



Diversity within species: interpreting strains in microbiomes

Thea Van Rossum¹, Pamela Ferretti¹, Oleksandr M. Maistrenko¹ and Peer Bork^{1,2,3,4}

Abstract | Studying within-species variation has traditionally been limited to culturable bacterial isolates and low-resolution microbial community fingerprinting. Metagenomic sequencing and technical advances have enabled culture-free, high-resolution strain and subspecies analyses at high throughput and in complex environments. This holds great scientific promise but has also led to an overwhelming number of methods and terms to describe infraspecific variation. This Review aims to clarify these advances by focusing on the diversity within bacterial and archaeal species in the context of microbiomics. We cover foundational microevolutionary concepts relevant to population genetics and summarize how within-species variation can be studied and stratified directly within microbial communities with a focus on metagenomics. Finally, we describe how common applications of within-species variation can be achieved using metagenomic data. We aim to guide the selection of appropriate terms and analytical approaches to facilitate researchers in benefiting from the increasing availability of large, high-resolution microbiome genetic sequencing data.

Conspecific

Belonging to the same species; for example, conspecific strains are strains that belong to the same species.

Metagenomics

The study of all genomes present in a sample from a microbial community. Often performed as shotgun metagenomics, in which extracted DNA is fragmented before sequencing.

¹European Molecular Biology Laboratory, Structural and Computational Biology Unit, Heidelberg, Germany.

²Max Delbrück Centre for Molecular Medicine, Berlin, Germany.

³Molecular Medicine Partnership Unit, University of Heidelberg and European Molecular Biology Laboratory, Heidelberg, Germany.

⁴Department of Bioinformatics, Biocenter, University of Würzburg, Würzburg, Germany.

e-mail: bork@embl.de

<https://doi.org/10.1038/s41579-020-0368-1>

For over a century, bacterial cultivation has enabled the isolation and classification of thousands of bacterial strains. Through these efforts, a species concept was translated in the bacterial context as a group of individuals who form a coherent genomic cluster¹ (see below for details and disagreements). Despite this genetic similarity, it was also established that a large magnitude of phenotypic variance is possible among strains from the same species (conspecific strains). The importance of variability within species has been particularly well studied in the context of pathogenicity, and many species have been found to have both pathogenic and commensal strains (for example, *Escherichia coli*² and *Bacteroides fragilis*³). Indeed, a classic example are *E. coli* strains, which can be pathogenic, commensal, host associated or environmental³. The relationship between strain identity and host health demonstrates how it can be insufficient to study microbial communities at species level resolution, and the same applies in many other areas such as drug response⁴, nutrient cycling⁵, nitrogen fixation⁶ and host association⁷.

Cultivation-based approaches have a fundamental⁸ and continued⁹ role in studying within-species variation but, despite their recent methodological progress¹⁰, they have important limitations. Few microorganisms can be easily cultivated under isolated, laboratory conditions, and cultivation is typically low throughput. Even when culturing is possible, organisms are then studied in isolation and not in their natural community setting. Culture-free, strain-level analysis of entire microbiomes

has been possible for over 15 years^{11–16}, but it has been limited due to shallow read depths and small sample sizes. Following the recent technological and algorithmic innovations in metagenomics (BOX 1) and the decreasing cost of sequencing, large-scale metagenomic analyses of variation within species have become feasible (BOX 2). There is great promise in these approaches^{17–20} and they have vastly increased the rate of discovery, but they are also leading to scientific and semantic challenges.

In the traditional cultivation approach, ‘strain’ refers to a pure culture or isolate, denoting a taxonomic entity rather than a natural concept²¹. This operational definition cannot be transferred directly to the modern culture-free approaches, and a widely accepted, biologically meaningful definition of strain remains elusive. Exacerbating this situation, and perhaps in response to the lack of generally accepted terminology, a plethora of overlapping terms have been coined in high-resolution microbiome studies, often being poorly defined. The resulting confusion impedes communication and synergy among researchers both in microbiome fields and beyond. To place new operational definitions in the correct context of existing conceptual definitions of within-species variation, it is essential to understand the microevolutionary processes that create and constrain variation within species.

