

Task Group	Task	Description	Number papers reporting task	Number papers not reporting software	Total number of software tools	Total number of software functions	Number of papers performing manually
Read preparation	<i>quality control</i>	<i>Generating a report of sequence quality information from a sample or set of samples - no modification is done to data</i>	19	0	4	4	0
	<i>adapter trimming</i>	<i>Trimming of sequencing adapters</i>	9	1	6	6	0
	<i>demultiplexing</i>	<i>Separation of sequences from a mixed pool into separate pools based on the occurrence of a unique set of bases (index or tag)</i>	55	17	16	19	0
	<i>pair merging</i>	<i>The assembly of mate pair reads into a single contig</i>	63	1	10	18	0
	<i>quality trimming</i>	<i>The removal of bases from either or both ends of sequences in a pool based on quality scores</i>	20	1	8	10	0
	<i>mate pairing</i>	<i>The identification and synchronisation of mate pair reads between two samples, often involving arranging reads in identical orders and/or removal of reads without a mate pair</i>	3	0	3	3	0
	<i>primer trimming</i>	<i>Trimming of PCR primers</i>	66	8	15	17	0
	<i>reverse complementation</i>	<i>Reverse complementing the sequences in a pool</i>	7	3	2	2	0
	<i>sequence conversion</i>	<i>Converting sequences from fastq to fasta</i>	3	0	2	3	0
	<i>length trimming</i>	<i>The removal of bases from either or both ends of sequences in a pool, either the removal of a fixed number of bases or the removal of a variable number of bases to reduce sequences to a standard length</i>	10	3	6	7	0
	<i>pair concatenation</i>	<i>Concatenating mate pair reads into a single contig (where reads don't overlap)</i>	8	4	4	4	0
	<i>assembly</i>	<i>The assembly of reads into contigs, applied when more than one pair of overlapping fragments have been metabarcoded</i>	6	0	4	4	0
	<i>degapping</i>	<i>Removal of gaps from sequences</i>	1	0	1	1	0
Sequence processing	<i>dereplication</i>	<i>The removal of duplicate reads to retain only unique sequences in a pool; often the total number of copies of a sequence is recorded in the header of the retained sequence</i>	58	10	11	19	0
	<i>size sorting</i>	<i>The sorting of a fasta file according to a size annotation in the header</i>	10	2	3	4	0
Filtering	<i>quality filtering</i>	<i>Removal and/or trimming of sequences from a pool based on quality information. Also often converts from fastq to fasta.</i>	81	11	20	27	0
	<i>similarity filtering</i>	<i>Removal of sequences based on similarity to an alignment, either based on sequence identity or alignment position</i>	9	1	4	4	0

	length filtering	<i>The removal of sequences from a pool that are less than, more than, or fall within or outside of a specified length threshold or thresholds</i>	54	21	17	23	0
	preclustering	<i>Reduction of sequence variation in a dataset prior to further processing - a form of denoising</i>	12	1	3	6	0
	denoising	<i>The removal of reads containing putative PCR or sequencing errors based on statistical assessment</i>	18	1	8	8	0
	normalisation	<i>A process by which the number of sequences for each of a set of samples is reduced where necessary such that the output set of samples all have the same number of sequences while maintaining the relative frequencies of OTUs</i>	2	0	1	1	1
	chimera filtering	<i>The filtering of putative chimeric assemblies from a pool of mate paired reads</i>	63	4	6	16	1
	translation filtering	<i>Removal of sequences from a set of sequence based on their translation, usually removing sequences with inframe stop codons or frameshifts due to erroneous indels or substitutions caused by sequencing errors</i>	22	3	11	12	0
	frequency filtering	<i>Removal of sequences based on their frequency in a pool</i>	51	37	11	15	1
	taxonomy filtering	<i>Removal of sequences based on an assigned taxonomy or a taxonomic classification</i>	9	5	1	1	1
	mistag filtering	<i>Removal of sequences based on putative tagging errors</i>	3	1	1	1	0
Data generation	OTU delimitation	<i>The grouping of a set of sequences into OTUs by some method</i>	84	5	12	22	0
	OTU mapping	<i>The mapping of sequences to OTUs to provide read counts for each OTU</i>	30	3	7	11	0
	uncurated taxonomic assignment	<i>The assignment (identification or classification) of taxonomy to OTUs using a global uncurated reference database (e.g. GenBank, BOLD)</i>	55	2	11	13	0
	reference taxonomic assignment	<i>The assignment (identification or classification) of taxonomy to OTUs using a purpose-built and/or specially curated reference set of sequences</i>	60	9	18	23	1

Table 1: Table of all bioinformatic tasks performed across the core papers set. Tasks are grouped into four groups by broad purposes, and a detailed definition of each task is given along with summary statistics of the implementation of each task across the 111 papers. For a list of the software used for each task, Table S1 is an expanded version of this table.