Microbiome

SOFTWARE

**Open Access**

# Cultivation-independent genomes greatly expand taxonomic-profiling capabilities of mOTUs across various environments

Hans-Joachim Ruscheweyh[1†], Alessio Milanese[1,2†], Lucas Paoli[1], Nicolai Karcher[2], Quentin Clayssen[1], Marisa Isabell Keller[2], Jakob Wirbel[2], Peer Bork[2,3,4], Daniel R. Mende[5], Georg Zeller[2*] and Shinichi Sunagawa[1*]

## Abstract

**Background:** Taxonomic profiling is a fundamental task in microbiome research that aims to detect and quantify the relative abundance of microorganisms in biological samples. Available methods using shotgun metagenomic data generally depend on the deposition of sequenced and taxonomically annotated genomes, usually from cultures of isolated strains, in reference databases (reference genomes). However, the majority of microorganisms have not been cultured yet. Thus, a substantial fraction of microbial community members remains unaccounted for during taxonomic profiling, particularly in samples from underexplored environments. To address this issue, we developed the mOTU profiler, a tool that enables reference genome-independent species-level profiling of metagenomes. As such, it supports the identification and quantification of both "known" and "unknown" species based on a set of select marker genes.

**Results:** We present mOTUs3, a command line tool that enables the profiling of metagenomes for >33,000 species-level operational taxonomic units. To achieve this, we leveraged the reconstruction of >600,000 draft genomes, most of which are metagenome-assembled genomes (MAGs), from diverse microbiomes, including soil, freshwater systems, and the gastrointestinal tract of ruminants and other animals, which we found to be underrepresented by reference genomes. Overall, two thirds of all species-level taxa lacked a reference genome. The cumulative relative abundance of these newly included taxa was low in well-studied microbiomes, such as the human body sites (6–11%). By contrast, they accounted for substantial proportions (ocean, freshwater, soil: 43–63%) or even the majority (pig, fish, cattle: 60–80%) of the relative abundance across diverse non-human-associated microbiomes. Using community-developed benchmarks and datasets, we found mOTUs3 to be more accurate than other methods and to be more congruent with 16S rRNA gene-based methods for taxonomic profiling. Furthermore, we demonstrate that mOTUs3 increases the resolution of well-known microbial groups into species-level taxa and helps identify new differentially abundant taxa in comparative metagenomic studies.

[†]Hans-Joachim Ruscheweyh and Alessio Milanese contributed equally to this work.

*Correspondence:  zeller@embl.de; ssunagawa@ethz.ch

[1] Department of Biology, Institute of Microbiology and Swiss Institute of Bioinformatics, ETH Zürich, 8093 Zürich, Switzerland
[2] Structural and Computational Biology Unit, European Molecular Biology Laboratory, 69117 Heidelberg, Germany
Full list of author information is available at the end of the article

**Conclusions:** We developed mOTUs3 to enable accurate species-level profiling of metagenomes. Compared to other methods, it provides a more comprehensive view of prokaryotic community diversity, in particular for currently underexplored microbiomes. To facilitate comparative analyses by the research community, it is released with >11,000 precomputed profiles for publicly available metagenomes and is freely available at: https://github.com/motu-tool/mOTUs.

**Keywords:** Metagenomics, Microbial community, Benchmarking, Taxonomic profiling, Marker gene, Metagenome-assembled genome, Single-cell genome, Reference genome

## Background

Identifying and quantifying the abundance of taxa (i.e., taxonomic profiling) is a critical step in linking the composition of microbial communities to environmental functions and host health-related phenotypes [1, 2]. Metagenomic sequencing of DNA directly extracted from an environmental or host-derived sample has enabled researchers to taxonomically profile microbial communities in an unbiased and cultivation-independent manner. The development of tools to generate accurate taxonomic profiles from metagenomic data has therefore become important to our understanding of microbial communities [3]. However, existing tools rely on the availability of informative sequences (such as k-mers or marker genes [4, 5]), which are predominantly extracted from taxonomically annotated reference genomes (RefGs).

In recent years, high-throughput culturing of microorganisms coupled with RefG sequencing (known as culturomics) [6] has substantially expanded the proportion of microbial taxa with whole genome sequences in data repositories (e.g., NCBI RefSeq) benefitting taxonomic profiling tools. However, there is a strong bias toward microorganisms from well-studied habitats (e.g., human body sites) and/or those that can be readily cultivated using standard laboratory methods. Thus, most microbes on Earth remain uncultivated and lack a representative RefG [7, 8], although they can be both globally prevalent [9] and numerically dominant in many environments [10–13]. As a result, the incorporation of RefGs from newly isolated microbes into taxonomic profiling tools can be slow and disproportional across environments. This poses an additional challenge for accurate taxonomic profiling, given that microorganisms that remain undetected bias the abundance estimates of those that are detected [14, 15].

To close the gap between the detectable and actual diversity present in microbial community samples, we developed mOTUs [14, 16], a software tool that uses universal, protein-coding, single-copy phylogenetic marker gene (MG) sequences to quantify the taxonomic composition of microbial communities from metagenomic sequence data (for further applications, see also Ruscheweyh et al. 2021 [17]). As these MGs are present in all organisms, they can be identified not only in RefGs, but also in metagenomic assemblies. Conceptually, mOTUs is based on clustering sets of MGs representing individual organisms by sequence similarity into species-level units. In the absence of a generalizable species concept for prokaryotes [18, 19], we refer to these units as MG-based operational taxonomic units (abbreviated as "mOTUs").

As an alternative to RefG sequencing, draft genomes are increasingly reconstructed by computational binning of metagenomic assemblies into metagenome-assembled genomes (MAGs [20]) or by sequencing amplified DNA from individual cells, resulting in single cell genomes (SAGs [21]). These cultivation-independent methods have provided genomic access to microbial diversity in previously underexplored environments. Several MAGs have been produced by different studies (>150,000 MAGs available in NCBI GenBank), and recently, some profiling tools had their databases extended with large-scale MAG collections from the human gut [22, 23]. The benefit of this approach was however limited to a single environment, for which comprehensive MAG datasets have recently been established [24]. Here, in addition to MGs found in RefG and metagenomic data, we now incorporate those found in MAGs and SAGs from various environments to more than double the number of taxa represented, adding >20,000 new mOTUs compared to the previous major release [14]. Our evaluations show that mOTUs3 outperforms other methods as assessed using metrics for taxonomic tool benchmarking developed independently from our study [3, 25]. Furthermore, we found mOTUs3 to provide an unprecedented view of the species-level diversity within the most dominant heterotrophic bacterial clade in the ocean and to greatly extend the number of detected and differentially abundant species in cross-sectional studies, as exemplified in a comparison between rumen microbiomes of high- and low-level methane-emitting sheep.
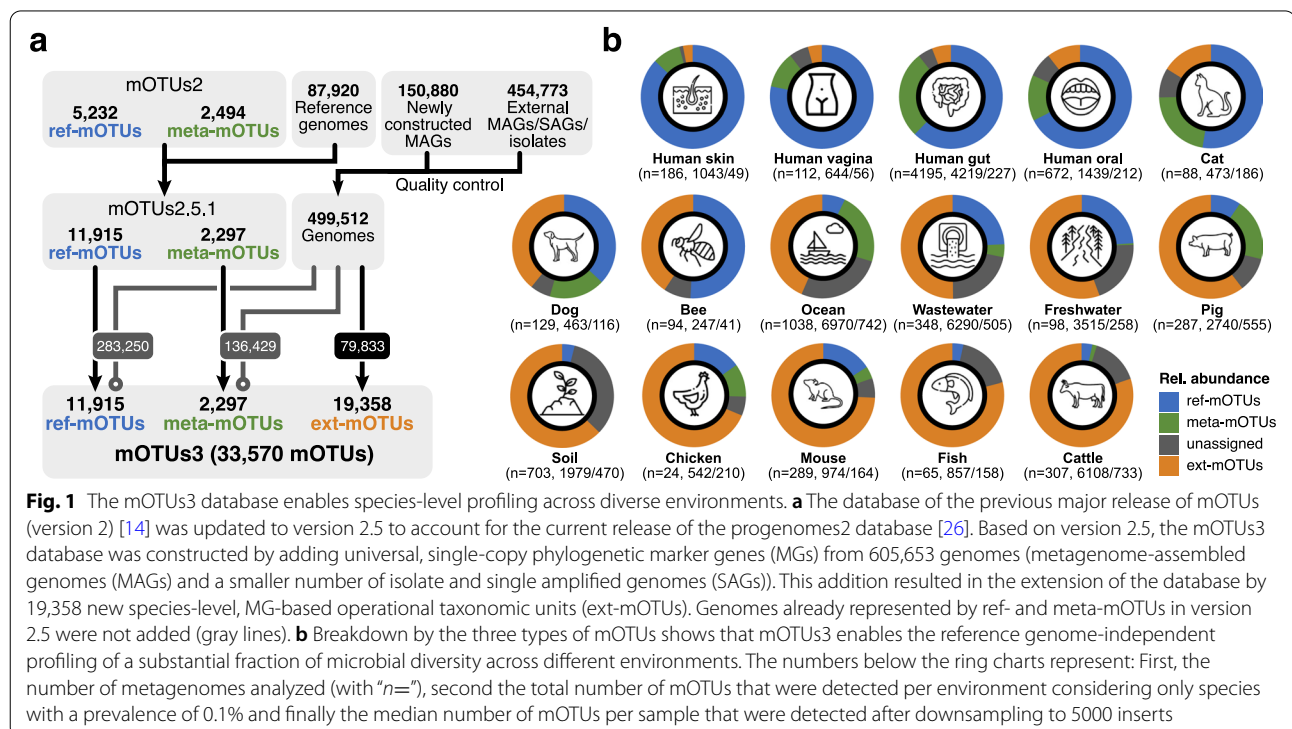
## Results

