



Cultivation-independent high-quality microbial genome reconstruction from environmental samples with midi-metagenomics

John Vollmers, Maximiano Correa Cassal and Anne-Kristin Kaster

Genome Res. published online May 15, 2026

Access the most recent version at doi:[10.1101/gr.280099.124](https://doi.org/10.1101/gr.280099.124)

P<P Published online May 15, 2026 in advance of the print journal.

Open Access Freely available online through the *Genome Research* Open Access option.

Creative Commons License This article, published in *Genome Research*, is available under a Creative Commons License (Attribution 4.0 International), as described at <http://creativecommons.org/licenses/by/4.0/>.

Email Alerting Service Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or [click here](#).

Advance online articles have been peer reviewed and accepted for publication but have not yet appeared in the paper journal (edited, typeset versions may be posted when available prior to final publication). Advance online articles are citable and establish publication priority; they are indexed by PubMed from initial publication. Citations to Advance online articles must include the digital object identifier (DOIs) and date of initial publication.

To subscribe to *Genome Research* go to:
<https://genome.cshlp.org/subscriptions>

Cultivation-independent high-quality microbial genome reconstruction from environmental samples with midi-metagenomics

John Vollmers,¹ Maximiano Correa Cassal,¹ and Anne-Kristin Kaster^{1,2}

¹*Institut für Biologische Grenzflächen 5, Karlsruhe Institute of Technology, 76344 Eggenstein-Leopoldshafen, Germany;*

²*Institute for Applied Biosciences, Karlsruhe Institute of Technology, 76131 Karlsruhe, Germany*

Because the majority of microbial organisms still evade cultivation attempts, genomic insights into many taxa are limited to cultivation-independent approaches. However, current methods of metagenomics and single-cell genome sequencing have individual drawbacks, which can limit the quality and completeness of the reconstructed genomes. Current attempts to combine both approaches still use whole-genome amplification techniques, which are prone to bias. Here, we propose a novel approach for the purpose of genome reconstructions that utilizes the potential of cell sorting for targeted enrichment and depletion of different cell types to create distinct cell fractions with sufficient DNA amounts, circumventing amplification. By distributing sequencing efforts over these fractions as well as the original sample, coassemblies become highly optimized for coabundance variation–based binning approaches. “Midi-metagenomics” enables accurate metagenome-assembled genome (MAG) reconstruction from individual sorted samples with higher quality than coassembly and binning of multiple distinct samples and therefore improves analyses of uncultivated microorganisms.

[Supplemental material is available for this article.]

The vast majority of prokaryotes still evade cultivation attempts under laboratory conditions and therefore cannot be subjected to direct analysis via classic culture-based microbiological and biochemical methods (Steen et al. 2019; Bodor et al. 2020). The discrepancy between the small number of cultured prokaryotic strains compared with the vast diversity and ubiquity of uncultured microbes, commonly referred to as “microbial dark matter” (MDM) (Marcy et al. 2007; Rinke et al. 2013; Solden et al. 2016; Zha et al. 2022), represents a large reservoir of biotechnological and/or pharmaceutical potential that is still untapped (Bodor et al. 2020; Kaster and Sobol 2020; Escudeiro et al. 2022).

Advances in cultivation-independent methodologies, such as metagenomics and single-cell genomics (SCGs), have enabled sequence-based predictions of phylogenetic and functional characteristics of uncultivated microorganisms (Fig. 1; Tyson et al. 2004; Solden et al. 2016; Vollmers et al. 2017), but with each method having distinct advantages and disadvantages. In metagenomics (Fig. 1A), bulk DNA from a mixed community, such as an environmental microbiome, is extracted and sequenced (Schmeisser et al. 2007). Subsequent analyses steps can then reveal the phylogenetic and functional diversity of a given community and even enable the reconstruction of so-called “metagenome-assembled genomes” (MAGs) from uncultivated individual community members via contig-based binning methods (Liu et al. 2025). Binning of metagenomic contigs is a challenging and error-prone process, especially for highly complex communities (Ma et al. 2023) and low abundant organisms (Qayyum et al. 2025). As a result, MAGs are highly susceptible to chimerism and can show varying degrees of fragmentation and completeness in addition to contamination (Orakov et al. 2021; Vollmers et al. 2022).

Furthermore, MAGs are often limited to the most abundant species present in the sample and may not reliably resolve strain variants or elements of horizontal gene transfer (Vollmers et al. 2017).

SCGs (Fig. 1B) circumvent this problem by directly targeting individual cells, thereby enabling reliable genome resolution on strain level (Xu and Zhao 2018; Kaster and Sobol 2020). However, because a single prokaryotic cell contains only a few femtograms of DNA, a whole-genome amplification (WGA) is required because the minimum requirement for high-throughput sequencing is typically in the nanogram range (Hutchison and Venter 2006; Kalisky and Quake 2011). This is a severe disadvantage, as WGA not only is expensive and prone to contamination but usually yields extremely uneven read coverage, constituting bias that is particularly pronounced for genomes with high GC content and usually results in single-cell amplified genomes (SAGs) that are typically even more fragmented and incomplete than MAGs (Lasken and Stockwell 2007; Alteio et al. 2020; Kaster and Sobol 2020).

To minimize these drawbacks and maximize the advantages of both methods, there is a strong interest in combining single-cell and metagenomic approaches. A current example for such an attempt is “mini-metagenomics” (Yu et al. 2017), which targets small groups of 10–1000 cells (Fig. 1C). These cells are then sequenced together and subsequently treated as a simplified metagenome to efficiently reduce random amplification bias (Yu et al. 2017; Alteio et al. 2020). The DNA yield of such a small amount of cells is, however, still not sufficient to circumvent amplification. The relatively low complexity of such mini-metagenomes should, in theory, allow for better genome reconstructions than the more complex metagenome of the original community; however, this approach is still affected by systematic WGA bias, which may be caused by, for example, variations in GC content (Marine et al.

Corresponding authors: john.vollmers@kit.edu, kaster@kit.edu

Article published online before print. Article, supplemental material, and publication date are at <https://www.genome.org/cgi/doi/10.1101/gr.280099.124>. Freely available online through the *Genome Research* Open Access option.

© 2026 Vollmers et al. This article, published in *Genome Research*, is available under a Creative Commons License (Attribution 4.0 International), as described at <http://creativecommons.org/licenses/by/4.0/>.

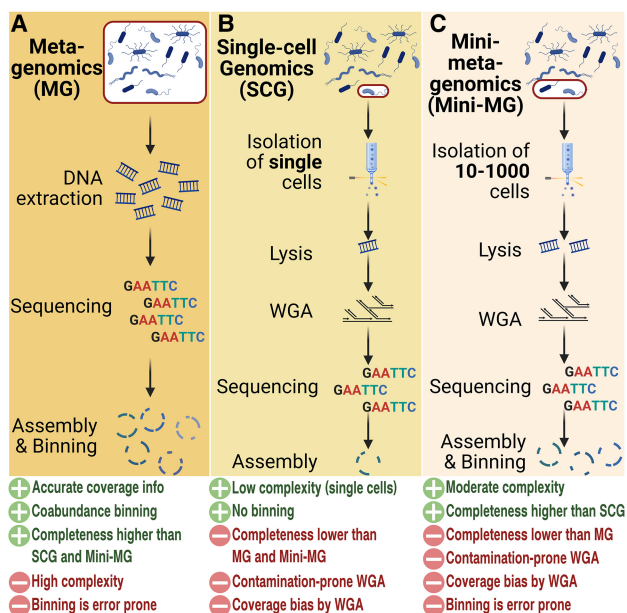


Figure 1. Current culture-independent methodologies. (A) In metagenomics, the entire DNA of an environmental community is sequenced. Assembled contigs are binned into metagenome-assembled genomes (MAGs). (B) In single-cell genomics (SCGs), individual cells are isolated, sequenced, and analyzed. Because of little DNA content per cell, whole-genome amplification (WGA) is required. (C) In mini-metagenomics, typically pools of five to 1000 cells are sorted and sequenced. Although complexity is lower than for standard metagenomics, binning is still required, and the low DNA content of small cell pools still necessitates WGA. Created with BioRender (<https://www.biorender.com>).