In this Review, we summarize the processes that produce and constrain variation within species and describe how the balance of these forces shapes the magnitude and structure of the variation. We provide an overview

Box 1 | Molecular approaches to characterize variation within species

A wide range of methods are available for studying within-species variation, either based on cultured isolates or directly in microbial communities.

Microbiome-based methods are less established but are not limited to culturable microbiota. Foundational community-fingerprinting methods like DGGE, TRFLP and ARISA^{181,182} enabled some species to be studied at high resolution without culturing. Owing to their low-throughput and limited resolution, these methods have largely been superseded by genetic sequencing approaches. Despite its origin as a low-resolution method, 16S rRNA gene amplicon analysis can sometimes now differentiate within some species using Oligotyping^{183,184}, amplicon sequence variants (ASVs)^{185–187} and single-nucleotide variants in full gene sequences¹⁸⁸. However, 16S rRNA approaches remain extremely limited in resolution for within-species analysis and can be confounded by multiple, non-identical copies of the 16S rRNA gene per genome¹⁸⁸.

Shotgun metagenomic sequencing provides more information by considering more marker genes or whole genomes. Many tools have been developed to analyse metagenomic data to describe variation within species^{19,20}. The major approaches include single-nucleotide variant-based profiling, either within predefined marker genes^{56,91,99,106} or across whole species-reference genomes^{104,105,120}, overall similarity to strain-reference genomes^{93,94}, sequence typing¹¹⁷ and gene content-based profiling⁹⁶. Metagenome-assembled genomes can be recovered by binning and assembling co-abundant genes¹⁸⁹; however, these come with important limitations (BOX 3).

Non-microbiomic but culture-free methods include microfluidics-based techniques that enable organism-specific enrichment prior to sequencing^{190,191} and single-cell sequencing, which produces single amplified genomes¹⁹². Culturing is becoming possible for a growing number of bacteria owing to methodological advances such as culturomics, which combines the use of multiple culture conditions with rapid bacterial identification¹⁰.

Non-genomic approaches, such as cryo-electron microscopy-based imaging and transcriptome-based, proteome-based and metabolome-based profiling methods, can capture phenotypic differences within species and can be used both separately and in conjunction with genomic approaches. These methods range from well established, such as serotyping and functional profiling, to more recent and high throughput such as thermal proteome profiling¹⁴¹.

of the major ways in which this variation can be studied and stratified into categories using metagenomic data and define the commonly used terminology, which we put into the context of applications. We use ‘within-species variant’ to refer to any grouping below the species level. Throughout this Review, we highlight the advances and challenges that are resulting from the use of metagenomic data to study within-species diversity.

Variation and cohesion within species

Processes leading to within-species variation. Diversity within species is the result of continuous processes of variation generation and subsequent selection and drift (FIG. 1). Mutations and gene flow introduce genetic variability into otherwise identical lineages of clonal daughter cells.

Mutations (that is, substitutions, insertions, deletions and inversions) arise continuously in the genome owing to errors in the DNA replication process, damages caused by mutagens, or errors in the DNA repair and recombination mechanisms²². Although the typical mutation rate for double-helix DNA-based organisms is approximately 1 nucleotide change per 10⁹ nucleotides per replication²³, mutation rates can vary across and within species by orders of magnitude²⁴. Selection for lower or higher rates balances the metabolic cost of reducing mutation frequency versus the impact of deleterious mutations²⁵. The direction of this balancing depends on habitat conditions, population size and mutator allele strength²⁵. The rate of accumulation of mutations within a lineage

of bacteria depends on the mutation rate as well as on natural selection and genetic drift, which act upon the mutations. This further diversifies the observed rates of mutation. For example, non-lethal rates of mutation from 10^{–9} to 10^{–3} mutations per genome per generation have been observed in *Vibrio* species^{26,27}. Further, not all portions of the bacterial genome are equally subject to mutations. Mutation accumulation rates are higher in accessory genes than in core genes, unless a core gene is located near accessory genes or mobile genetic elements, and higher in secondary chromosomes than in primary chromosomes^{28,29}. In general, deletions are more frequent than insertions, and non-functional sequences are readily lost from bacterial genomes^{30,31}. Mutations that arise in one genome can be passed vertically to descendants or horizontally to neighbouring cells.