### Taxonomic profiling of diverse environments with mOTUs3

We developed mOTUs3 to facilitate the metagenomic profiling of 33,570 mOTUs, which is a 4.3-fold increase

compared to mOTUs2 (Fig. 1a). Among all mOTUs, 35% were represented by a RefG (*n*=11,915; ref-mOTUs), while an additional 21,655 were derived using MGs from either metagenomic contigs (*n*=2297; meta-mOTUs) or extended sources, such as MAGs (de novo-assembled or imported) and a smaller number of SAGs and isolate genomes (*n*=19,358; ext-mOTUs), to substantially extend the database coverage for reference genome-independent taxonomic profiling of diverse environments (Fig. 1a). MGs not assigned to any mOTU were additionally added to the database and merged into a single "unassigned" group to improve the quantification accuracy of taxonomic profiles, as previously demonstrated [14].

The newly established database allowed us to determine and systematically compare the fraction of taxa currently not represented by RefGs in various environments. These environments include extensively studied human-associated ones, for which metagenomic studies are complemented by several culturomics efforts (e.g., Lagier et al. [27]). Furthermore, we included data from >20 environmental and animal-associated microbiomes (Supplementary Tables 1 and 2) that have been primarily studied by metagenomic approaches. Overall, we found that more than half (11,882) of all meta/ext-mOTUs (i.e., mOTUs not represented by any RefG) could not be assigned to any known family (Supplementary Table 3; Methods), illustrating the taxonomic novelty covered by mOTUs3. The distribution of the newly included data into ref/ meta/ext-mOTUs was highly variable across the different environments (Supplementary Fig. 1). As expected, 97% of the ~400,000 MAGs from human microbiome samples (Supplementary Table 1) had already been represented by 2360 pre-existing (i.e., ref/meta-)mOTUs (Supplementary Table 4). Notably, the remaining 3% represented 2750 new ext-mOTUs, showing that novel species can still be uncovered by studying underrepresented populations, dietary habits, and/or disease states [28, 29]. By contrast, we found that only ~25% of the 6479 MAGs from mouse gut metagenomes (Supplementary Table 1) corresponded to pre-existing mOTUs (*n*=72, 68 ref-mOTUs and 4 meta-mOTUs), despite ongoing cultivation efforts [6]; the remaining 75% were grouped into 587 ext-mOTUs (Supplementary Table 4), meaning that 90% of the mOTUs represent novel species, which is in accordance with recently published studies [11, 30]. However, the vast majority of ext-mOTUs (*n*=16,021) resulted from the inclusion of other animal-associated (e.g., ruminants, fish, chicken, pig, bee, dog, cat) and environmental (e.g., soil, freshwater, wastewater, ocean, air) microbiomes (Supplementary Table 1) for which the generation of representative RefGs is lagging.

We used mOTUs3 to profile 10,541 available shotgun metagenomic data sets across the 23 environments covered by its database (Supplementary Table 1). For comparative analyses, we subset the data to 5756 high-quality samples (Methods; Supplementary Table 5)



**Fig. 1** The mOTUs3 database enables species-level profiling across diverse environments. **a** The database of the previous major release of mOTUs (version 2) [14] was updated to version 2.5 to account for the current release of the progenomes2 database [26]. Based on version 2.5, the mOTUs3 database was constructed by adding universal, single-copy phylogenetic marker genes (MGs) from 605,653 genomes (metagenome-assembled genomes (MAGs) and a smaller number of isolate and single amplified genomes (SAGs)). This addition resulted in the extension of the database by 19,358 new species-level, MG-based operational taxonomic units (ext-mOTUs). Genomes already represented by ref- and meta-mOTUs in version 2.5 were not added (gray lines). **b** Breakdown by the three types of mOTUs shows that mOTUs3 enables the reference genome-independent profiling of a substantial fraction of microbial diversity across different environments. The numbers below the ring charts represent: First, the number of metagenomes analyzed (with "*n*="), second the total number of mOTUs that were detected per environment considering only species with a prevalence of 0.1% and finally the median number of mOTUs per sample that were detected after downsampling to 5000 inserts

from 16 environments and found the overall number of detected mOTUs to range from 247 (honey bee) to >6000 (ocean, wastewater and cattle microbiomes). To illustrate the proportion of quantifying taxa currently not represented by RefGs (Fig. 1b), we summarized the cumulative relative abundances of unassigned taxa and the different types of mOTUs (ref-mOTUs, meta-mOTUs, ext-mOTUs). The fraction of unassigned taxa was highest for soil samples (33%; s.d. 8%), which reflects the high microbial diversity in soil as well as challenges in reconstructing genomes from this environment [31] (Supplementary Fig. 2). By contrast, more than 87% (s.d. 0.7%) of the relative abundance was represented by ref-mOTUs in human skin samples (in agreement with previous reports [32]) mainly due to the dominance of few taxa with cultivated representatives [33]. Similarly, the fraction of relative abundance assigned to ext-mOTUs varied considerably between environments: on average, only ~6% of the bacterial abundance in human-associated samples was assigned to newly added taxa, while this fraction was as high as ~80% in cattle rumen microbiomes.

### Comparison with other taxonomic profilers

As in other fields of bioinformatics, there is broad consensus that the performance of analysis tools needs to be carefully evaluated. However, best practices (e.g., balancing precision and recall, selecting criteria for "best" performance) are often debated [34, 35], and in microbiome research, an agreement on some fundamental concepts (e.g., sequence vs. taxonomic abundance, representation of unknown taxa in ground truth data) is still lacking [36, 37]. In an attempt to address some of these issues in a community-driven effort, modeled after successful examples in other fields [38, 39], the Critical Assessment of Metagenome Interpretation (CAMI) has provided curated ground truth datasets along with a tool (OPAL) to reproducibly evaluate metagenomic analysis tools [3, 25].