2014). Most importantly though, effective binning criteria are limited because contig abundance information is not available owing to the uneven read coverage, a severe drawback that also obstructs the currently most effective binning strategy: coabundance variation across samples (Alneberg et al. 2014; Mattock and Watson 2023). Therefore contigs will likely have to be binned exclusively based on nucleotide signatures, which is less reliable, especially for short contigs of highly fragmented genomes (Vollmers et al. 2022; Mattock and Watson 2023).

We here present an alternative approach, termed “midi-metagenomics,” that utilizes cell sorting to create custom community fractions of sufficient cell count to circumvent the need for amplification entirely. Fluorescence-activated cell sorting (FACS) is used for targeted enrichment and depletion of different cell types to create fractions that are highly optimized for coabundance variation-based binning approaches. This way, the quality of genome reconstructions can be maximized, even if only individual samples without spatial or temporal parallels are available.

Results

Established workflow

In midi-metagenomics, the original sample population is divided into multiple fractions, in which different community members are selectively enriched or depleted via FACS (Fig. 2A). Possible strategies for selectively fractionating a complex community into distinct subpopulations are manifold (Woyke et al. 2017; Sturm et al. 2023). However, in this proof-of-principle study, sorting

was based on relatively simple gating strategies exploiting only cell characteristics easily detectable via FACS: cell size as determined by forward scatter and “complexity” representing cell structures and granularity as defined by side scatter gating (Supplemental Fig. S1; Lavigne et al. 1997). Because soil represents one of the most complex and challenging microbial communities for metagenomic analyses (Vollmers et al. 2017; Jansson and Hofmöckel 2018), it was chosen as a test environment instead of controlled artificial consortia, which may not faithfully represent the diversity and dynamics of natural microbiomes.

In contrast to standard single-cell and “mini-metagenomic” approaches, which require an amplification step (Rinke et al. 2013; Yu et al. 2017; Alteio et al. 2020), the midi-metagenomic methodology utilizes bulk sorts of several hundred thousand to millions of cells into the same fraction. Preliminary trial DNA extractions performed on bulk sorts of bacterial cultures and soil samples indicated the presence of DNA predominantly in the supernatant and not the pellet (Wiegand et al. 2021) of centrifuged cell suspensions after FACS, especially in the case of soil samples (Supplemental Fig. S2). This observation indicates possible cell damage caused by the sorting process and subsequent release of

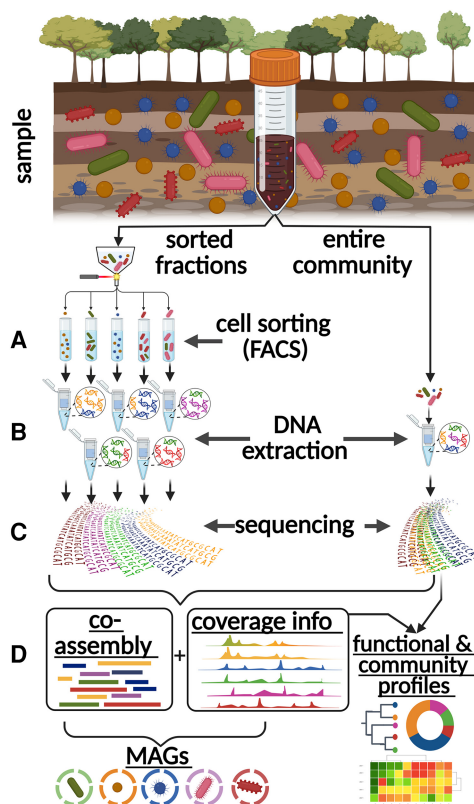


Figure 2. Midi-metagenomics workflow. (A) Part of the sample community is fractionated into distinct groups of several hundred thousand to millions of cells by cell sorting. (B) Different cell types are not separated with absolute stringency, but differentially enriched DNA is extracted separately from each fraction, as well as the original unsorted sample. (C) Extracted DNA is sequenced directly without whole-genome amplification (WGA). (D) Because the resulting read data sets represent different enrichments based on the same original community, they are optimal for coassembly as well as coabundance variation-based binning approaches. An unbiased representation of the source community is achieved by also including the original unsorted sample in the analyses. Created with BioRender (<https://www.biorender.com>).

cellular DNA (Mollet et al. 2008; Blainey 2013; Binek et al. 2019; Wiegand et al. 2021). Therefore, an adapted DNA extraction protocol was used for midi-metagenomic fractions, which includes an alcohol precipitation step directly from sorted cell suspensions rather than centrifuged cell pellets, thereby ensuring maximized DNA yields (Fig. 2B). In preliminary trial runs, DNA yields ranged between 5 and 30 ng DNA for up to 5 million sorted cells (Supplemental Table S1), which is more than sufficient for direct sequencing (Sato et al. 2019; Ribarska et al. 2022). Therefore, based on the low input requirements of <1 ng (Parkinson et al. 2012; Bowman et al. 2013; Tvedte et al. 2021) for modern sequencing library preparation techniques, sorting efforts may be reduced to 10,000–1,00,000 cells per sorted fraction in the future. The separate sequencing of genomic DNA for each fraction, as well as the original unfractionated sample (Fig. 2C), resulted in multiple distinct data sets for each sample.

Efficiency of cell sorting–based community fractionation

The relationship between sorted fractions and corresponding unsorted samples was analyzed based on 16S rRNA gene diversity within the assembled metagenomic and midi-metagenomic fractions (Fig. 3; Supplemental Table S2) as well as 16S rRNA amplicons with increased sequencing depth (Supplemental Fig. S3; Supplemental Table S3). Weighted UniFrac scores calculated from these analyses show higher beta-diversities between sorted fractions and their respective nonfractionated communities than between nonfractionated samples taken at different time points (Fig. 3). This increased beta-diversity represents a strong shift in relative taxon abundances within the respective microbial communities, which can be exploited for distinguishing different organisms based on differential coverage information during downstream binning attempts.

At the same time, all sorted fractions show decreased alpha-diversity values and, therefore, lower community complexity compared with their respective nonfractionated counterparts

(Supplemental Fig. S4). One sample exclusively sorted on forward scatter signals (which roughly indicate cell size) clusters closer to the unsorted sample when analyzed on high-depth amplicon sequencing level (Supplemental Fig. S3). This indicates that the use of the forward scatter alone provides a less systematic separation of the community compared with sorting based on combined forward and side scatter, possibly owing to the reduced resolution of morphological differences. Additional sorting metrics such as (auto)fluorescence should therefore even further improve binning efficiency.

Assembly and binning performance

Coassemblies of standard bulk metagenomics and midi-metagenomics were compared using the same total sequencing depth of 15 Gbp (averaging at 70 million read pairs per coassembly), equally distributed across the combined samples and fractions (Supplemental Table S4). Based on maximum contig length and N50 metrics, midi-metagenomic coassemblies of sorted fractions originating from the same original sample were significantly less fragmented than were coassemblies of distinct bulk metagenomic samples (Fig. 4), with $P < 0.001$ as determined via Mann–Whitney U tests (MacFarland and Yates 2016).

Improved assembly metrics also affect the distribution of quality categories among the produced MAGs, as defined by the “minimum information about a metagenome-assembled genome” (MIMAG) standard, developed by the Genomic Standards Consortium (Bowers et al. 2017): Even before contamination filtering with MDMcleaner (Vollmers et al. 2022), midi-metagenomic approaches produced far more high-quality MAGs (completeness >90%, contamination <5%) compared with standard metagenomic assemblies, which predominantly consisted of only moderate-quality genomes (completeness >50%, contamination <10%) or low-quality genomes (either completeness <50% or contamination >10%) with contamination values typically >5% (Fig. 5A–C; Supplemental Tables S5, S7). These trends persist across different

midi-metagenomic samples and coassembly subset groups of different sizes, despite varying community complexities, proving the robustness of the approach.

