

VSEARCH and Swarm: Tools for analysis of microbiome sequencing data

BI09905MERG1 Course, UiO, 3 May 2021

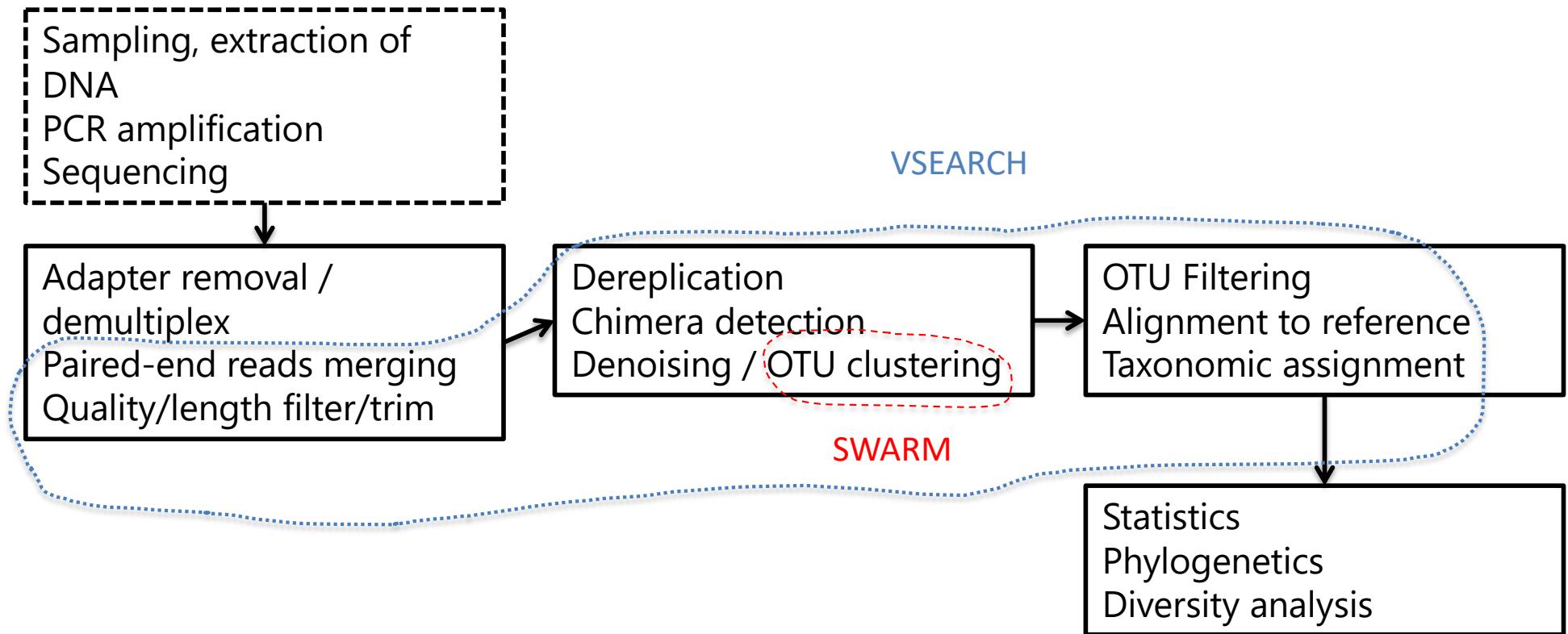
Torbjørn Rognes
Dept. of Informatics, UiO & Dept. of Microbiology, OUS
torognes@ifi.uio.no



UiO • University of Oslo



Amplicon sequence analysis pipeline



VSEARCH: a versatile analysis tool

- Versatile command-line tool to analyse DNA sequence data in microbiome projects
- Free of charge
- Open source software
- Robust, well-tested
- Available on many platforms and architectures
- Able to handle very large databases (>4GB)
- A drop-in replacement for USEARCH in many cases
- In general, equal or better accuracy and speed as USEARCH

Motivation: Problems with USEARCH

Frustration in the metagenomics community with the **USEARCH** tool (Edgar, 2010) often used for many purposes in metagenomics projects, including the QIIME pipeline.

Main problems:

- Inner workings only rudimentary described: we don't really know how it works
- Source code is not publicly available: we cannot improve on it, adapt it, or fully understand it
- Only a 32-bit version freely available for academic use: it cannot be used with large datasets
- 64-bit version requires payment, per process (USD 885 for academic / USD 1485 commercial)

Features of VSEARCH v2.17.0

- ~45 commands with ~180 command line options:
 - Chimera detection (*de novo* (uchime, uchime2 or uchime3), and reference)
 - Clustering (after sorting by abundance or length, or using initial order, or similar to unoise)
 - Dereplication of sequences (full length, prefix & id) and rereplication
 - Extraction of sequences and sub-sequences from large FASTA files
 - FASTQ encoding detection and conversion, SFF to FASTQ conversion
 - Masking of low-complexity regions with Dust
 - Orienting of sequences in same direction as database sequences
 - Paired-end reads merging and joining
 - Pairwise global sequence alignment (all vs all)
 - Restriction site cutting
 - Reverse complementation of sequences
 - Searching (global alignment and exact matches)
 - Sequence quality statistics and filtering
 - Shuffling and sorting (abundance and length)
 - Subsampling of sequences
 - Taxonomic classification (SINTAX)
 - UDB files: Building, extraction and statistics of fast-loading pre-indexed database files
- Automatic detection and reading of compressed files (.gz, .bz2)