Using the latest CAMI datasets with disclosed results [40], we compared mOTUs3 to its prior major release version (mOTUs2) [14] and other selected metagenomic profiling tools (MetaPhlAn3 [5] and Bracken [4, 41], Methods) representing conceptually different, well-performing approaches to taxonomic profiling [36]. Using the OPAL tool for scoring and evaluation, we first evaluated presence/absence ($F_1$-score) and relative abundance predictions (L1 norm error) at the species level. For the different datasets, which represented samples from five human body sites and the mouse gut microbiome, mOTUs3, and MetaPhlAn3 performed generally better than Bracken and mOTUs2 (Fig. 2a/b). At higher taxonomic ranks, mOTUs3 had similar or higher scores than the other tools. For some datasets, taxonomic ranks,
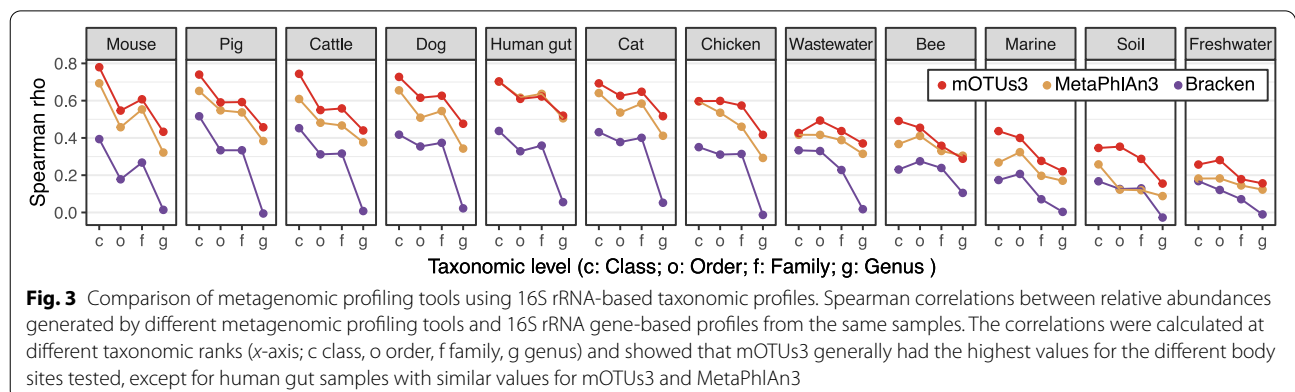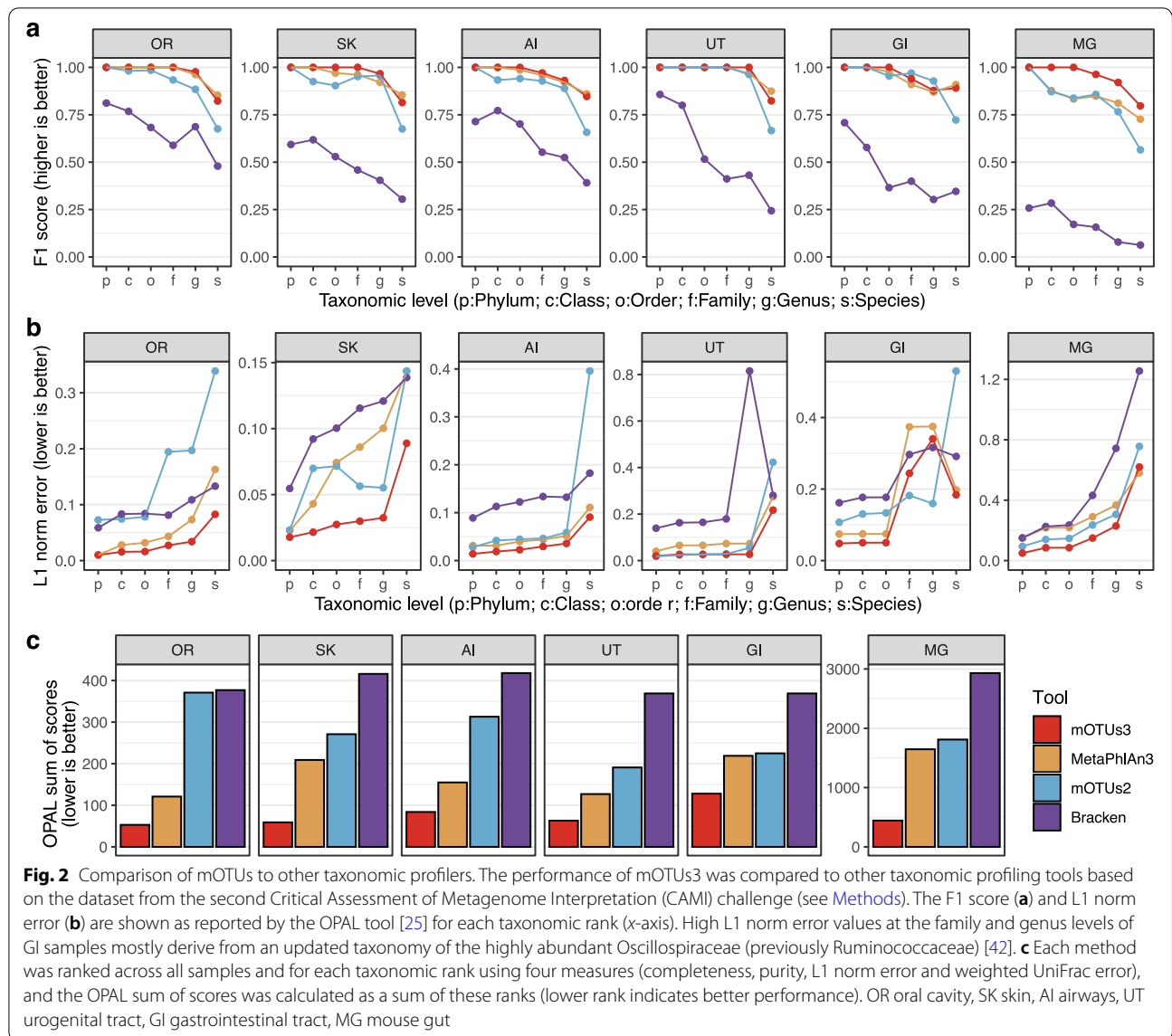
and tools, there was little to no room for improvements of the $F_1$-score or L1 norm error. This may be due to the simulated datasets being mainly based on taxa for which RefGs are available and/or result from incongruencies of taxonomic annotations used by the different profilers compared to the ground truth. In addition to the L1 norm error, OPAL computes additional metrics for profiling quality (completeness, purity, weighted UniFrac error) and summarizes them across taxonomic ranks into a composite score. Based on this evaluation criterion, mOTUs3 outperformed the other tools (Fig. 2c), as well as additional tools assessed in the CAMI challenge (Methods; Supplementary Fig. 3).

In the absence of independent ground truth data sets to benchmark taxonomic profiling tools for less well-studied environments, we correlated taxonomic profiles obtained by mOTUs3 and other tools to those obtained by analyzing 16S rRNA gene (16S) fragments. This approach leverages both the availability of comprehensive 16S databases for taxonomic classification [43] and the possibility of estimating taxonomic abundances based on 16S-based data from metagenomes [44]. Briefly, we extracted 16S fragments from the same datasets we used for metagenomic profiling and generated relative abundance profiles for them (Methods). To ensure comparability between 16S and metagenomic profiles, the analysis was performed at the genus and higher taxonomic ranks (for discussion, see Salazar et al. [44]). We found that mOTUs3 had consistently higher correlations with 16S profiles than the other tools across all environments, except for the human gut for which MetaPhlAn3 showed correlation coefficients similar to those of mOTUs3 (Fig. 3).

### Resolving the diversity of Pelagibacterales with mOTUs3

In addition to the broader taxonomic coverage by mOTUs3 across environments, we sought to investigate the capability of mOTUs3 to resolve microbial clades into more fine-grained taxonomic units. To this end, we focused on Pelagibacterales (also referred to as the SAR11 clade), which is the most abundant heterotrophic bacterial group in the global oceans [45]. Members of the Pelagibacterales have previously been shown to display high genomic variability while maintaining highly conserved 16S sequences [46]. This prompted us to evaluate the species-level resolution of mOTUs3 and to compare the diversity represented by mOTUs to the diversity represented by operational taxonomic units (OTUs) defined by 16S sequence similarity.

For this analysis, we selected from all mOTUs annotated as Pelagibacterales ($n$=1029; 2063 genomes) those that were represented by genomes with complete 16S sequences ($n$=602; 1105 genomes). The number of mOTUs was comparable to the number resulting from
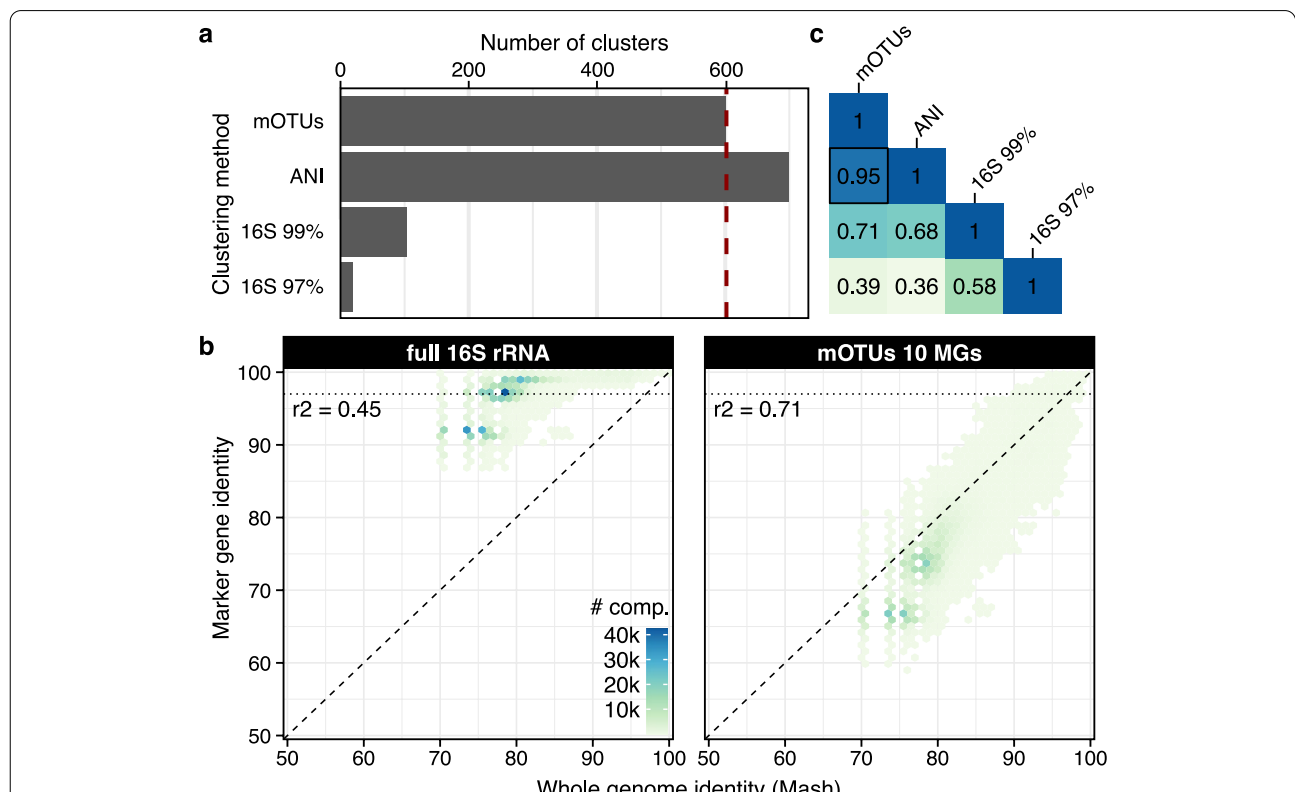
**Fig. 2** Comparison of mOTUs to other taxonomic profilers. The performance of mOTUs3 was compared to other taxonomic profiling tools based on the dataset from the second Critical Assessment of Metagenome Interpretation (CAMI) challenge (see Methods). The F1 score (**a**) and L1 norm error (**b**) are shown as reported by the OPAL tool [25] for each taxonomic rank (*x*-axis). High L1 norm error values at the family and genus levels of GI samples mostly derive from an updated taxonomy of the highly abundant Oscillospiraceae (previously Ruminococcaceae) [42]. **c** Each method was ranked across all samples and for each taxonomic rank using four measures (completeness, purity, L1 norm error and weighted UniFrac error), and the OPAL sum of scores was calculated as a sum of these ranks (lower rank indicates better performance). OR oral cavity, SK skin, AI airways, UT urogenital tract, GI gastrointestinal tract, MG mouse gut