Interestingly, midi-metagenomic MAGs represented almost twice as many distinct phyla compared with standard MAGs (Fig. 5D), illustrating another important aspect of improved assembly and binning metrics, namely improved representation of original sample diversity. This is further corroborated by detailed phylogenomic analyses of low-contaminated MAGs (<5% contamination estimate) with at least moderate (50%) completeness (Fig. 6), which indicate a far broader and, owing to increased MAG qualities, also more reliable phylogenomic representation by midi-metagenomic MAGs compared with standard metagenomics. In addition, in the midi-metagenomic approach, a higher diversity of closely related but still distinct genomes could be reconstructed (as shown in Fig. 6 in the case of Acidobacteriota,

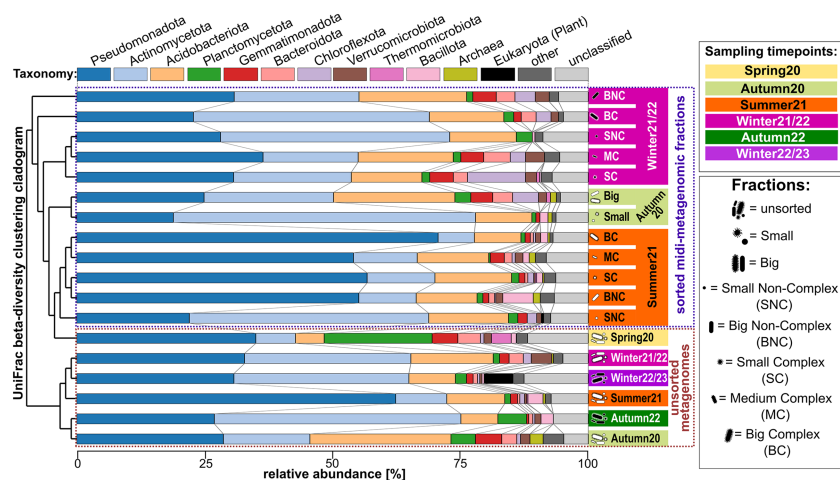


Figure 3. 16S rRNA gene-based diversity among different (midi-)metagenomic fractions of different samples. Clustering is based on weighted UniFrac (Lozupone et al. 2011) beta-diversity scores and is shown as a cladogram on the left. The background coloring of y-axis labels on the right side indicates the respective origin-sample. Stacked bar charts indicate the community composition of each sample and fraction, with different phyla being indicated by a distinct color code as indicated above the plot, and relative abundances being indicated by bar heights according to the x-axis below the plot. Sorted midi-metagenomic fractions are indicated by pictograms, and abbreviations are given in the legend on the right.

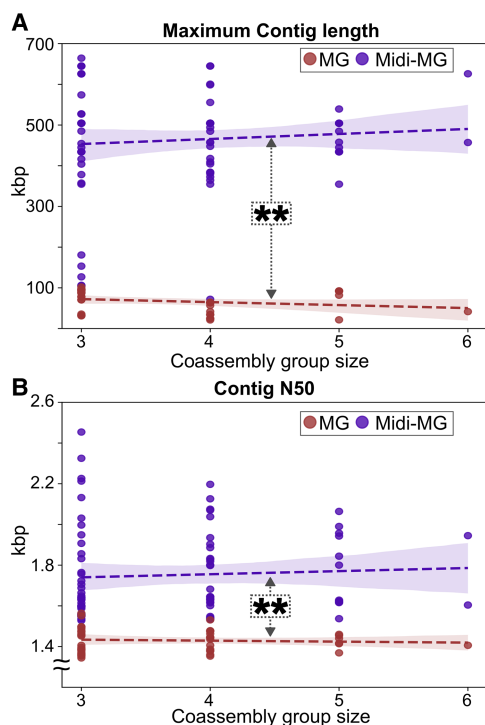


Figure 4. Comparison of assembly metrics for metagenomic and midi-metagenomic approaches in dependence on coassembled samples and fractions. Scatterplots showing metagenomic results as red, and midi-metagenomic results as purple dots. Trendlines and corresponding confidence areas were determined by regression analysis and are indicated by dashed lines and background coloring, respectively. (A) Maximum contig lengths. (B) Contig N50 values. The significance of differences in the distribution between metagenomic and midi-metagenomic assembly metrics were determined via Mann-Whitney *U* tests (MacFarland and Yates 2016). (MG) Metagenome, (Midi-MG) midi-metagenome.

and to a lesser extent also for Alphaproteobacteria and Actinomycetota), indicating a better resolution of sequence homologies between closely related organisms with this new method.

Comparison with alternative assembly and coverage distribution strategies

Although the coassembly of different fractions or samples enhances read coverage for assembly and coabundance based binning, there is some debate about its applicability for standard metagenomics, as seasonal or regional variations may introduce complexities that may obstruct optimal assembly (Olm et al. 2017). To present an objective comparison of the performance of midi-metagenomics versus metagenomics, we therefore tested alternative assembly strategies for the same overall sequencing depth of 15 Gbp (Supplemental Fig. S6). A common strategy is to perform separate single assemblies for each sample and then map every read data set against each individual assembly (Olm et al. 2017). This strategy greatly reduced the performance of both metagenomic and midi-metagenomic approaches, likely owing to the reduced read depth negatively affecting assembly, thereby causing an increase of “low-quality” MAGs and reducing the yield of “moderate-quality” to “high-quality” MAGs (Supplemental Fig. S6A). Nonetheless, midi-metagenomics still performed at least as well as standard metagenomics under these conditions.

Another potential alternative sequence coverage distribution strategy is to focus sequencing and assembly efforts only on one “main” sample and to supplement with additional lower depth “auxiliary” read data sets meant exclusively for mapping and binning purposes. However, in our trials, this strategy could not mitigate the negative effects of the “single assembly” strategy for standard metagenomics, instead only showing a slightly positive effect for midi-metagenomic data sets, which thereby outperformed standard metagenomics also under these conditions (Supplemental Fig. S6B).

Comparison with mini-metagenomics

For comparison purposes, mini-metagenomics was also applied to one sample. This approach is designed to reduce MDA bias by supplying higher amounts of input DNA by sorting and lysing multiple cells (Fig. 1C). Accordingly, we encountered fewer negative MDA reactions and more complete genomes using this approach compared with standard SCGs. However, only two moderate-quality mini-metagenomic MAGs could be recovered, both displaying high contamination estimates close to the MIMAG cutoff of 10% (Supplemental Table S5; Supplemental Fig. S6C).

Discussion

Midi-metagenomics integrates cell sorting and metagenome sequencing approaches into a new workflow that is optimized for high-quality MAG reconstruction. The key step of this approach is the separation of the sampled microbiome into highly

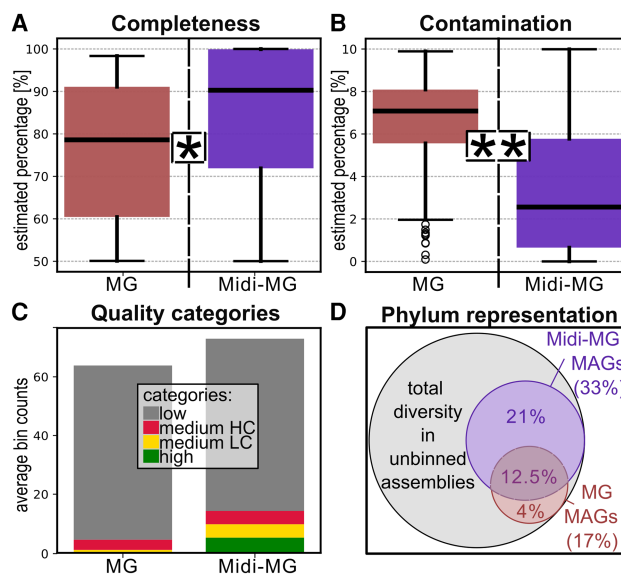


Figure 5. Comparison of quality metrics and diversity of MAGs obtained from standard metagenomic and midi-metagenomic coassemblies. (A, B) Boxplots showing the distribution of CheckM2 completeness and contamination estimates, respectively. The difference between metagenomic and midi-metagenomic results was found to be statistically significant ($P < 0.01$ based on the Moods median test) (Gibbons 2005) in both cases. A more detailed plot showing individual results is given in Supplemental Figure S5. (C) Average number of MAGs belonging to different quality categories obtained by standard metagenomics and midi-metagenomics. (D) Relative fractions of total phylum level diversity detected in the unbinned metagenomic coassemblies that are represented by metagenomic and midi-metagenomic MAGs, respectively. (MG) Metagenome, (Midi-MG) midi-metagenome.