Getting help with VSEARCH

Commands:

vsearch

vsearch --help | less

man vsearch

Manual also available as PDF:
vsearch_manual.pdf

```
vsearch(1)                               USER COMMANDS                               vsearch(1)
NAME
    vsearch — chimera detection, clustering, dereplication, masking, pairwise alignment, searching, shuffling
    and sorting of amplicons from metagenomic projects.

SYNOPSIS
    Chimera detection:
        vsearch --uchime_denovo fastaf file (-chimeras | -nonchimeras | -uchimealns | -uchimeout) out-
        putfile [options]
        vsearch --uchime_ref fastaf file (-chimeras | -nonchimeras | -uchimealns | -uchimeout) outputfile
        -db fastaf file [options]

    Clustering:
        vsearch (-cluster_fast | -cluster_size | -cluster_smallmem) fastaf file (-alnout | -blast6out | -cen-
        troids | -clusters | -missout | -uc | -userout) outputfile -id real [options]

    Dereplication:
        vsearch --derep_fulllength fastaf file (-output | -uc) outputfile [options]

    Masking:
        vsearch --maskfasta fastaf file -output outputfile [options]

    Pairwise alignment:
        vsearch --allpairs_global fastaf file (-alnout | -blast6out | -matched | -notmatched | -uc |
        -userout) outputfile (-acceptall | -id real) [options]

    Searching:
        vsearch --usearch_global fastaf file -db fastaf file (-alnout | -blast6out | -uc | -userout) outputfile
        -id real [options]

    Shuffling:
        vsearch --shuffle fastaf file -output outputfile [options]

    Sorting:
        vsearch (-sortbylength | -sortbysize) fastaf file -output outputfile [options]

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

    Nucleotide sequence comparisons is at the core of vsearch. To speed up comparisons, vsearch implements
    an extremely fast implementation of the Needleman-Wunsch algorithm, making use of the Streaming SIMD
    Extensions (SSE2) of modern x86-64 CPUs. If SSE2 instructions are not available, vsearch exits with an
    error message. For comparisons involving sequences longer than 5,000 nucleotides, vsearch uses a slower
    alignment method with smaller memory requirements.

Input
    vsearch input is a fasta file containing one or several nucleotide sequences. For each sequence, the
    sequence identifier is defined as the string comprised between the '>*' symbol and the first space, or the end
    of the line, whichever comes first. Additionally, if the line starts with '>[j]size=integer;label', contains
    '>label;size=integer;label' or ends with '>label;size=integer;', vsearch will remove the pattern
    [j]size=integer[] from the header and 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, dereplica-
    tion, sorting and searching.

    The nucleotide 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 non-blocking

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

Nucleotide sequence comparisons is at the core of vsearch. To speed up comparisons, vsearch implements
an extremely fast implementation of the Needleman-Wunsch algorithm, making use of the Streaming SIMD
Extensions (SSE2) of modern x86-64 CPUs. If SSE2 instructions are not available, vsearch exits with an
error message. For comparisons involving sequences longer than 5,000 nucleotides, vsearch uses a slower
alignment method with smaller memory requirements.

Input
    vsearch input is a fasta file containing one or several nucleotide sequences. For each sequence, the
    sequence identifier is defined as the string comprised between the '>*' symbol and the first space, or the end
    of the line, whichever comes first. Additionally, if the line starts with '>[j]size=integer;label', contains
    '>label;size=integer;label' or ends with '>label;size=integer;', vsearch will remove the pattern
    [j]size=integer[] from the header and 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, dereplica-
    tion, sorting and searching.

    The nucleotide 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 non-blocking

version 1.0.10                               January 23, 2015                               1
```

Merging paired-end reads

```
vsearch --fastq_mergepairs forward.fastq \
    --reverse reverse.fastq \
    --fastqout merged.fastq \
    --minovlen 10 \
    --fastq_allowmergestagger
```

Merges overlapping forward and reverse reads into one sequence, if possible. Unmerged reads may be written to separate files.

Merging paired-end reads



Read 1 (green): Forward sequence

Read 2 (red): Reverse complementary sequence

- Find best overlap between the two sequences
- Take the quality score (error probability) of each base into account

Filtering reads

```
vsearch --fastq_filter input.fastq \
--fastq_maxee 1.0 \
--fastq_maxns 0 \
--fastq_minlen 100 \
--fastqout filtered.fastq \
--fastaout filtered.fasta \
--relabel abc
```

Removes or truncates reads that does not satisfy given requirements.

Dereplication

```
vsearch --derep_fulllength input.fasta \
--output derep.fasta \
--minuniquesize 2 \
--sizeout
```

Groups identical sequences into one entry in the FASTA file.
Adds a "size" attribute to the FASTA header with the
abundance (number of identical copies) of the sequence.