**Fig. 3** Comparison of metagenomic profiling tools using 16S rRNA-based taxonomic profiles. Spearman correlations between relative abundances generated by different metagenomic profiling tools and 16S rRNA gene-based profiles from the same samples. The correlations were calculated at different taxonomic ranks (*x*-axis; c class, o order, f family, g genus) and showed that mOTUs3 generally had the highest values for the different body sites tested, except for human gut samples with similar values for mOTUs3 and MetaPhlAn3

a 95% average nucleotide identity (ANI)-based cluster-
ing of the 1105 genome sequences into species-level
groups ($n$=700; Fig. 4a), which is common practice in
the field of microbial phylogenomics [7, 47]. Moreo-
ver, we found sequence identities of mOTU-repre-
senting MGs to linearly correlate with those of whole
genomes across the whole range of observed values
($r^2$=0.71; Fig. 4b). By contrast, 16S sequence-based
OTUs using a 97% or 99% sequence similarity cutoff
resulted in a 31.7-fold ($n$=19) or 5.8-fold ($n$=104) lower
number of taxonomic units, respectively, compared
to mOTUs (Fig. 4a). This discrepancy is also reflected
by a weaker correlation ($r^2$=0.45; Fig. 4b) of identi-
ties between 16S sequences and corresponding whole
genome sequences. The minimum 16S identities were
ca. 87% and started saturating at approximately 97%
at which point genome identities were still as low as
~70–80% (Fig. 4b). Similar findings were reported pre-
viously albeit on smaller datasets [46]. Finally, compar-
ing the grouping of genomes by mOTUs and ANI into
species-level clusters, we found almost perfect congru-
ence (Fig. 4c, Methods).

## Differential abundance of novel archaea in low-/ high-methane-emitting sheep rumen metagenomes

High-resolution taxonomic profiling of metagenomes
from underexplored environments can be achieved by
custom-made marker gene or genome databases selected
for the microbial community under study [12, 48]. How-
ever, this approach is often labor- and resource-intensive
and requires specialized expertise, and its results cannot
easily be compared across studies and communities. To
demonstrate the utility of mOTUs3 to address these chal-
lenges, we reanalyzed rumen metagenomes from high-
and low-methane-emitting (HME and LME) sheep [48].
Importantly, these data were not used for the database
construction of mOTUs3.

Based on mOTUs3 taxonomic profiles, we identified
131 microbial species that differed significantly in abun-
dance between HME and LME samples and showed an



**Fig. 4** Species-level diversity of Pelagibacterales as resolved by mOTUs3. **a** The number of taxonomic units within the Pelagibacterales order varies depending on the clustering method used, which was based on using marker gene (MG) sequences (used by mOTUs), average nucleotide identity (ANI) of whole genomes, and full-length 16S rRNA gene sequences. **b** mOTU marker gene distances better capture whole-genome distances compared to full-length 16S, explaining the patterns observed in **a**. In particular, 16S rRNA gene sequence identity saturates while whole-genome similarity can be as low as 70–80%. **c** The different clustering approaches vary in their agreement with each other as determined by the V-measure, which captures both the completeness and homogeneity of the clusterings. The highest agreement was found between mOTUs and with whole genome clustering by ANI
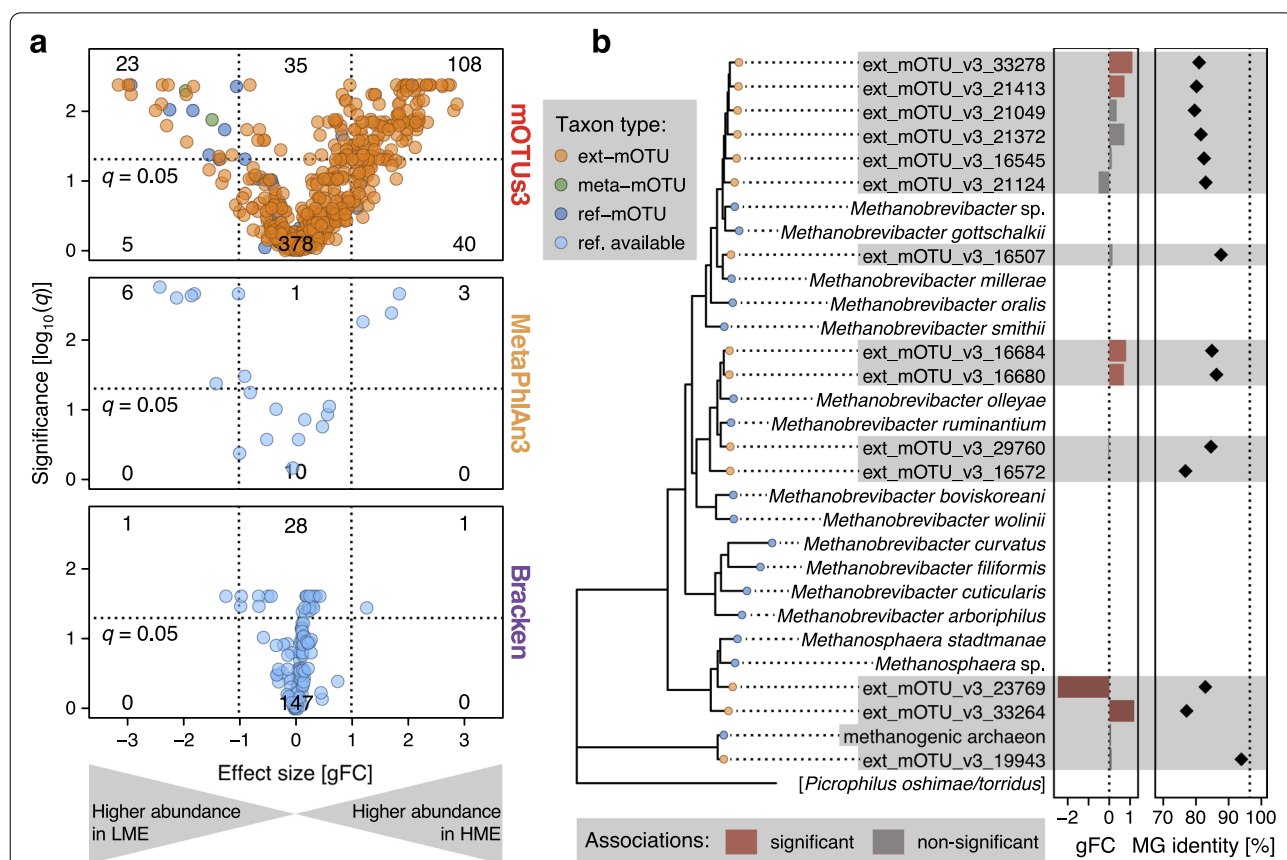
at least tenfold increase or decrease in relative abundance (corresponding to a generalized fold change of $\geq 1$ [49]). Among these differentially abundant species, 92% were represented by ext-mOTUs. These were therefore not expected to be detectable by reference-based profilers. To test this, we applied the same workflow using MetaPhlAn3 and Bracken (see Methods), which yielded only 10 and 30 differentially abundant species for the respective tools (Fig. 5a).