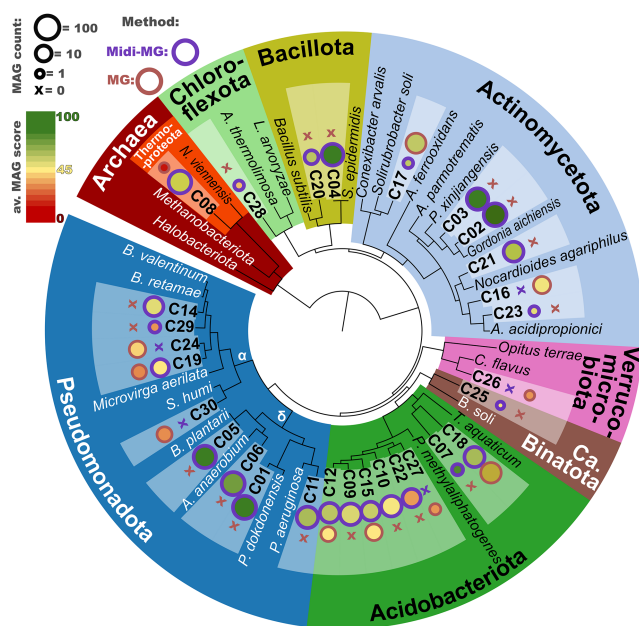


Figure 6. Multilocus sequence analysis (MLSA)-based phylogeny of representative MAGs and related reference genomes. Maximum likelihood phylogenetic clustering based on 61 single-copy orthologs shared by all comparison genomes, concatenated to a total length of 7178 amino acids. The software tool dRep (Olm et al. 2017) was used to group all MAGs on species level based on ANI comparisons and to rank the members of each group based on genome quality. Only groups with representatives showing >50% completeness and $\leq 5\%$ contamination were considered, resulting in 30 groups labeled C1–C30, for each of which only the best representative is compared. Bubble plots next to each group designation indicate the number and average dRep score of each group, as indicated by the legend on the upper left. (α) Alphaproteobacteria, (δ) gammaproteobacterial, (MG) metagenomics, (Midi-MG) midi-metagenomics.

comparable but nevertheless distinct fractions of reduced complexity (Fig. 3; Supplemental Fig. S2) that can be sequenced directly without involving WGA methods (Fig. 2).

Because the resulting fractions represent subsets of the same community at the exact same time point, they are optimally suited for a coassembly strategy, maximizing the available sequencing depth for assembly as well as binning purposes. Accordingly, midi-metagenomics yielded significantly better coassembly metrics and MAG qualities compared with standard metagenomics (Figs. 4, 5), the latter having produced significantly shorter contigs and lower-quality MAGs, likely owing to intersample heterogeneities, which are known to generally affect coassemblies of temporally or spatially distinct samples (Olm et al. 2017). The common alternative strategy of separately assembling each sample for metagenomics did also not improve assembly and binning outcomes (Supplemental Fig. S6), because it severely limits the read depth available for assembly per sample (Hofmeyr et al. 2020). Compensating this effect requires substantially increased sequencing efforts and costs to reach optimal coverage for each individual sample (Ma et al. 2023; Supplemental Table S7). This limitation is especially critical given that coabundance-based binning performs best across large numbers of parallel samples (Alneberg et al. 2014; Han et al. 2025), theoretically requiring deep sequencing of each individual sample.

Systematic comparisons with multiple independently published metagenome studies also confirm a generally lower efficien-

cy of MAG reconstruction for metagenome approaches compared with midi-metagenomics, which yielded 411 high-quality MAGs across all midi-metagenomic samples, with an average of three per coassembly subgroup (Supplemental Table S7). A large-scale analysis by Ma et al. (2023) of several thousands of publicly available as well as newly generated soil metagenomes established a direct correlation between sequencing depth and the number and quality of generated MAGs. The sequencing depths in that study ranged from 0.03 to 146 Gbp, with an average depth of 14 Gbp (Supplemental Table S7), which is comparable to the soil metagenome control used in our study. As to be expected, the number of high-quality genomes produced per metagenome varied drastically from zero to 72, averaging at three, with larger numbers being generated at extreme sequencing depths of ≥ 100 Gbp. Interestingly, $\sim 41\%$ of the 3304 data sets in that study did not yield any high-quality MAGs, indicating a generally relatively low efficiency of standard metagenomic approaches. Similar results were observed for two independent metagenome studies of agricultural soil samples (Supplemental Table S7; Nelkner et al. 2019; Hu et al. 2025), analyzing eight to 24 samples at depths of 10–12 Gbp each. These two studies yielded considerably lower counts of high-quality MAGs compared with midi-metagenomics despite surpassing the overall sequencing effort three- to 16-fold. These comparisons also showcase that a large fraction of metagenome samples are restricted to moderate- or lower-quality MAGs, mostly likely owing to interspecies homologies affecting assembly and binning. In contrast, midi-metagenomics yielded high-quality MAGs in almost every combination of samples and fractions and in larger proportions than any of the aforementioned approaches, indicating a far more efficient distribution of sequencing and sampling efforts. The only study to reach an equivalent quality of MAGs for soil samples via standard metagenomics was a comparison of metagenomics and mini-metagenomics by Alteio et al. (2020), which compared the genome reconstructions obtained from four deep sequenced soil metagenomes and 359 mini-metagenomes. This study, however, used a considerably higher overall sequencing depth of ~ 200 Gbp, averaging at ~ 50 Gb per sample, illustrating that standard metagenomics require extreme sequencing depths to guarantee high-quality genome reconstructions from soil samples, which is not financially feasible in most cases.

Our comparisons also indicate a higher effectiveness (Supplemental Fig. S6) as well as cost efficiency of the midi-metagenomic approach compared with mini-metagenomics, which showed at least nine times higher sequencing costs per MAG (Supplemental Table S6). A major issue with mini-metagenomics is that it targets multiple cells but still relies on WGA, which introduces bias and makes coabundance variation-based binning infeasible, thereby foregoing the main advantages of both SCGs and metagenomics. These conclusions are corroborated by the results of a more thorough mini-metagenomic sampling effort by Alteio et al. (2020), mentioned above. Although the mini-metagenomic MAGs obtained during that study did show lower contamination values than those of standard metagenomics, they also showcased a low overall efficiency of the approach, with few high-quality MAGs being obtained from only 0.8% of the samples, despite a high combined sequencing effort of >190 Gbp (Supplemental Table S7).

The coassembly of midi-metagenomic fractions is therefore not only the most effective but also the most cost-efficient approach compared with current alternatives (Supplemental Tables S6, S7), because sequencing efforts can be distributed across multiple fractions for optimal binning but can still be fully utilized for a

combined assembly. The increased sequencing efficiency may be of particular interest for applications of long-read sequencing technologies, which can provide more coherent sequence context at the cost of lower read throughput (Hu et al. 2021). Furthermore, sequencing efforts do not need to be evenly distributed across midi-metagenomic fractions. Thus, the deep sequencing of an unsorted “main” fraction supplemented by auxiliary sorted fractions of lower sequencing depth can further reduce costs in cases in which the research focus lies on the original community composition and the reconstruction of MAGs is only considered a secondary goal.

Importantly, significantly reduced contamination rates are observed in midi-metagenomic MAGs, which is especially noteworthy considering the growing complaints about reference database contaminations caused by insufficiently screened subquality MAGs and SAGs (Breitwieser et al. 2019; Arkhipova 2020; Vollmers et al. 2022).