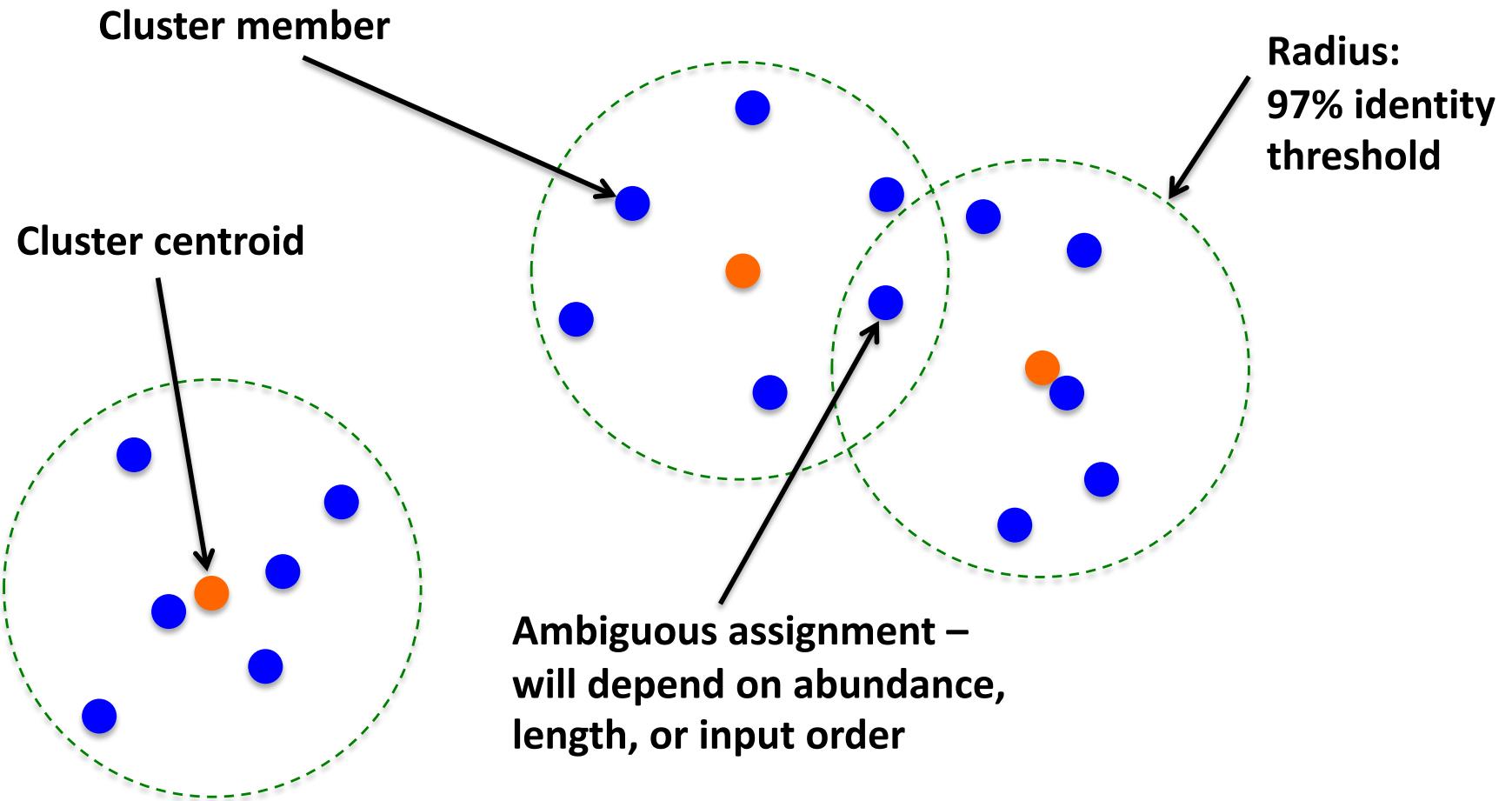
```
>abc;size=123
acgttagtcagactgtcagactgtg
```

Sequence similarity clustering

```
vsearch --cluster_size input.fasta \
--id 0.97 \
--centroids centroids.fasta \
--sizein \
--sizeout
```

Groups similar sequences into clusters (OTUs) using a heuristic and greedy centroid-based algorithm.

Heuristic centroid-based clustering



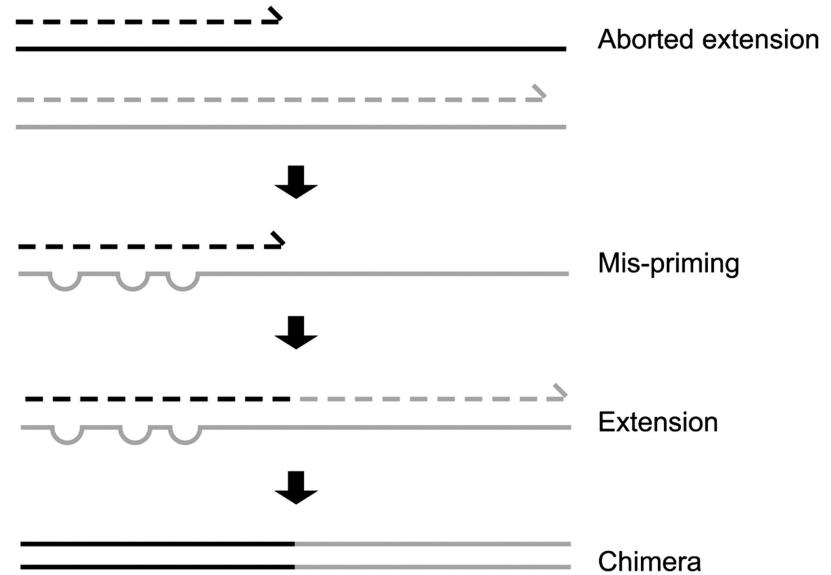
Chimera detection

```
vsearch --uchime_denovo input.fasta \
--nonchimeras nonchimeras.fasta \
--chimeras chimeras.fasta \
--sizein \
--sizeout
```

Detects potential chimeric sequences in the input file and writes the classified sequences into separate files.