Given the metabolic importance of methanogenic archaea in ruminants as well as previous evidence of uncharted archaeal diversity in the sheep rumen [12], we further investigated the species-level diversity of known and unknown archaeal species. To this end, we reconstructed a phylogenetic tree of the archaeal mOTUs detected in the sheep rumen metagenomes ($n$=15) and contextualized them with reference genomes from members of the genera *Methanobrevibacter* and

*Methanosphaera* (Fig. 5b). This analysis revealed that all six differentially abundant archaea in the sheep rumen corresponded to ext-mOTUs. Two of them, which were significantly more abundant in high-methane emitters, were most closely related to *Methanobrevibacter gottschalkii*, which itself was not detected. Notably, the MG sequence similarity between these ext-mOTUs and *M. gottschalkii* was <85% (Fig. 5b), which is well below the species-level cutoff of 96.5% used by mOTUs [16] and therefore suggests that these ext-mOTUs represent novel *Methanobrevibacter* spp.

## Discussion

With mOTUs3, we have developed a taxonomic profiler that combines state-of-the-art accuracy, as demonstrated in competitive benchmarks based on simulated datasets, with an innovative database construction approach to detect and quantify underrepresented microbes from



**Fig. 5** Detection of differentially abundant taxa in low-/high-level-methane-emitting sheep rumen microbiomes. **a** A comparison between metagenomic profilers shows that mOTUs3 detected 131 differentially abundant species (*q* value <0.05 and an absolute generalized fold change > 1; indicated by dotted lines) between low- and high-level methane-emitting sheep, while MetaPhlAn3 and Bracken detected nine and two species, respectively. Most of the species detected by mOTUs were represented by ext-mOTUs only, demonstrating the added value of reference genome-independent profiling enabled by mOTUs3. **b** Archaeal mOTUs present in the sheep rumen microbiome (highlighted in gray) were phylogenetically contextualized with *Methanobrevibacter* spp. and *Methanosphaera* spp. represented by ref-mOTUs. All differentially abundant ext-mOTUs (middle panel) correspond to distinct yet undescribed *Methanobrevibacter* spp. as supported by all MG sequence identities ("MG identity [%]", right panel) to the closest known reference genome being below the species-level cutoff of 96.5% (dotted vertical line)

diverse environments at high (i.e., species-level) taxonomic resolution. Its core database will be updated with high-quality genomes (e.g., from proGenomes 3). Furthermore, the possibility to incorporate MG sequences from any MAG and SAG to generate mOTUs de novo and independently from the availability of RefGs and/or prior existence of taxonomic annotations (such as NCBI or GTDB species names) allows users to continuously extend the core database of mOTUs to represent microbial diversity from newly explored microbiomes. Such future extensions could also target eukaryotic microorganisms, as these are an integral part of many microbial communities, but are not well represented in databases of existing taxonomic profiling tools.

However, the flexibility in defining operational taxonomic units de novo comes with a need for taxonomic annotation, as is also the case for 16S rRNA-based de novo clustered OTUs. Despite the calibration of MG sequence identity cutoffs to maximize congruence with the NCBI taxonomy [16], this procedure can lead to conflicts with existing taxonomies. Irrespective of the ongoing debate on whether prokaryotic species should be consistent with genomic similarity-based criteria, delineating species by sequence identity puts mOTUs at a disadvantage in benchmarks, such as CAMI, which rely on rigid matching of taxonomic labels. The high performance of mOTUs [40] despite this disadvantage is likely due to the higher number of quantified taxa and the resulting reduction in compositionality-related biases.

## Conclusions

The present work introduces mOTUs3 as a reference-genome independent tool that allows for charting the taxonomic landscape of many environments at species-level resolution. Its independence from taxonomically annotated reference genomes makes it generally applicable also beyond well-studied environments to quantify and reveal yet uncharacterized microbial species of potential biological relevance. To support the research community, mOTUs3 is documented and available as open source software at https://github.com/motu-tool/mOTUs.

## Methods

### Collection and processing of data to compile the mOTUs3 database

To extend the taxonomic coverage of the mOTUs3 database, 4531 publicly available metagenomic datasets from 23 environments (Supplementary Table 1) were processed to generate 150,880 MAGs as previously described [50]. Briefly, BBMap (v.38.71) was used to quality control sequencing reads from all samples by removing adapters from the reads, removing reads that mapped to quality

control sequences (PhiX genome) and discarding low-quality reads (*trimq=14*, *maq=20*, *maxns=1*, and *minlength=45*). For metagenomic data of human origin, human genome-derived reads were removed using the masked human reference genome provided by BBMap. Quality-controlled reads were merged using bbmerge.sh with a minimum overlap of 16 bases, resulting in merged, unmerged paired, and single reads. The reads were assembled into scaffolded contigs (hereafter scaffolds) using the SPAdes assembler (v3.14 or v3.12) [51] in metagenomic mode. Genes were predicted on length-filtered ($\geq$ 500 bp) scaffolded contigs (hereafter scaffolds) using Prodigal (v2.6.3) [52]. Universal single-copy phylogenetic marker genes (MGs) were extracted using fetchMGs (v1.2; *-m extraction*) [16].

Scaffolds were length-filtered ($\geq$ 1000 bp) and within each study, quality-controlled reads from each sample were mapped against the scaffolds of each sample. Mapping was performed using BWA (v0.7.17-r1188; *-a*) [53]. Alignments were filtered to be at least 45 bp in length, with an identity of $\geq$ 97% and a coverage of $\geq$ 80% of the read sequence. The resulting BAM files were processed using the *jgi_summarize_bam_contig_depths* script of MetaBAT2 (v2.12.1) [20] to compute within- and between-sample coverages for each scaffold. The scaffolds were binned by running MetaBAT2 on all samples individually (*--minContig 2000* and *--maxEdges 500* for increased sensitivity). Metagenomic bins were annotated with Anvio (v5.5.0) [54], quality-controlled using the CheckM (v1.0.13) [55] lineage workflow (completeness $\geq$ 50% and contamination < 10%) to generate 150,880 MAGs. Complete genes were predicted using Prodigal (v2.6.3; *-c -m -g 11 -p single*), and MGs were extracted using fetchMGs (v1.2) *(-m extraction -v -i)*. These MAGs were complemented with 454,773 external genomes (~96% MAGs; ~4% isolate and single-cell genomes) from previous work (Supplementary Table 1), for which MGs were extracted using the same settings we used for MAGs.

All genomes containing at least six out of the 10 MGs used by mOTUs [16] were kept to produce the dataset of MGs from a total of 499,512 de novo-generated MAGs and external genomes for the construction of the mOTUs3 database. To evaluate the quality of these genomes, we calculated the agreement of the taxonomic annotation of the marker genes within each genome (Supplementary Fig. 4).

### Construction of the mOTUs3 database

MGs from 499,512 genomes were mapped against the latest mOTUs database (v2.5.1), which was an update of version 2.0 to account for a more recent release of the progenomes2 database [26] (Fig. 1a) using vsearch [56]

(v2.14.1; *--usearch_global --strand both --id 0.8 --max-accepts 10000 --maxrejects 10000*). MGs from a total of 283,250 and 136,429 genomes were assigned to existing ref-mOTUs and meta-mOTUs, respectively. These genomes were removed since they were already represented. The remaining 79,833 genomes resulted in an extension of the mOTUs database by 19,358 new mOTUs (ext-mOTUs). For consistency with the taxonomic annotation of ref-mOTUs, ext-mOTUs were annotated using the STAG classifier (https://github.com/zellerlab/stag, version 0.7; default parameters) trained on genomes in the proGenomes2 database [26] (NCBI taxonomy, version: 8 January 2019). MGs identified on scaffolds that were not binned into MAGs were used to update the "unassigned" mOTU, which contain unbinned MGs that are used to estimate the quantity of unknown species, by aligning these MGs against the extended database using vsearch (v2.14.1; *usearch_global --maxaccepts 1000 --maxrejects 1000 --strand both*). MGs that did not align within MG-specific cutoffs [57] were clustered using vsearch (v2.14.1; *--cluster_fast*) using MG-specific cutoffs and the representative sequence was added to the unassigned mOTU.

**Computation of mOTUs3 profiles for comparative analyses**
A total of 11,164 metagenomic and metatranscriptomic samples (Supplementary Tables 1 and 2) were quality controlled and merged as described above and profiled with mOTUs3 using default parameters and the *-c* option to build a community resource of taxonomic profiles. For comparative analyses across environments, 5756 of these samples were used after removing all (*n*=623) metatranscriptomic samples, metagenomic samples from environments with too few samples (termite, panda, aerosols, and bioreactor) or from studies comprising samples from different environments and samples with less than 5000 mapped inserts. To calculate the total number of detected mOTUs for a given environment, we counted the number of mOTUs with a prevalence greater than 0.1% (Supplementary Table 5). To compare the median number of detected mOTUs across different environments, we downsampled the insert counts to 5000 using the *rrarefy* function of the vegan package [58].