We could here show that even a simple sorting setup is sufficient for substantial improvements in both the yield and quality of binned MAGs, as long as partial enrichments or depletions of different community members can be achieved. The fractionation of the sampled community can be done using FACS based on many different cell properties, of which the here-utilized cell size and morphology are only the most simple examples (Sturm et al. 2023). In fact, just the act of FACS sorting itself, independent of applied criteria, already represents a general depletion of large multi-cell aggregates, extracellular DNA, and potential stress susceptible cell types (Wiegand et al. 2021). However, more stringent sorting criteria may further improve the efficiency of the midi-metagenomic approach. Possible criteria could be labeling with fluorescence in situ hybridization (FISH) probes targeting 16S rRNA genes of specific taxonomic groups (Pratscher et al. 2018; Dam et al. 2020) or with function-based enrichments using specific gene or mRNA targeting probes (Kaster and Sobol 2020; Takahashi et al. 2020), radioactive substrate labeling (Lo et al. 2023), fluorescent-labeled antibodies (Müller and Nebe-von-Caron 2010), or sorting based on different autofluorescence spectra caused by species-specific membrane protein compositions (Kang et al. 2020).

However, it needs to be kept in mind that, owing to inherent biases of the sorting process itself, namely, the different enrichment and depletion of specific strains, any definite conclusions on relative abundances within the original community must be based on an unsorted bulk metagenome data set. Such a bulk fraction should therefore generally be included as a reference in the fractionation workflow for each sample (Fig. 2). Because this unsorted bulk metagenome can seamlessly be integrated into the analyses as a dedicated fraction, this does not increase the overall sequencing efforts. Consequently, the total species richness derived by midi-metagenomic assembly can be expected to represent at least the same complexity as a standard bulk metagenome but with the potential addition of strains that are strongly enriched in some of the sorted fractions but sequenced at depths below detection threshold in the bulk metagenome (Figs. 3, 5D; Supplemental Fig. S4; Supplemental Table S8). Consequently, midi-metagenomics may also serve to boost binning efforts in cases in which the variability between samples may turn out not to be sufficient for coabundance-based binning, especially for sample locations that are hard or expensive to access for additional sampling, namely, deep-sea sediments. The exact sorting criteria do not even need to be decided beforehand as a glycerol stock of frozen sample can be revisited for sorting after a preliminary whole-community metagenome analyses.

Although the fractionation process does require access to dedicated equipment such as a FACS- or microfluidic-based cell sorter, which is typically priced at around \$230.00 (€200.00), the MDA-free workflow allows for much more streamlined and cost-effective setups compared with mini-metagenomics: Without the highly contamination sensitive MDA-step, the necessity of clean-room standards is removed, greatly simplifying the required infrastructure as well as expertise and potentially even allowing outsourcing the sorting process to external facilities that already possess adequate sorting devices. Consequently, midi-metagenomics represent a novel and improved technique that is widely accessible to the scientific community and can significantly enhance the quality and reliability of prokaryotic genome reconstruction from environmental samples.

Methods

Microbial samples

To evaluate midi-metagenomic performance compared with metagenomics, soil samples were collected at the Karlsruhe Institute of Technology (KIT), Campus North, Eggenstein-Leopoldshafen, Germany (49°5′48.8″N, 8°25′55.6″E), during four different periods of time: May 25, 2020, October 7, 2020, August 10, 2021, and February 15, 2022. From each sample, several grams were directly frozen at –80°C immediately after collection for subsequent standard metagenome DNA extraction and sequencing.

Five grams of each sample was then prepared for cell sorting by adding 30 mL of filtered, autoclaved, and UV-sterilized phosphate buffer saline (PBS) solution; brief vortexing to disrupt aggregates and dislocate cells attached to debris; and subsequent pelleting and removal of debris by brief centrifugation at 2000g. Sterile glycerol was added to a final concentration of 30% as an antifreezing agent, and the samples were stored at –80°C until further processing. An overview of all samples is given in Table 1.

Fluorescence-activated cell sorting

Prior to sorting, the samples aliquoted for midi-metagenomics were centrifuged for 1 min at 15,871g and 20°C. The supernatant was discarded, and after resuspension of the pellet in 1 mL PBS, 5 µL SYBR green I was added to all samples. The samples were then vortexed, incubated for 20 min at 4°C, and subsequently pelleted again by centrifugation for 1 min at 15,871g. Each pellet was then washed twice with 1 mL PBS.

Before loading the sample into the FACS (FACSMelody, BD Biosciences), an unlabeled negative control was filtered into a 5 mL FACS tube using a sterile SYSMEX CellTrics filter with a 20 µM mesh size and then diluted with PBS. The negative control

Table 1. Overview of samples and fractions

Sample	Sampling date	Fractions produced
Spring20	May 15, 2020	Only unsorted
Autumn20	October 7, 2020	Unsorted, big & small
Summer21	October 8, 2021	Unsorted, BC, MC, SC, BNC & SNC
Winter21/22	February 15, 2022	Unsorted, BC, MC, SC, BNC & SNC
Autumn22	October 7, 2022	Only unsorted
Winter22/23	March 7, 2023	Unsorted, BC, MC, SC, BNC & SNC

Sample location was at KIT, Campus North, 49°5′48.8″N, 8°25′55.6″E. (BC) Big complex, (MC) medium complex, (SC) small complex, (BNC) big noncomplex, (SNC) small noncomplex.

was used to compare the difference of fluorescence signals for a correct gating that included only labeled cells. Subsequently, the same procedure was applied to the SYBR-labeled samples. A threshold was set up in order to disregard smaller particles such as debris during the sorting process, and an excitation wavelength of 488 nM was used.

For samples “summer21,” “winter21/22,” and “winter22/23,” cells were sorted into five different groups via gatings based on plotting fluorescence intensity against the forward scatter signal (FSC) and side scatter signal (SSC), which are roughly proportional to cell size and complexity, respectively (Supplemental Table S1; Supplemental Fig. S1). For sample “autumn20,” only two groups were sorted, according to size measured by differences in FSC (Supplemental Table S1). Configurations for fluidic, optical, and electronic settings were kept constant for all sorting runs, as specified in Supplemental Table S1. No compensation was applied, as only one fluorochrome was used. After sorting, the cells were stored at -80°C until further processing. An overview of the fractions produced per sample is included in Table 1.

DNA extraction

For metagenomics of the unsorted sample, DNA was extracted with the DNeasy PowerSoil kit (Qiagen) following the manufacturer’s instructions. For midi-metagenomic community fractions, DNA was extracted directly from FACS-sorted cell suspensions consisting of 4×10^6 cells. First, the cells were freeze-thawed three times using liquid nitrogen and a 60°C water bath. Then, bead beating was performed three times for 30 sec at 6 m/sec using one tube of lysing matrix for each fraction (MP Biomedicals 6914-800) and an FastPrep-24 homogenizer (MP Biomedicals). Beads and cell debris were pelleted by centrifugation at 14,000g for 5 min, and the supernatant was subjected to standard alcohol precipitation using 1 volume of 80% isopropanol, 0.1 volume 3 M sodium acetate, and 340 μg linear polyacrylamide. After a subsequent wash step with ice-cold 70% ethanol, the resulting DNA pellet was resuspended with 100 μL PCR-grade water followed by further purification via solid-phase reversible immobilization using 1.5 volume of AMPure XP beads (Beckman Coulter) and final elution in 20 μL 1 \times TE. All extracted DNA was immediately stored at -20°C until use.

Polymerase chain reaction for amplicons

Amplicon sequencing was performed using a nested polymerase chain reaction (PCR) approach. Almost full-length PCR products were obtained in a preliminary PCR using 1.25 U OneTaq Quick-Load DNA polymerase (New England Biolabs), 200 μM mixed dNTPs, 500 μM biology-grade bovine serum albumin (BSA; Thermo Fisher Scientific), and 0.2 μM of each universal bacterial forward and reverse primer 27F (5'-AGRGTTYGATYMTGGCTCAG-3') and 1492R(5'-AGRGTTYGATYMTGGCTCAG-3'). PCR products were purified using DNA Clean & Concentrator-5 columns (Zymo Research) according to the manufacturer’s instructions. The purified product was then used as template for a subsequent amplicon PCRs using 0.5 U Q5 high-fidelity DNA polymerase (New England Biolabs), 0.5 U 200 μM dNTP solution mix (New England Biolabs), Q5 high GC enhancer, 0.1 $\mu\text{g}/\mu\text{L}$ BSA (Thermo Fisher Scientific), and 0.2 μM of each universal bacterial primer 341F (5'-AGRGTTYGATYMTGGCTCAG-3') and 518R (5'-AGRGTTYGATYMTGGCTCAG-3'), targeting the V3 hypervariable region.