Chimeras

- Chimeric DNA sequences may form during PCR amplification
- Chimeras contain segments from 2 or more parent sequences
- Will inflate the apparent diversity of organisms in the sample unless removed
- VSEARCH implements the algorithms UCHIME (Edgar et al., 2011), UCHIME2 and UCHIME3
- Use either the dataset itself (`uchime_denovo`) or a reference database (`uchime_ref`) during analysis



Chimera detection example

Query (250 nt) ch31_m2_90_95/sp8:0-149/sp62:149-249/N=2/top=sp8:92.4%

ParentA (250 nt) sp8/name=Clostridiummethylpentosum_0_1_2

ParentB (250 nt) sp62/name=Clostridiumsporogenes_0_3_2

A 1 TGCTGCCTCCGTAGGAGTCTGGCCGTGtctcagtcCCAATGTGCCGT-TAACCTCTCAGTCGGCTA-CTGATCGt 78

Q 1 TGCTGCCTCCGTAGGAGTCTGGCCGTGTTAGTCGCCAATGTGCCGTCCAACCTCTCAGTCGGCTAGCTGATCG- 79

B 1 TGCTGCCTCCGTAGGAGTCTGGaCCGTGTctcagttCCAATGTGCCGat-CAcCCTCTCAGgtCGGCTAcgcacatcg- 78

Diffs A NNNNNN? A A AA AAAAAAA

Votes + 0000000 + + ++ ++++

Model AAA

A 79 CGACTTGGTGAGCCATTACCTCACCAACTATctaATCAGA-CGCGAGCCCATTaCAGCGATATAATCTTGAT-AAcA 156

Q 80 CGACTTGGTGAGCCATTACCTCACCAACTAT-TAATCAGACCGCGAGCCCATT-CAGCGATATAATCTTGATAAAAAA 157

B 79 tGcCTTGGTaAGCCgTTACCTtACCAACTAg-ctAatgcgCCGCGgGtCCATCTc-aAagcAataAACTTTGAT-AAAA 155

Diffs A A A A A AA AAA A A AAA AAA B

Votes + + + + + + + + + + + + + + + +

Model AAxxxxxxxxxxBB

A 157 AAACATGCGATTCCgTTATgTTATGCGGTATTAgcgTTCgTTTCC--AAacGtTATTCCCctcTgtAAGGCAGGTTgCt 234

Q 158 AAATCATGCGATTCTCTTATATTGCGGTATTAAATCTCCTTCG--AA--GCTATCCCCACTTGAAGGCAGGTTACC 233

B 156 AAATCATGCGATTCTCTTATATTGCGGTATTAAATCTCCTTCGgaAg--GCTATCCCCcacTTtgAGGCAGGTTACC 233

Diffs B BB B BBB B B A B B N?N BNa B B

Votes + ++ + +++ + + + + 000 +0! + +

Model BBBBxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx

A 235 CACGTGTTACTCACCC- 250

Q 234 CACGTGTTACTCACCCG 250

B 234 CACGTGTTACTCACCCG 250

Diffs

Votes

Model BBBBxxxxxxxxxxxxxxxxxxxx

Ids. QA 88.9%, QB 82.7%, AB 80.5%, QModel 94.7%, Div. +6.5%

Diffs Left 34: N 0, A 7, Y 27 (79.4%); Right 19: N 1, A 4, Y 14 (73.7%), Score 0.8952

Sequence similarity search

```
vsearch --usearch_global query.fasta \
--db database.fasta \
--id 0.99 \
--maxaccepts 3 \
--maxrejects 1000 \
--alnout alignments.txt \
--otutabout otutable.txt
```

Generic heuristic sequence similarity search using global alignment with many adjustable parameters.

Taxonomic classification (SINTAX)

```
vsearch --sintax otus.fasta \
--db silva.fasta \
--sintax_cutoff 0.8 \
--tabbedout tax.txt
```

Rapid classification of sequences using a taxonomy given in a specially formatted database. Uses the SINTAX algorithm.
Results may vary slightly due to randomness in the algorithm.

```
>AY232296.1.1383;tax=k:Archaea,p:Euryarchaeota,c:Halobacteria,o:Halobacteriales,f:Halobacteriaceae,g:Natrinema,s:Natrinema_sp._HM06;
```

General heuristic search algorithm

Used during search, chimera detection, and clustering

For each query sequence:

- Sort target (database) sequences by decreasing number of k -mers shared with the query sequence
- Consider target sequences in order and align the query to each candidate target sequence
- Stop when A targets (*maxaccepts*, default 1) have been accepted, i.e. they satisfy the accept criteria (e.g. id > 97%)
- Stop when R targets (*maxrejects*, default 32) have been rejected, i.e. they do not satisfy the accept criteria.

Sequence database indexing

- The sequence database is indexed
- Which database sequences contain at least one occurrence of a given k-mer?
- 4^k possible k-mers
- Quickly identify which database sequences share the most k-mers with the query
- By default $k=8$ (*wordlength*)
- By default at least 12 shared k-mers are required (*minwordmatches*), but never more than those that exist in the query
- Repeated k-mers counts only once

Example with $k=2$

Query:

AACGTCCTGCATCGATC

k-mers:

AA AC CG GT TC CC CT TG GC CA AT TC CG GA AT TC

Sorted:

AA AC AT AT CA CC CG CG CT GA GC GT TC TC TC TG

Remove duplicates:

AA AC AT -- CA CC CG -- CT GA GC GT TC -- -- TG

Target:

AACGGCCTGCATAAGATC

k-mers:

AA AC CG GG GC CC CT TG GC CA AT TA AG GA AT TC

Sorted:

AA AC AG AT AT CA CC CG CT GA GC GC GG TA TC TG

Remove duplicates:

AA AC AG AT -- CA CC CG CT GA GC -- GG TA TC TG

Shared k-mers:

AA AC AT CA CC CG CT GA GC TC TG

11 shared k -mers

Alignment differences

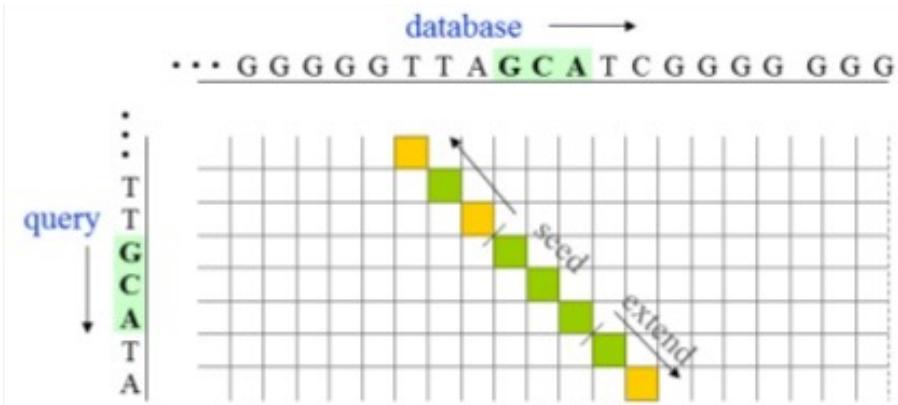
VSEARCH:

- optimal global alignment algorithm with backtracking (Needleman-Wunsch + Gotoh), 8-way SIMD parallelised

	_	A	G	A	C	T	A	G	T	T	A	C
_	0	-5	-10	-15	-20	-25	-30	-35	-40	-45	-50	-55
C	-5	-3	-8	-13	-6	-11	-16	-21	-26	-31	-36	-41
G	-10	-6	4	-1	-6	-9	-12	-9	-14	-19	-24	-29
A	-15	0	-1	14	9	4	1	-4	-9	-14	-9	-14
G	-20	-5	7	9	9	6	3	8	3	-2	-7	-12
A	-25	-10	2	17	12	7	16	11	6	1	8	3
C	-30	-15	-3	12	20	21	16	11	11	6	3	17
G	-35	-20	-8	7	21	23	20	25	20	15	10	12
T	-40	-25	-13	2	16	29	24	20	33	28	23	18

USEARCH:

- Seed-and-extend heuristic gapped alignment, as in BLAST



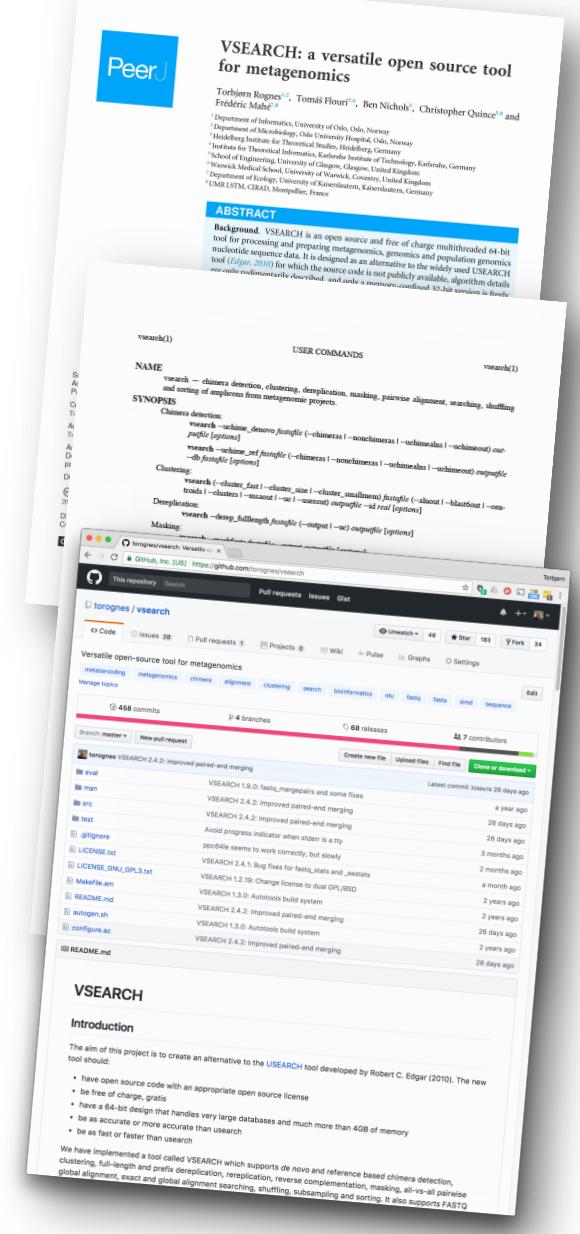
Some disadvantages of VSEARCH

Compared to USEARCH:

- Slower on long sequences due to full optimal alignment
- no support for amino acid sequences
- no UPARSE pipeline support
- newer USEARCH commands not supported