**Comparison of taxonomic profilers using the CAMI framework**
The performance of mOTUs3 was evaluated and compared to mOTUs2 and other taxonomic profilers by analyzing 113 publicly available samples (49 human-associated, 63 mouse gut metagenomes) provided by the second CAMI challenge (https://cami-challenge.org/participate). The samples were profiled with mOTUs3 (v3.0.1; *-C precision*), mOTUs2 (v2.1.1; *-C*

*precision*), MetaPhlAn3 (v3.0.7; *--CAMI_format_output --index mpa_v30_CHOCOPhlAn_201901*) [5], and Kraken/Bracken (v2.1.2; *--db=k2_standard_20201202 --paired* / v2.6.1; *--db=k2_standard_20201202 -r 100 -l S|G|F|O|C|P|D*) [4, 41]. Kraken/Bracken reports were further translated into the CAMI format ed files using the *tocami.py* script provided at https://github.com/hzi-bifo/cami2_pipelines. For comparative analyses, the OPAL framework (v1.0.9) [25] was used with default parameters providing the gold standard with the parameter *--gold_standard_file,* the names of the tools with *--labels,* the description with *-d*, the output with *--output_dir,* and the taxonomic profiles files as positional arguments.

**Comparison of metagenomic profiles with 16S rRNA gene-based profiles**
The 16S rRNA-based taxonomic profiler mTAGs [44] (v1.0.1; *-ma 1000 -mr 1000*) was used to generate relative abundance profiles for metagenomic samples (Supplementary Table 1). The output of mTAGs was mapped to the NCBI taxonomy to facilitate comparative analysis. The same samples were profiled with MetaPhlAn3 (v3.0.7; *--index mpa_v30_CHOCOPhlAn_201901*) and Kraken/Bracken (v2.1.2; *--db=k2_standard_20201202 --paired* / v2.6.1; *--db=k2_standard_20201202 -r 100 -l S*). Samples with small read/insert coverages (mTAGs<10,000, mOTUs<1000, Kraken/Bracken<10,000, no filtering was done on MetaPhlAn3 as profiles contain relative abundances) were removed, leaving 6119 samples for comparative analysis. Spearman correlations were calculated for each taxonomic rank based on concatenated relative abundances between mTAGs and the metagenomic profiling tools.

**Comparison of Pelagibacterales genome clusters with marker gene and 16S rRNA gene sequences**
Out of 2063 genomes belonging to 1029 mOTUs annotated as Pelagibacterales, 1105 genomes (from 602 mOTUs) that contained a complete copy of the 16S rRNA gene were selected. These genomes were also clustered based on average nucleotide identity using dRep [59] (v2.5.4; *-comp 0 -con 1000 -sa 0.95 -nc 0.2*) using a 95% cutoff as part of the OMD [50]. In addition, these genomes were clustered based on their 16S rRNA gene identity (99% and 97%) using vsearch [56] (v2.14.1; *--cluster_smallmem --id 0.97 / 0.99*). The consistency between the different clustering approaches was evaluated using the V-measure, which combines both the homogeneity and completeness metrics [60].

To correlate distances of the 1105 genomes between the different clustering techniques, we performed exhaustive distance calculations at the whole-genome level, the 10 MGs used by mOTUs, and the 16S rRNA gene. Whole

genome distances were computed using MASH [61] as implemented in dRep (v2.5.4). MG- and 16S rRNA gene-based distances were computed using vsearch (v2.14.1; *--allpairs_global --id 0.0*), and MG distances were averaged across the 10 genes prior to computing correlations.

## Differential abundance of mOTUs between low-/high-methane-emitting sheep

Samples from sheep rumen metagenomes (*n*=16) [48] were profiled with mOTUs3 (v3.0.1; *-c*), MetaPhlAn3 (v3.0.7; *--index mpa_v30_CHOCOPhlAn_201901*), and Kraken/Bracken (v2.1.2; *--db=k2_standard_20201202 --paired* / v2.6.1; *--db=k2_standard_20201202 -r 100 -l S*). To test for differentially abundant species between low-methane emitters (LMEs) and high-methane emitters (HMEs), the respective profiles were analyzed using SIAMCAT default workflows [49]. This workflow includes filtering of species/mOTUs with a relative abundance of >0.1% in at least one sample [49]. Wilcoxon test results were corrected for multiple testing using the Benjamini–Hochberg method [62] at 5% FDR. The reported effect size measure is the generalized fold change (gFC), calculated as the log10 of the geometric mean of quantile differences between groups as defined in SIAMCAT [49].

A phylogeny was constructed for all archaeal mOTUs belonging to the *Methanobrevibacter* and *Methanosphaera* genera or the *Thermoplasmata* class that passed the relative abundance filtering (14 ext-mOTUs, 1 ref-mOTU) together with ref-mOTUs from *Methanobrevibacter* and *Methanosphaera* (*n*=15) and a randomly selected Thermoplasmata ref-mOTU as an outgroup. Representative genomes from these 31 mOTUs were selected either by picking the centroid genome (for ext-mOTUs) or the reference genome (for ref-mOTUs). Marker genes were individually aligned (*mafft* [63], v7.458), the alignments were concatenated and a maximum-likelihood phylogeny was calculated using RAxML [64] (v8.2.12; *raxmlHPC -p 12345 -m PROTGAMMAAUTO*). The distance between the 14 ext-mOTUs and their closest ref-mOTU was calculated based on averaged marker gene distances across the 10 genes (v2.14.1; *vsearch --allpairs_global --id 0.0*).

## Availability and requirements

Project name: mOTUs

Project home page: https://github.com/motu-tool/mOTUs

Operating systems: Linux, MacOS

Programming language: Python 3

License: GNU General Public License v3.0

Any restrictions to use by non-academics: None

## Abbreviations

mOTUs3: A tool for marker gene-based OTU (mOTU) profiling; mOTU: Marker gene-based OTU; MAG: Metagenome-assembled genome; 16S rRNA: 16S ribosomal RNA; K-mer: A substring of a longer sequence of length k; RefG: Reference genome; NCBI RefSeq: National Center of Biotechnology Information (NCBI) Reference Sequence database; MG: Universal, protein-coding, single-copy phylogenetic marker gene; SAG: Single-cell amplified genome; NCBI GenBank: NIH genetic sequence database; ref-mOTU: A mOTU containing at least one reference genome; meta-mOTU: A mOTU built from co-abundance-binned marker genes across metagenomes; ext-mOTU: A mOTU built from genomes added in mOTUs3; CAMI: Critical Assessment of Metagenome Interpretation; OPAL: Open-community Profiling Assessment tooL; OTU: Operational taxonomic unit; ANI: Average nucleotide identity; HME/LME: High- and low-methane-emitting; GTDB: Genome Taxonomy Database; OMD: Ocean Microbiomics Database; FDR: False discovery rate.

## Supplementary Information

The online version contains supplementary material available at https://doi.org/10.1186/s40168-022-01410-z.

**Additional file 1: Supplementary Figure 1.** Environment-specific membership of genomes in ref-, meta- and ext-mOTUs. A total of 499,512 genomes derived from 23 environments (environments with few genomes are grouped as 'Other', see Supplementary Tables 1 and 3) were used for the extension. The number of genomes was normalized by environments. The proportions of genomes per environment that are either associated with ref- and meta-mOTUs or were used to build ext-mOTUs are shown in the colors blue, green or orange, respectively. For example, the majority of genomes from the human gut match ref-mOTUs, whereas the vast majority of genomes from the fish environment are used to build ext-mOTUs. **Supplementary Figure 2.** Comparison of Shannon index from profiling using mOTUs and 16S rRNA gene OTUs. In order to improve our understanding, the Shannon index evaluated with mOTUs (y-axis) and 16S rRNA OTUs (x-axis) reconstructed from the cattle and soil samples. Pearson correlation of indices generated from cattle profiles show a high agreement between mOTUs and 16S rRNA based methods whereas mOTUs underestimates species diversity for soil samples. The mOTUs profiles were generated using default parameters. For the 16S RNA profiles we extracted the first 100 bp from reads containing the V4 primer sequence and clustered at 97% identity using vsearch (*--derep_fulllength, --cluster_size --id 0.97, --usearch_global --id 0.97*). **Supplementary Figure 3.** OPAL score broken down to individual metrics for the 63 mouse gut metagenomic samples. The evaluation was performed using the OPAL tool [1] on 63 simulated mouse gut metagenomes [2], which also provided taxonomic profiles for seven different taxonomic profiling tools, and to which we have added mOTUs3 profiling results. The OPAL tool ranks the tools for each sample and for each taxonomic level. The measures considered are completeness, purity, L1 norm error and weighted UniFrac error, shown individually in the bottom 4 plots. Tools with a lower score perform better, as the OPAL score is a sum over rank. The top plot represents the OPAL sum of scores, which is the sum over the four individual measures. mOTUs3 scored best in all categories, including the OPAL sum of scores. **Supplementary Figure 4.** Taxonomic consistency of marker genes from the 499,512 genomes used to extend the mOTUs database. Marker genes from each genome were taxonomically annotated to evaluate taxonomic consistency. Agreeing, all marker genes have the same annotation; Majority agreeing, more than half of the marker genes agree to one taxonomic annotation; Not agreeing, there is no taxonomic annotation that agrees in more than 50% of the marker genes; Not annotated, there is no taxonomic annotation for this taxonomic level. Below the graph, we show a table with the percentage of Not agreeing annotations per taxonomic level, either as the percentage of all genomes (top) or of the genomes that have an annotation at that taxonomic level (bottom). **Supplementary Table 1.** Included studies and associated environments. Data from 91 studies from 23 environments were included in the extension and/or profiling of the mOTUs database. Of these, 39 studies were selected for in-house MAG reconstruction and 11,164 sequencing samples from 67 studies were used for taxonomic profiling. **Supplementary Table 2.** Sequencing samples included in the