Sequencing

All libraries were prepared using the NEBNext Ultra II FS DNA library prep kit for Illumina (New England Biolabs), according to the manufacturer’s instructions. Libraries were sequenced on an Illumina NextSeq 550 (New England Biolabs) device using 300 cycles and a paired-end approach.

Read processing and assembly

Reads were quality-trimmed and adapter-clipped using Trimmomatic v.0.36, bbdutk v.35.69, and cutadapt v.1.14 successively (Martin 2011; Bolger et al. 2014; <https://bbmap.org>). Overlapping read pairs were identified and merged using FLASH v.1.2.11 (Magoč and Salzberg 2011). For amplicon data sets, reads were clustered into amplicon sequence variants (ASVs) using the Qiime2 pipeline with dada2 (Callahan et al. 2016; Bolyen et al. 2019). To account for different sequencing depths, the produced ASVs were rarefied to a value of 42,000, which was determined via preliminary rarefaction analysis. Rarefied ASVs were subsequently taxonomically classified using RDP Classifier v1.24 (Wang and Cole 2024) and SINA v1.7.2 (Pruesse et al. 2012). Shotgun data sets were arranged into 81 coassembly subgroups representing all possible subsets of three to six data sets per metagenome or midi-metagenome (Supplemental Table S4). To enable standardized assemblies with 15 Gbp total read input each, shotgun data sets were randomly subsampled down to 2.5, 3.75, 3, and 5 Gbp when possible and coassembled at equal amounts for each coassembly subgroup. Additional assemblies with nonuniform sequencing depth distribution were also performed using unsorted metagenome samples as “main” data sets with 12–13 Gbp sequencing depth and 5 “auxiliary” data sets subsampled to 0.4–0.5 Gbp each. Assemblies were performed using MEGAHIT v1.2.9 (Li et al. 2015).

For 16S rRNA gene diversity analysis in shotgun assemblies, 16S rRNA genes were extracted from all assemblies using barnap (<https://github.com/tseemann/barnap>) and clustered at 99% sequence identity level using VSEARCH v2.21.1 (Rognes et al. 2016), and cluster representatives were aligned using SINA v1.7.2. Alignments were then filtered in order to retain only those that fully overlap a defined sequence window of 600 bp roughly representing the hypervariable regions V3–V6. For beta-diversity analyses, OTU tables were generated based the filtered sequences, and read coverages were determined via coverm v0.7.0 (Aroney et al. 2025). Beta-diversity analyses were then performed using Qiime2, and taxonomic classifications were performed as described above.

MAG reconstruction and analyses

For each coassembly, three different binning tools were used in parallel: Metabat2 v.2.15 (Kang et al. 2019), CONCOCT (Alneberg et al. 2014), and Rosella v.0.4.1 (Newell 2023). For midi-metagenomic approaches, Rosella was substituted for Maxbin (Wu et al. 2016), as Rosella did not function without coabundance information and Maxbin utilizes additional taxonomic and marker-gene criteria that may optimize results for mini-metagenomic and SCG assemblies that lack reliable coverage information (Marine et al. 2014; Kaster and Sobol 2020). Resulting bins were preassessed and filtered using MDMcleaner. Quality categories were then determined based on reassessments using CheckM2 (Chklovski et al. 2023). Taxonomic classifications were based on GTDB-TK v2.1.1 (Chaumeil et al. 2020).

dRep v.3.4.0 (Olm et al. 2017) was employed to identify groups of redundant MAGs created by different assemblies or binning tools and to select the respective most representative MAG.

Similarities between MAGs were additionally determined and visualized based on gene content as previously described elsewhere (Howat et al. 2018).

Data access

All raw sequencing data and processed dereplicated high-quality MAG sequences have been submitted to the NCBI BioProject database (<https://www.ncbi.nlm.nih.gov/bioproject/>) under accession number PRJNA900514. The complete set of processed sequencing data, including redundant and low-quality MAGs with completeness estimates of at least 50% and contamination estimates <25%, has been submitted to Zenodo (<https://doi.org/10.5281/zenodo.13150466>).

Competing interest statement

The authors declare no competing interests.

Acknowledgments

We acknowledge support by the state of Baden-Württemberg through bwHPC. This work was financially supported through the Helmholtz Association program “Materials Systems Engineering” under the topic “Adaptive and Bioinspired Materials Systems” (project ID 43.33.11) and by the German government, through Bundesministerium für Bildung und Forschung (BMBF) project MicroMatrix (project ID 161L0284A).

Author contributions: Study conception and design were by J.V. and A.-K.K. Data collection was by M.C.C. and J.V. Analysis and interpretation of results were by J.V. Manuscript preparation was by J.V., M.C.C., and A.-K.K. Funding was by A.-K.K.

References

- Alneberg J, Bjarnason BS, De Bruijn I, Schirmer M, Quick J, Ijaz UZ, Lahti L, Loman NJ, Andersson AF, Quince C. 2014. Binning metagenomic contigs by coverage and composition. *Nat Methods* **11**: 1144–1146. doi:10.1038/nmeth.3103
- Alteio LV, Schulz F, Seshadri R, Varghese N, Rodriguez-Reillo W, Ryan E, Goudeau D, Eichorst SA, Malmstrom RR, Bowers RM, et al. 2020. Complementary metagenomic approaches improve reconstruction of microbial diversity in a forest soil. *mSystems* **5**: e00768-19. doi:10.1128/mSystems.00768-19
- Arhipova IR. 2020. Metagenome proteins and database contamination. *mSphere* **5**: e00854-20. doi:10.1128/mSphere.00854-20
- Aronoy STN, Newell RJP, Nissen JN, Camargo AP, Tyson GW, Woodcroft BJ. 2025. CoverM: read alignment statistics for metagenomics. *Bioinformatics* **41**: btaf147. doi:10.1093/bioinformatics/btaf147
- Binek A, Rojo D, Godzien J, Rupérez FJ, Nuñez V, Jorge I, Ricote M, Vázquez J, Barbas C. 2019. Flow cytometry has a significant impact on the cellular metabolome. *J Proteome Res* **18**: 169–181. doi:10.1021/acs.jproteome.8b00472
- Blainey PC. 2013. The future is now: single-cell genomics of bacteria and archaea. *FEMS Microbiol Rev* **37**: 407–427. doi:10.1111/1574-6976.12015
- Bodor A, Bounedjoum N, Vincze GE, Erdeiné Kis Á, Laczi K, Bende G, Szilágyi Á, Kovács T, Perei K, Rákhely G. 2020. Challenges of unculturable bacteria: environmental perspectives. *Rev Environ Sci Biotechnol* **19**: 1–22. doi:10.1007/s11157-020-09522-4
- Bolger AM, Lohse M, Usadel B. 2014. Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics* **30**: 2114–2120. doi:10.1093/bioinformatics/btu170
- Bolyen E, Rideout JR, Dillon MR, Bokulich NA, Abnet CC, Al-Ghalith GA, Alexander H, Alm EJ, Arumugam M, Asnicar F, et al. 2019. Reproducible, interactive, scalable and extensible microbiome data science using QIIME 2. *Nat Biotechnol* **37**: 852–857. doi:10.1038/s41587-019-0209-9
- Bowers RM, Kyrpides NC, Stepanauskas R, Harmon-Smith M, Doud D, Reddy TBK, Schulz F, Jarett J, Rivers AR, Eloe-Fadrosh EA, et al. 2017. Minimum information about a single amplified genome (MISAG) and