VSEARCH availability & documentation

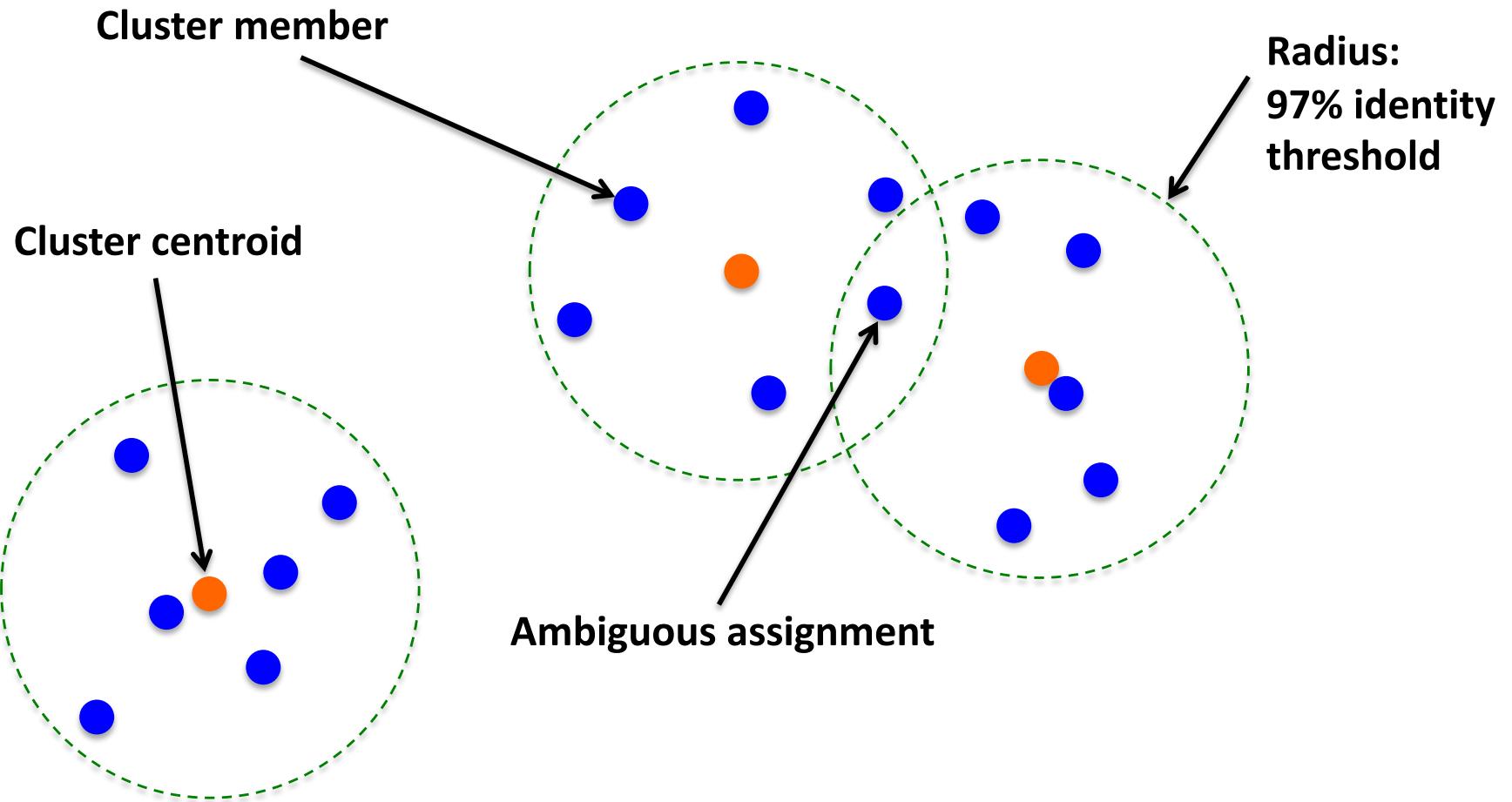
- VSEARCH source code: <https://github.com/torognes/vsearch>
- 64-bit binary binaries for Linux (x86, ARM, PPC), macOS (x86, ARM) and Windows available
- Dual open-source license: GNU AGPL v3 or BSD 2-clause
- May be used directly for clustering and chimera detection in mothur
- QIIME 2 plugin; Conda, Homebrew, Debian packages; Galaxy wrapper
- Publication:
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](https://doi.org/10.7717/peerj.2584) (>3000 citations)
- User manual with command and option details (over 40 pages)
- Extensively tested with >2000 unit tests
- Wiki, issue tracker, online support forum etc
- Suggest new features! Report bugs! Ask for help!



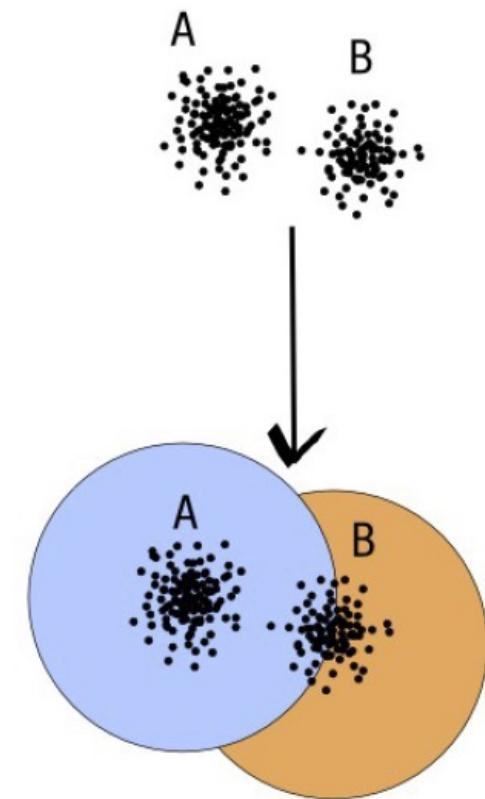
Swarm: a new clustering method

- New method for clustering amplicon sequences into Operational Taxonomic Units (OTUs)
- OTUs are groups of sequences that could approximately represent e.g. a species or a genus as we define it
- Single linkage hierarchical clustering approach
- Avoids two problems with traditional methods:
 - input order dependence
 - global clustering threshold
- High resolution
- Very fast: Linear time and space complexity: fast and limited memory demands. Allows huge datasets to be clustered in reasonable time