taxonomic profile. A total of 11,164 samples were taxonomically profiled. Sample names are connected to public repositories by biosample and sequencing run ids. The project name column links the sample name to the study name used in Supplementary Table 1. **Supplementary Table 3.** Breakdown of taxonomic novelty in ext-mOTUs. Taxonomic novelty increases with higher ranks, i.e., more than 50% of ext-mOTUs were assigned to previously unknown families. **Supplementary Table 4.** Contribution of genomes to ref-, meta- or ext-mOTUs. Genomes/MAGs from different studies and environments contribute in varying proportions to the extension of the database. **Supplementary Table 5.** Data for Fig. 1. For each sample that passed the filter (total 5,756), we reported the relative abundance for each mOTU type. Additionally, we added the total number of detected mOTUs and the habitat. **Supplementary Table 6.** Data for Fig. 5a. Generalized fold change and adjusted p-value for species detected in 20 sheep rumen metagenomes when profiled with mOTUs3, Bracken or MetaPhlAn3.

## Authors' contributions
GZ and SS conceived and supervised the work. HJR and AM developed the code, generated the database with support from DRM, and performed the benchmark analysis. LP and NK performed the taxonomic diversity analysis of the SAR11 clade and the comparative metagenomic analysis, respectively. QC supported the collection and processing of the data. MIK and JW contributed to the taxonomic annotation of mOTUs. HJR, AM, LP, NK, PB, DRM, GZ, and SS wrote the manuscript. The authors read and approved the final manuscript.

## Availability of data and materials
The updated mOTUs3 database can be found at Zenodo (https://doi.org/10.5281/zenodo.5140350) and contains all MGs used in this study and the public profiles generated with mOTUs3. All genomes used to build the mOTUs3 database are deposited at the ETH Research Collection (https://doi.org/10.3929/ethz-b-000563099). A tool for programmatic access to all genome files is available at: https://github.com/motu-tool/motus_v3_genomes. Metadata of individual genomes, such as quality metrics and mOTU associations, are deposited at Zenodo (https://doi.org/10.5281/zenodo.6975138). The mOTUs extender software is publicly available at https://github.com/motu-tool/mOTUs-extender. A complete list with all sequencing samples used for building the database and/or for profiling can be found in Supplementary Tables 1 and 2.

## Declarations

### Ethics approval and consent to participate
Not applicable.

### Consent for publication
Not applicable.

### Competing interests
The authors declare that they have no competing interests.

## Author details
[1]Department of Biology, Institute of Microbiology and Swiss Institute of Bioinformatics, ETH Zürich, 8093 Zürich, Switzerland. [2]Structural and Computational Biology Unit, European Molecular Biology Laboratory, 69117 Heidelberg, Germany. [3]Max Delbrück Centre for Molecular Medicine, Robert-Rössle-Str. 10, 13092 Berlin, Germany. [4]Department of Bioinformatics, Biocenter, University of Würzburg, Am Hubland, 97074 Würzburg, Germany. [5]Department of Medical Microbiology, Amsterdam UMC, University of Amsterdam, Amsterdam, The Netherlands.

## References
1. Fuhrman JA. Microbial community structure and its functional implications. Nature. 2009;459:193–9 Nature Publishing Group.
2. Wang J, Jia H. Metagenome-wide association studies: fine-mining the microbiome. Nat Rev Microbiol. 2016;14:508–22.
3. Sczyrba A, Hofmann P, Belmann P, Koslicki D, Janssen S, Dröge J, et al. Critical assessment of metagenome interpretation—a benchmark of metagenomics software. Nat Methods. 2017;14:1063–71 Nature Publishing Group.
4. Wood DE, Salzberg SL. Kraken: ultrafast metagenomic sequence classification using exact alignments. Genome Biol. 2014;15:1–12 BioMed Central.
5. Beghini F, McIver LJ, Blanco-Míguez A, Dubois L, Asnicar F, Maharjan S, et al. Integrating taxonomic, functional, and strain-level profiling of diverse microbial communities with bioBakery 3. Elife. 2021:10. https://doi.org/10.7554/eLife.65088.
6. Lagkouvardos I, Pukall R, Abt B, Foesel BU, Meier-Kolthoff JP, Kumar N, et al. The Mouse Intestinal Bacterial Collection (miBC) provides host-specific insight into cultured diversity and functional potential of the gut microbiota. Nat Microbiol. 2016;1:16131.
7. Konstantinidis KT, Rosselló-Móra R. Classifying the uncultivated microbial majority: a place for metagenomic data in the Candidatus proposal. Syst Appl Microbiol. 2015;38:223–30.
8. Lloyd-Price J, Mahurkar A, Rahnavard G, Crabtree J, Orvis J, Hall AB, et al. Strains, functions and dynamics in the expanded Human Microbiome Project. Nature. 2017;550:61–6.
9. Hug LA, Baker BJ, Anantharaman K, Brown CT, Probst AJ, Castelle CJ, et al. A new view of the tree of life. Nat Microbiol. 2016;1:16048.
10. Salazar G, Paoli L, Alberti A, Huerta-Cepas J, Ruscheweyh H-J, Cuenca M, et al. Gene expression changes and community turnover differentially shape the global ocean metatranscriptome. Cell. 2019;179:1068–83.e21.
11. Lesker TR, Durairaj AC, Gálvez EJC, Lagkouvardos I, Baines JF, Clavel T, et al. An integrated metagenome catalog reveals new insights into the murine gut microbiome. Cell Rep. 2020;30:2909–22.e6.
12. Stewart RD, Auffret MD, Warr A, Walker AW, Roehe R, Watson M. Compendium of 4,941 rumen metagenome-assembled genomes for rumen microbiome biology and enzyme discovery. Nat Biotechnol. 2019;37:953–61 Nature Publishing Group.
13. Wilhelm RC, Cardenas E, Leung H, Maas K, Hartmann M, Hahn A, et al. A metagenomic survey of forest soil microbial communities more than a decade after timber harvesting. Sci Data. 2017;4:170092.
14. Milanese A, Mende DR, Paoli L, Salazar G, Ruscheweyh H-J, Cuenca M, et al. Microbial abundance, activity and population genomic profiling with mOTUs2. Nat Commun. 2019;10:1014.
15. Gloor GB, Macklaim JM, Pawlowsky-Glahn V, Egozcue JJ. Microbiome datasets are compositional: and this is not optional. Front Microbiol. 2017;8:2224.
16. Sunagawa S, Mende DR, Zeller G, Izquierdo-Carrasco F, Berger SA, Kultima JR, et al. Metagenomic species profiling using universal phylogenetic marker genes. Nat Methods. 2013;10:1196–9.
17. Ruscheweyh H-J, Milanese A, Paoli L, Sintsova A, Mende DR, Zeller G, et al. mOTUs: profiling taxonomic composition, transcriptional activity and strain populations of microbial communities. Curr Protoc. 2021;1:e218.
18. Rosselló-Mora R, Amann R. The species concept for prokaryotes. FEMS Microbiol Rev. 2001;25:39–67.

19. Staley JT. The bacterial species dilemma and the genomic-phylogenetic species concept. Philos Trans R Soc Lond B Biol Sci. 2006;361:1899–909.

20. Kang DD, Li F, Kirton E, Thomas A, Egan R, An H, et al. MetaBAT 2: an adaptive binning algorithm for robust and efficient genome reconstruction from metagenome assemblies. PeerJ. 2019;7:e7359.

21. Woyke T, Doud DFR, Schulz F. The trajectory of microbial single-cell sequencing. Nat Methods. 2017;14:1045–54.

22. Nayfach S, Shi ZJ, Seshadri R, Pollard KS, Kyrpides NC. New insights from uncultivated genomes of the global human gut microbiome. Nature. 2019;568:505–10.