- a metagenome-assembled genome (MIMAG) of bacteria and archaea. *Nat Biotechnol* **35**: 725–731. doi:10.1038/nbt.3893
- Bowman SK, Simon MD, Deaton AM, Tolstorukov M, Borowsky ML, Kingston RE. 2013. Multiplexed Illumina sequencing libraries from picogram quantities of DNA. *BMC Genomics* **14**: 466. doi:10.1186/1471-2164-14-466
- Breitwieser FP, Perteza M, Zimin AV, Salzberg SL. 2019. Human contamination in bacterial genomes has created thousands of spurious proteins. *Genome Res* **29**: 954–960. doi:10.1101/gr.245373.118
- Callahan BJ, McMurdie PJ, Rosen MJ, Han AW, Johnson AJA, Holmes SP. 2016. DADA2: high-resolution sample inference from Illumina amplicon data. *Nat Methods* **13**: 581–583. doi:10.1038/nmeth.3869
- Chaumeil P-A, Mussig AJ, Hugenholtz P, Parks DH. 2020. GTDB-Tk: a toolkit to classify genomes with the genome taxonomy database. *Bioinformatics* **36**: 1925–1927. doi:10.1093/bioinformatics/btz848
- Chklovski A, Parks DH, Woodcroft BJ, Tyson GW. 2023. CheckM2: a rapid, scalable and accurate tool for assessing microbial genome quality using machine learning. *Nat Methods* **20**: 1203–1212. doi:10.1038/s41592-023-01940-w
- Dam HT, Vollmers J, Sobol MS, Cabezas A, Kaster A-K. 2020. Targeted cell sorting combined with single cell genomics captures low abundant microbial dark matter with higher sensitivity than metagenomics. *Front Microbiol* **11**: 1377. doi:10.3389/fmicb.2020.01377
- Escudero P, Henry CS, Dias RPM. 2022. Functional characterization of prokaryotic dark matter: the road so far and what lies ahead. *Curr Res Microb Sci* **3**: 100159. doi:10.1016/j.crmicr.2022.100159
- Gibbons JD. 2005. Median test, brown–mood. In *Encyclopedia of statistical sciences* (ed. Kotz S, et al.), Vol. 7, pp. 4710–4712. Wiley, Hoboken, NJ. doi:10.1002/0471667196.ess0181.pub2
- Han H, Wang Z, Zhu S. 2025. Benchmarking metagenomic binning tools on real datasets across sequencing platforms and binning modes. *Nat Commun* **16**: 2865. doi:10.1038/s41467-025-57957-6
- Hofmeyr S, Egan R, Georganas E, Copeland AC, Riley R, Clum A, Eloe-Fadrosh E, Roux S, Goltsman E, Buluç A, et al. 2020. Terabase-scale metagenome coassembly with MetaHipMer. *Sci Rep* **10**: 10689. doi:10.1038/s41598-020-67416-5
- Howat AM, Vollmers J, Taubert M, Grob C, Dixon JL, Todd JD, Chen Y, Kaster A-K, Murrell JC. 2018. Comparative genomics and mutational analysis reveals a novel XoxF-utilizing methylotroph in the roseobacter group isolated from the marine environment. *Front Microbiol* **9**: 766. doi:10.3389/fmicb.2018.00766
- Hu T, Chitnis N, Monos D, Dinh A. 2021. Next-generation sequencing technologies: an overview. *Hum Immunol* **82**: 801–811. doi:10.1016/j.humimm.2021.02.012
- Hu X, Liu J, Liang A, Gu H, Liu Z, Jin J, Wang G. 2025. Soil metagenomics reveals reduced tillage improves soil functional profiles of carbon, nitrogen, and phosphorus cycling in bulk and rhizosphere soils. *Agric Ecosyst Environ* **379**: 109371. doi:10.1016/j.agee.2024.109371
- Hutchison CA, Venter JC. 2006. Single-cell genomics. *Nat Biotechnol* **24**: 657–658. doi:10.1038/nbt0606-657
- Jansson JK, Hofmockel KS. 2018. The soil microbiome—from metagenomics to metaproteomics. *Curr Opin Microbiol* **43**: 162–168. doi:10.1016/j.mib.2018.01.013
- Kalisky T, Quake SR. 2011. Single-cell genomics. *Nat Methods* **8**: 311–314. doi:10.1038/nmeth0411-311
- Kang DD, Li F, Kirton E, Thomas A, Egan R, An H, Wang Z. 2019. MetaBAT 2: an adaptive binning algorithm for robust and efficient genome reconstruction from metagenome assemblies. *PeerJ* **7**: e7359. doi:10.7717/peerj.7359
- Kang S-M, de Josselin de Jong E, Higham SM, Hope CK, Kim B-I. 2020. Fluorescence fingerprints of oral bacteria. *J Biophotonics* **13**: e201900190. doi:10.1002/jbip.201900190
- Kaster A-K, Sobol MS. 2020. Microbial single-cell omics: the crux of the matter. *Appl Microbiol Biotechnol* **104**: 8209–8220. doi:10.1007/s00253-020-10844-0
- Lasken RS, Stockwell TB. 2007. Mechanism of chimera formation during the multiple displacement amplification reaction. *BMC Biotechnol* **7**: 19. doi:10.1186/1472-6750-7-19
- Lavigne S, Bossé M, Boulet LP, Laviolette M. 1997. Identification and analysis of eosinophils by flow cytometry using the depolarized side scatter-aponin method. *Cytometry* **29**: 197–203. doi:10.1002/(sici)1097-0320(19971101)29:3<197::aid-cyto2>3.0.co;2-9
- Li D, Liu C-M, Luo R, Sadakane K, Lam T-W. 2015. MEGAHIT: an ultra-fast single-node solution for large and complex metagenomics assembly via succinct de Bruijn graph. *Bioinformatics* **31**: 1674–1676. doi:10.1093/bioinformatics/btv033
- Liu S, Rodriguez JS, Munteanu V, Ronkowski C, Sharma NK, Alser M, Andreae F, Blekhan R, Blaszczyk D, Chikhi R, et al. 2025. Analysis of metagenomic data. *Nat Rev Methods Primer* **5**: 5. doi:10.1038/s43586-024-00376-6