Heuristic centroid-based clustering

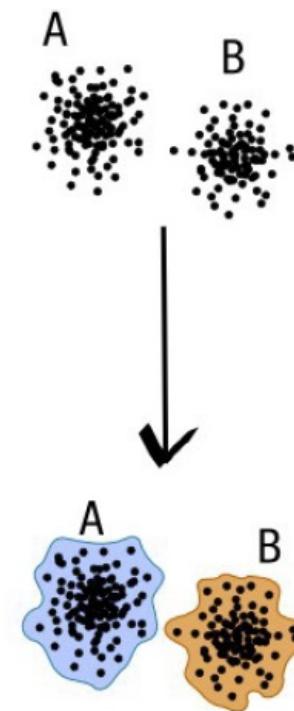


Clustering thresholds



compromise threshold
unadapted threshold

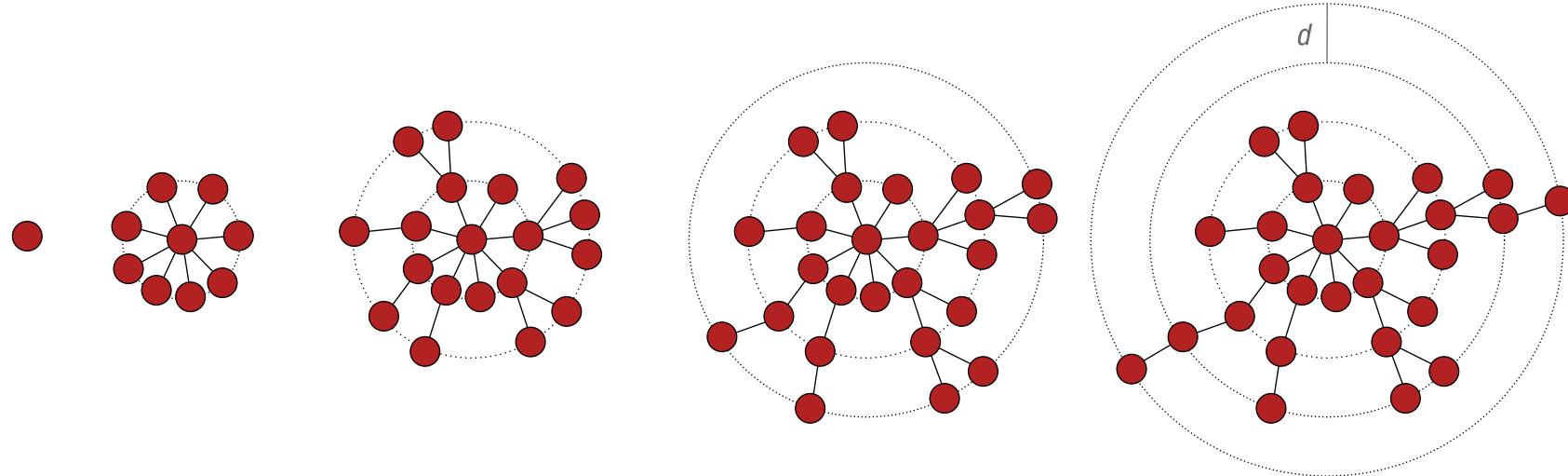
defined by maximum
distance from centroid



natural limits of OTUs

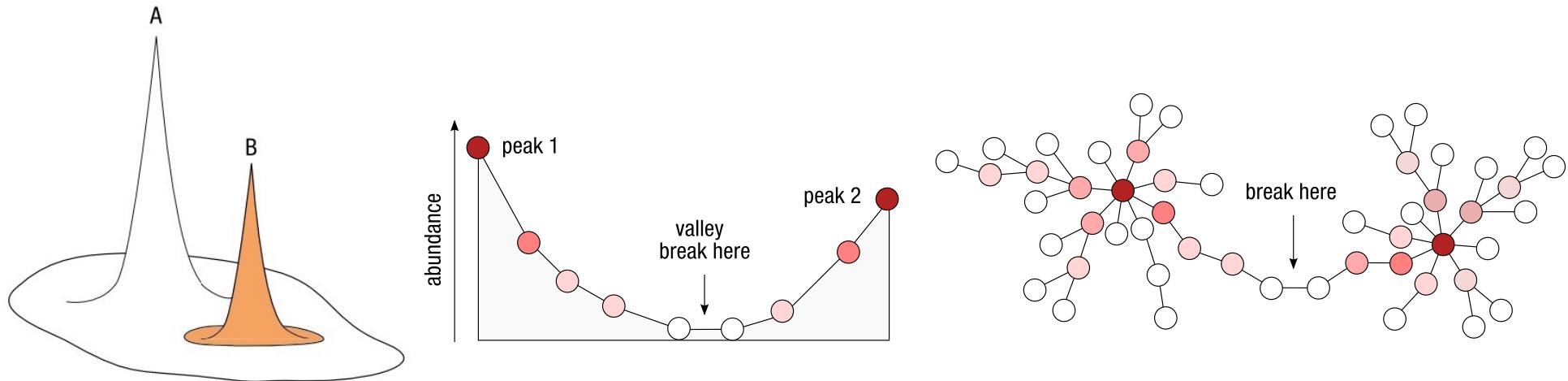
defined by minimum
distance separating clusters

Linking amplicons



- Amplicons (nodes) are linked when distance is less than or equal to d .
- Similar to single linkage hierarchical clustering
- Forms a spanning tree
- The adjustable parameter d is 1 by default.

