength, --sortbysize, --shuffle, --rereplicate, --cut, --fastq\_join, --  
fastq\_chars),  
(1,210 new  
tests, 4,753 in total)

v2.29.1 released October 24th, 2024  
exam-  
Fix a segmentation fault that could occur during alignment in version 2.29.0, for  
ple with --uchime\_ref. Some improvements to code and documentation.

v2.29.2 released December 20th, 2024  
Initial  
Fix a segmentation fault during clustering when the set of clusters is empty.  
documentation in markdown format available on GitHub Pages.

v2.29.3 released February 3rd, 2025  
compiling the  
This version is released in order to mitigate a bug that occurs when  
later  
cause bad  
for now  
Dockerfiles have  
include this  
`align\_simd.cc` file on x86\_64 systems with the GNU C++ compiler version 9 or  
with the `-O3` optimization option. It results in incorrect code that may  
alignments in some circumstances. We are investigating this issue further, but  
we recommend compiling with the `-O2` flag. The README.md file and the  
been updated to reflect this. The binaries released with this version will  
fix.

v2.29.4 released February 14th, 2025  
window size  
Adjust the window size used for chimera detection down from 64 to 32. The

fewer was by accident increased from 32 to 64 in version 2.23.0, leading to somewhat chimeras being predicted. In addition, a compiler pragma has been included in align\_simd.cc to further protect the compiler from generating wrong code.

v2.30.0 released February 27th, 2025  
Add options `--n\_mismatch`, `--fastq\_minqual`, and `--fastq\_truncee\_rate`. The  
`--n\_mismatch` option will count N's as mismatches in alignments, which may be useful  
to get sensible alignments for sequences with lots of N's. By default N's are counted  
as matches. Both the scoring and the counting of matches are affected.  
The new `--fastq\_minqual` option for the `fastq\_filter` and `fastx\_filter` commands  
will discard sequences with any bases with a quality scores below the given value.  
The new `--fastq\_truncee\_rate` option for the same commands will truncate sequences  
at the first position where the number of expected errors per base is above the given value.

v2.30.1 released October 3rd, 2025  
This version incorporates many code improvements, more extensive testing, better documentation and some minor bug fixes. List of changes:

- fix: use-after-free introduced in commit de6c1d8 (Jun 13, 2024),  
8a0a508b),
  - fix: (harmless) out-of-bounds memory issue in --derep\_prefix (commit a9c42713),
    - fix: (harmless) memory leak in --fastx\_getseqs --label\_field (commit a9c42713),
      - fix: (harmless) memory leak when using option --userfields (--allpairs\_global 03b95bcf; --cluster\_\* commit 2fde5472; --search\_exact commit 45cd56d6; arch\_global commit d83bfee9),
        - fix: (harmless) valgrind error, use of uninitialized values (commit 8bab2444), also eliminates a pesky compilation warning,
- change: passing a negative value to --fastq\_truncee\_rate is now an error (commit a120f371),
  - change: passing a negative value to --fastq\_minqual is now an error (commit ff5b0c99),
    - change: passing a non-ASCII symbol to --join\_padgap or --join\_padgapq is now an error (commit a708f5b3),
      - change: when using --gapopen "\*" to forbid gap opening, the penalty is now set to

relied INT\_MAX (rather than 1,000). This might change alignment results for users who  
on the old behavior (thanks to Denis Filloux, issue #602, commit 96e9cf9e),  
the only - change: when using command --chimeras\_denovo, --tabbedout or --alnout can be  
output files specified (commit d51f0a456300a0ea69c035ba3de23c5ecf3da348)  
accepted - change: when using command --chimeras\_denovo, option --lengthout is now  
(com- (commit 3f55cc6b))  
5cf4a6c1, - change: when using command --chimeras\_denovo, option --xlength is now accepted  
option (commit 0fd346cd)  
headers - add: compilation option GLIBCXX\_DEBUG when compiling for debugging (commit  
option activates more runtime checks),  
head- - add: official support for clang 20,  
ers (commit 98a851ed),  
(commit - add: initial support for clang 21, initial support for GCC 15,  
- add: experimental support for clang 22,  
- improve: more accurate line number when reporting illegal characters in fastq  
(commit 539084e9),  
- improve: more accurate line number when reporting non-ASCII characters in fastq  
ers (commit 98a851ed),  
- improve: remove checks for unneeded libraries during compilation  
249bb5d5276b1be6c9add60a538a16e1b9d0ebfb),  
- improve: code quality (8,536 clang-tidy warnings eliminated),  
- improve: documentation and help messages (issues #604),  
- improve: complete refactoring and modernization of the command --fastq\_stats,  
- improve: command --fastq\_stats is now up to twice faster (tested on x86-64),  
- improve: extensive test-suites for --fastq\_stats and --sff\_convert,  
- improve: code coverage of our test-suite