23. Parks DH, Rigato F, Vera-Wolf P, Krause L, Hugenholtz P, Tyson GW, et al. Evaluation of the microba community profiler for taxonomic profiling of metagenomic datasets from the human gut microbiome. Front Microbiol. 2021;12:643682.

24. Almeida A, Nayfach S, Boland M, Strozzi F, Beracochea M, Shi ZJ, et al. A unified catalog of 204,938 reference genomes from the human gut microbiome. Nat Biotechnol. 2021;39:105–14.

25. Meyer F, Bremges A, Belmann P, Janssen S, McHardy AC, Koslicki D. Assessing taxonomic metagenome profilers with OPAL. Genome Biol. 2019;20:51.

26. Mende DR, Letunic I, Maistrenko OM, Schmidt TSB, Milanese A, Paoli L, et al. proGenomes2: an improved database for accurate and consistent habitat, taxonomic and functional annotations of prokaryotic genomes. Nucleic Acids Res. 2020;48:D621–5.

27. Lagier J-C, Khelaifia S, Alou MT, Ndongo S, Dione N, Hugon P, et al. Culture of previously uncultured members of the human gut microbiota by culturomics. Nat Microbiol. 2016;1:16203.

28. Wirbel J, Pyl PT, Kartal E, Zych K, Kashani A, Milanese A, et al. Meta-analysis of fecal metagenomes reveals global microbial signatures that are specific for colorectal cancer. Nat Med. 2019;25:679–89.

29. Pasolli E, Asnicar F, Manara S, Zolfo M, Karcher N, Armanini F, et al. Extensive unexplored human microbiome diversity revealed by over 150,000 genomes from metagenomes spanning age, geography, and lifestyle. Cell. 2019;176:649–62.e20.

30. Kieser S, Zdobnov EM, Trajkovski M. Comprehensive mouse microbiota genome catalog reveals major difference to its human counterpart. PLoS Comput Biol. 2022;18:e1009947.

31. Alteio LV, Schulz F, Seshadri R, Varghese N, Rodriguez-Reillo W, Ryan E, et al. Complementary metagenomic approaches improve reconstruction of microbial diversity in a forest soil. mSystems. 2020;5. https://doi.org/10.1128/mSystems.00768-19.

32. Nayfach S, Rodriguez-Mueller B, Garud N, Pollard KS. An integrated metagenomics pipeline for strain profiling reveals novel patterns of bacterial transmission and biogeography. Genome Res. 2016;26:1612–25.

33. Byrd AL, Belkaid Y, Segre JA. The human skin microbiome. Nat Rev Microbiol. 2018;16:143–55.

34. Buchka S, Hapfelmeier A, Gardner PP, Wilson R, Boulesteix A-L. On the optimistic performance evaluation of newly introduced bioinformatic methods. Genome Biol. 2021;22:152.

35. Marx V. Bench pressing with genomics benchmarkers. Nat Methods. 2020;17:255–8.

36. Sun Z, Huang S, Zhang M, Zhu Q, Haiminen N, Carrieri AP, et al. Challenges in benchmarking metagenomic profilers. Nat Methods. 2021;18:618–26.

37. Ye SH, Siddle KJ, Park DJ, Sabeti PC. Benchmarking metagenomics tools for taxonomic classification. Cell. 2019;178:779–94.

38. Stolovitzky G, Monroe D, Califano A. Dialogue on reverse-engineering assessment and methods: the DREAM of high-throughput pathway inference. Ann N Y Acad Sci. 2007;1115:1–22.

39. Moult J, Pedersen JT, Judson R, Fidelis K. A large-scale experiment to assess protein structure prediction methods. Proteins. 1995;23:ii–v.

40. Meyer F, Lesker T-R, Koslicki D, Fritz A, Gurevich A, Darling AE, et al. Tutorial: assessing metagenomics software with the CAMI benchmarking toolkit. Nat Protoc. 2021;16:1785–801.

41. Lu J, Breitwieser FP, Thielen P, Salzberg SL. Bracken: estimating species abundance in metagenomics data. PeerJ Comput Sci. 2017;3:e104 PeerJ Inc.

42. Zhang X, Tu B, Dai L-R, Lawson PA, Zheng Z-Z, Liu L-Y, et al. Petroclostridium xylanilyticum gen. nov., sp. nov., a xylan-degrading bacterium isolated from an oilfield, and reclassification of clostridial cluster III members into four novel genera in a new Hungateiclostridiaceae fam. nov. Int J Syst Evol Microbiol. 2018;68:3197–211.

43. Quast C, Pruesse E, Yilmaz P, Gerken J, Schweer T, Yarza P, et al. The SILVA ribosomal RNA gene database project: improved data processing and web-based tools. Nucleic Acids Res. 2013;41:D590–6.

44. Salazar G, Ruscheweyh H-J, Hildebrand F, Acinas SG, Sunagawa S. mTAGs: taxonomic profiling using degenerate consensus reference sequences of ribosomal RNA genes. Bioinformatics. 2021. https://doi.org/10.1093/bioinformatics/btab465.

45. Giovannoni SJ. SAR11 bacteria: the most abundant plankton in the oceans. Ann Rev Mar Sci. 2017;9:231–55.

46. Grote J, Thrash JC, Huggett MJ, Landry ZC, Carini P, Giovannoni SJ, et al. Streamlining and core genome conservation among highly divergent members of the SAR11 clade. MBio. 2012:3. https://doi.org/10.1128/mBio.00252-12.

47. Chaumeil P-A, Mussig AJ, Hugenholtz P, Parks DH. GTDB-Tk: a toolkit to classify genomes with the Genome Taxonomy Database. Bioinformatics. 2019. https://doi.org/10.1093/bioinformatics/btz848.

48. Shi W, Moon CD, Leahy SC, Kang D, Froula J, Kittelmann S, et al. Methane yield phenotypes linked to differential gene expression in the sheep rumen microbiome. Genome Res. 2014:1517–25. https://doi.org/10.1101/gr.168245.113.

49. Wirbel J, Zych K, Essex M, Karcher N, Kartal E, Salazar G, et al. Microbiome meta-analysis and cross-disease comparison enabled by the SIAMCAT machine learning toolbox. Genome Biol. 2021;22:93.

50. Paoli L, Ruscheweyh H-J, Forneris CC, Hubrich F, Kautsar S, Bhushan A, et al. Biosynthetic potential of the global ocean microbiome. Nature. 2022;607:111–8.

51. Nurk S, Meleshko D, Korobeynikov A, Pevzner PA. metaSPAdes: a new versatile metagenomic assembler. Genome Res. 2017;27:824–34.

52. Hyatt D, Chen G-L, Locascio PF, Land ML, Larimer FW, Hauser LJ. Prodigal: prokaryotic gene recognition and translation initiation site identification. BMC Bioinformatics. 2010;11:119.

53. Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform. Bioinformatics. 2009;25:1754–60.

54. Eren AM, Esen ÖC, Quince C, Vineis JH, Morrison HG, Sogin ML, et al. Anvi'o: an advanced analysis and visualization platform for 'omics data. PeerJ. 2015;3:e1319.

55. Parks DH, Imelfort M, Skennerton CT, Hugenholtz P, Tyson GW. CheckM: assessing the quality of microbial genomes recovered from isolates, single cells, and metagenomes. Genome Res. 2015;25:1043–55.

56. Rognes T, Flouri T, Nichols B, Quince C, Mahé F. VSEARCH: a versatile open source tool for metagenomics. PeerJ. 2016;4:e2584.

57. Mende DR, Sunagawa S, Zeller G, Bork P. Accurate and universal delineation of prokaryotic species. Nat Methods. 2013;10:881–4 Nature Publishing Group.

58. Oksanen J, Kindt R, Legendre P, O'Hara B, Stevens MHH, Oksanen MJ, et al. The vegan package. Community Ecol Package. 2007;10:719.

59. Olm MR, Brown CT, Brooks B, Banfield JF. dRep: a tool for fast and accurate genomic comparisons that enables improved genome recovery from metagenomes through de-replication. ISME J. 2017;11:2864–8.

60. Hirschberg JB, Rosenberg A. V-Measure: a conditional entropy-based external cluster evaluation: Columbia University; 2007. Available from: https://academiccommons.columbia.edu/doi/10.7916/D80V8N84

61. Ondov BD, Treangen TJ, Melsted P, Mallonee AB, Bergman NH, Koren S, et al. Mash: fast genome and metagenome distance estimation using MinHash. Genome Biol. 2016;17:132.

62. Benjamini Y, Hochberg Y. Controlling the false discovery rate: a practical and powerful approach to multiple testing. J R Stat Soc. 1995;57:289–300 Wiley.

63. Katoh K, Standley DM. MAFFT multiple sequence alignment software version 7: improvements in performance and usability. Mol Biol Evol. 2013;30:772–80.

64. Stamatakis A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies. Bioinformatics. 2014;30:1312–3.

## Publisher's Note