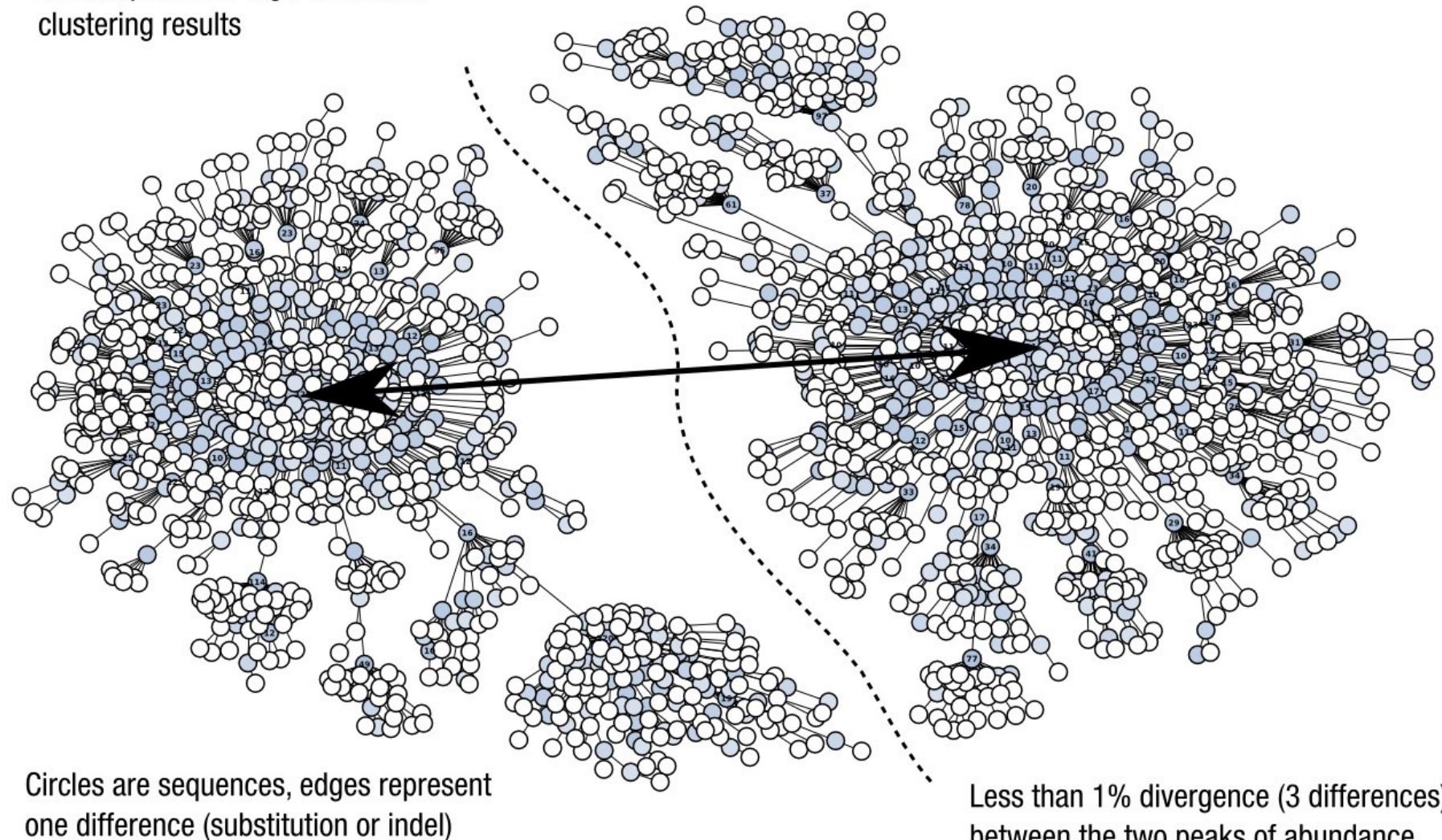
Breaking long chains of amplicons



- OTUs can be considered mountains in sequence space with peaks corresponding to the most abundant sequences
- Chains of linked amplicons can be broken to avoid problems with over-grouping
- Breakpoint in valleys between abundance peaks.
- Start at highest abundance amplicons and only link amplicons with lower or equal abundance

Visualising the clustering results

Swarm produces high-resolution clustering results



Availability & documentation

- Available on GitHub: <https://github.com/torognes/swarm>
- Free of charge, Open source license (GPL)
- Extensive documentation
- Precompiled 64-bit binaries for Linux, macOS and Windows
- Available for *de novo* OTU picking in QIIME 1.9
- Publications
 - Mahé F, Rognes T, Quince C, de Vargas C, Dunthorn M. (2014) **Swarm: robust and fast clustering method for amplicon-based studies.** PeerJ 2:e593.
 - Mahé F, Rognes T, Quince C, de Vargas C, Dunthorn M. (2015) **Swarm v2: highly-scalable and high-resolution amplicon clustering.** PeerJ 3:e1420.
- Publication on Swarm 3 recently submitted

Main international collaborators

Frédéric Mahé
CIRAD, Montpellier, France



Lucas Czech
Heidelberg Inst. for Theoretical Studies, Germany



Tomáš Flouri
Heidelberg Inst. for Theoretical Studies, Germany



Christopher Quince
Warwick & Glasgow Univ., UK





Thank you!