- Lo H-Y, Wink K, Nitz H, Kästner M, Belder D, Müller JA, Kaster A-K. 2023. scMAR-Seq: a novel workflow for targeted single-cell genomics of microorganisms using radioactive labeling. *mSystems* **8**: e0099823. doi:10.1128/mSystems.00998-23
- Lozupone C, Lladser ME, Knights D, Stombaugh J, Knight R. 2011. UniFrac: an effective distance metric for microbial community comparison. *ISME J* **5**: 169–172. doi:10.1038/ismej.2010.133
- Ma B, Lu C, Wang Y, Yu J, Zhao K, Xue R, Ren H, Lv X, Pan R, Zhang J, et al. 2023. A genomic catalogue of soil microbiomes boosts mining of biodiversity and genetic resources. *Nat Commun* **14**: 7318. doi:10.1038/s41467-023-43000-z
- MacFarland TW, Yates JM. 2016. Mann–Whitney *U* test. In *Introduction to nonparametric statistics for the biological sciences using R* (ed. MacFarland TW, Yates JM), pp. 103–132. Springer International Publishing, Cham, Switzerland.
- Magoč T, Salzberg SL. 2011. FLASH: fast length adjustment of short reads to improve genome assemblies. *Bioinformatics* **27**: 2957–2963. doi:10.1093/bioinformatics/btr507
- Marcy Y, Ouverney C, Bik EM, Lösekann T, Ivanova N, Martin HG, Szeto E, Platt D, Hugenholtz P, Relman DA, et al. 2007. Dissecting biological “dark matter” with single-cell genetic analysis of rare and uncultivated TM7 microbes from the human mouth. *Proc Natl Acad Sci* **104**: 11889–11894. doi:10.1073/pnas.0704662104
- Marine R, McCarren C, Vorrasane V, Nasko D, Crowgey E, Polson SW, Wommack KE. 2014. Caught in the middle with multiple displacement amplification: the myth of pooling for avoiding multiple displacement amplification bias in a metagenome. *Microbiome* **2**: 3. doi:10.1186/2049-2618-2-3
- Martin M. 2011. Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnetj* **17**: 10–12. doi:10.14806/ej.17.1.200
- Mattock J, Watson M. 2023. A comparison of single-coverage and multi-coverage metagenomic binning reveals extensive hidden contamination. *Nat Methods* **20**: 1170–1173. doi:10.1038/s41592-023-01934-8
- Mollet M, Godoy-Silva R, Berdugo C, Chalmers JJ. 2008. Computer simulations of the energy dissipation rate in a fluorescence-activated cell sorter: implications to cells. *Biotechnol Bioeng* **100**: 260–272. doi:10.1002/bit.21762
- Müller S, Nebe-von-Caron G. 2010. Functional single-cell analyses: flow cytometry and cell sorting of microbial populations and communities. *FEMS Microbiol Rev* **34**: 554–587. doi:10.1111/j.1574-6976.2010.00214.x
- Nelkner J, Henke C, Lin TW, Pätzold W, Hassa J, Jaenicke S, Grosch R, Pühler A, Sczyrba A, Schlüter A. 2019. Effect of long-term farming practices on agricultural soil microbiome members represented by metagenomically assembled genomes (MAGs) and their predicted plant-beneficial genes. *Genes (Basel)* **10**: 424. doi:10.3390/genes10060424
- Newell R. 2023. “Bioinformatic methods for genome-centric metagenomics.” PhD thesis, Queensland University of Technology, Brisbane, Australia. doi:10.5204/thesis.eprints.237953
- Olm MR, Brown CT, Brooks B, Banfield JF. 2017. dRep: a tool for fast and accurate genomic comparisons that enables improved genome recovery from metagenomes through de-replication. *ISME J* **11**: 2864–2868. doi:10.1038/ismej.2017.126
- Orakov A, Fullam A, Coelho LP, Khedkar S, Szklarczyk D, Mende DR, Schmidt TSB, Bork P. 2021. GUNC: detection of chimerism and contamination in prokaryotic genomes. *Genome Biol* **22**: 178. doi:10.1186/s13059-021-02393-0
- Parkinson NJ, Maslau S, Ferneyhough B, Zhang G, Gregory L, Buck D, Ragoussis J, Ponting CP, Fischer MD. 2012. Preparation of high-quality next-generation sequencing libraries from picogram quantities of target DNA. *Genome Res* **22**: 125–133. doi:10.1101/gr.124016.111
- Pratscher J, Vollmers J, Wiegand S, Dumont MG, Kaster A-K. 2018. Unravelling the identity, metabolic potential and global biogeography of the atmospheric methane-oxidizing upland soil cluster α . *Environ Microbiol* **20**: 1016–1029. doi:10.1111/1462-2920.14036
- Pruesse E, Peplies J, Glöckner FO. 2012. SINA: accurate high-throughput multiple sequence alignment of ribosomal RNA genes. *Bioinformatics* **28**: 1823–1829. doi:10.1093/bioinformatics/bts252
- Qayyum H, Talib MS, Ali A, Kayani MUR. 2025. Evaluating the potential of assembler-binner combinations in recovering low-abundance and strain-resolved genomes from human metagenomes. *Heliyon* **11**: e41938. doi:10.1016/j.heliyon.2025.e41938
- Ribarska T, Bjørnstad PM, Sundaram AYM, Gilfillan GD. 2022. Optimization of enzymatic fragmentation is crucial to maximize genome coverage: a comparison of library preparation methods for Illumina sequencing. *BMC Genomics* **23**: 92. doi:10.1186/s12864-022-08316-y
- Rinke C, Schwientek P, Sczyrba A, Ivanova NN, Anderson IJ, Cheng J-F, Darling A, Malfatti S, Swan BK, Gies EA, et al. 2013. Insights into the phylogeny and coding potential of microbial dark matter. *Nature* **499**: 431–437. doi:10.1038/nature12352
- 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
- Sato MP, Ogura Y, Nakamura K, Nishida R, Gotoh Y, Hayashi M, Hisatsune J, Sugai M, Takehiko I, Hayashi T. 2019. Comparison of the sequencing bias of currently available library preparation kits for Illumina sequencing of bacterial genomes and metagenomes. *DNA Res* **26**: 391–398. doi:10.1093/dnares/dsz017
- Schmeisser C, Steele H, Streit WR. 2007. Metagenomics, biotechnology with non-culturable microbes. *Appl Microbiol Biotechnol* **75**: 955–962. doi:10.1007/s00253-007-0945-5
- Solden L, Lloyd K, Wrighton K. 2016. The bright side of microbial dark matter: lessons learned from the uncultivated majority. *Curr Opin Microbiol* **31**: 217–226. doi:10.1016/j.mib.2016.04.020
- Steen AD, Crits-Christoph A, Carini P, DeAngelis KM, Fierer N, Lloyd KG, Cameron Thrash J. 2019. High proportions of bacteria and archaea across most biomes remain uncultured. *ISME J* **13**: 3126–3130. doi:10.1038/s41396-019-0484-y
- Sturm G, Mojarrad M, Kaster A-K. 2023. Targeted cell labeling and sorting of prokaryotes for cultivation and omics approaches. *Microb Physiol* **33**: 63–84. doi:10.1159/000532088
- Takahashi H, Horio K, Kato S, Kobori T, Watanabe K, Aki T, Nakashimada Y, Okamura Y. 2020. Direct detection of mRNA expression in microbial cells by fluorescence in situ hybridization using RNase H-assisted rolling circle amplification. *Sci Rep* **10**: 9588. doi:10.1038/s41598-020-65864-7
- Tvedte ES, Michalski J, Cheng S, Patkus RS, Tallon LJ, Sadzewicz L, Bruno VM, Silva JC, Rasko DA, Dunning Hotopp JC. 2021. Evaluation of a high-throughput, cost-effective Illumina library preparation kit. *Sci Rep* **11**: 15925. doi:10.1038/s41598-021-94911-0
- Tyson GW, Chapman J, Hugenholtz P, Allen EE, Ram RJ, Richardson PM, Solovyev VV, Rubin EM, Rokhsar DS, Banfield JF. 2004. Community structure and metabolism through reconstruction of microbial genomes from the environment. *Nature* **428**: 37–43. doi:10.1038/nature02340
- Vollmers J, Wiegand S, Kaster A-K. 2017. Comparing and evaluating metagenome assembly tools from a microbiologist’s perspective—not only size matters! *PLoS One* **12**: e0169662. doi:10.1371/journal.pone.0169662
- Vollmers J, Wiegand S, Lenk F, Kaster A-K. 2022. How clear is our current view on microbial dark matter? (Re-)assessing public MAG & SAG datasets with MDMcleaner. *Nucleic Acids Res* **50**: e76. doi:10.1093/nar/gkac294
- Wang Q, Cole JR. 2024. Updated RDP taxonomy and RDP classifier for more accurate taxonomic classification. *Microbiol Resour Announc* **13**: e01063-23. doi:10.1128/mra.01063-23
- Wiegand S, Dam HT, Riba J, Vollmers J, Kaster A-K. 2021. Printing microbial dark matter: using single cell dispensing and genomics to investigate the Patescibacteria/Candidate Phyla Radiation. *Front Microbiol* **12**: 635506. doi:10.3389/fmicb.2021.635506
- Woyke T, Doud DFR, Schulz F. 2017. The trajectory of microbial single-cell sequencing. *Nat Methods* **14**: 1045–1054. doi:10.1038/nmeth.4469
- Wu Y-W, Simmons BA, Singer SW. 2016. MaxBin 2.0: an automated binning algorithm to recover genomes from multiple metagenomic datasets. *Bioinformatics* **32**: 605–607. doi:10.1093/bioinformatics/btv638
- Xu Y, Zhao F. 2018. Single-cell metagenomics: challenges and applications. *Protein Cell* **9**: 501–510. doi:10.1007/s13238-018-0544-5
- Yu FB, Blainey PC, Schulz F, Woyke T, Horowitz MA, Quake SR. 2017. Microfluidic-based mini-metagenomics enables discovery of novel microbial lineages from complex environmental samples. *eLife* **6**: e26580. doi:10.7554/eLife.26580
- Zha Y, Chong H, Yang P, Ning K. 2022. Microbial dark matter: from discovery to applications. *Genomics Proteomics Bioinformatics* **20**: 867–881. doi:10.1016/j.gpb.2022.02.007

Received October 9, 2024; accepted in revised form April 23, 2026.