The transfer of genetic variation from one population to another (gene flow) can cause rapid and large-scale additions and rearrangements of genomic regions³². DNA can be transferred between cells by horizontal gene transfer (HGT) via transformation, transduction, conjugation, gene transfer agents and membrane vesicles^{33,34}. Newly acquired donor DNA can stay separate within the acceptor cell (for example, as a plasmid or lytic phage) or can be incorporated into the genome of the acceptor through a number of mechanisms³⁴, including homologous recombination³⁵. HGT is more frequent within species, but it can also occur between species³⁶. HGT can result in the replacement of genetic segments with donor homologues, often within species via homologous recombination, or in the acquisition of new genetic material. In terms of impact on within-species variation, the most important factor of HGT is not the mechanism (for example, homologous recombination) but rather whether or not the genetic material being transferred is novel to the recipient population or species (discussed below). The main processes limiting HGT include a lack of surface compatibility for the conjugative process, CRISPR-mediated microbial immunity³⁷ and restricted host specificity of bacteriophages³⁴. Notorious examples of HGT between conspecific variants include two cases where toxin genes were transferred from toxigenic to non-toxigenic strains in *Clostridioides difficile*³⁸ and in *E. coli*³⁹, with the latter causing 54 deaths in 2011 in Germany.

Natural selection and genetic drift determine the fate of within-species variation introduced through mutation and gene flow. Genetic drift randomly eliminates genetic variations within a population, whereas natural selection maintains or eliminates variations that respectively confer a fitness advantage or disadvantage. In this context, the effect of natural selection is limited by the background noise of genetic drift⁴⁰. Natural selection is driven by a multitude of biotic and abiotic factors that differentially influence the survival and replicative capability of species subpopulations (FIG. 1). These factors can shape the composition of microbial communities at the species and within-species levels through community assembly⁴¹ and classic evolutionary forces. Selective pressure factors vary from habitat to habitat and can include pH, temperature, the concentration of oxygen and other gases, nutrient availability, direct competition or commensalism with other bacteria, predation by phages and

Population

A set of individuals who occupy a particular spatial area.

Mutator allele

Genetic variation (allele) that results in an increased mutation rate.

Genetic drift

Change of allele frequencies in a population caused by stochastic factors.

Horizontal gene transfer

(HGT). The movement of genetic information between organisms, in contrast to vertical gene transfer from parent to offspring.

Homologous recombination

(HR). Type of genetic recombination in which genetic material is exchanged between two similar or identical regions of DNA.

Marker genes

In microbiome context: genes or genetic segments, the presence or specific DNA sequence of which is distinctive of a category of interest such as a species or clade.

eukaryotes, and the presence of stress-inducing xenobiotics such as drugs, antimicrobial compounds and heavy metals.

Species definitions and mechanisms of species cohesion.

Through the vertical accumulation of mutations and the horizontal acquisition of genes, variation among the descendants of one cell could constantly increase, creating a continuous landscape of genetic variation across bacterial genomes. However, when genomic similarities are compared across bacteria, distinct clusters are observed. These clusters are thought of as species in bacteria⁴², though the applicability of a 'species' concept is contested⁴³. In this Review, we use the word 'species' to reflect these clusters of genetic similarity.

For many decades, bacterial species delineation based on genome similarity has been measured using DNA–DNA hybridization (DDH). According to the

bacterial nomenclature code, conspecific genomes have $\geq 70\%$ similarity by DDH. Increasingly, DDH is complemented or replaced by DNA sequencing of isolates and average nucleotide identity (ANI) comparisons^{8,44}, with approximately $\geq 70\%$ similarity in DDH corresponding to $\geq 94\%$ ANI in the core genome and to $\geq 96\%$ ANI in universal marker genes^{7,45–49}. The approximation in these correspondences can affect classification, as in the case of *Fusobacterium nucleatum*, for which subspecies were defined based on DDH⁵⁰ but were then suggested to be reclassified as separate species after reassessment with in silico measurements of ANI⁵¹. As suggested by early studies^{13,52,53}, the presence of a distinctive bacterial species boundary is identifiable using metagenomic data. This was recently confirmed by large-scale studies identifying this boundary at ANI thresholds based on whole genomes ($\sim 95\%$)^{54,55} and on marker genes (96.5%)^{48,56} as well as describing a drastic drop in gene flow in core genomes.

Box 2 | Culturing isolates versus metagenomics for analysis of variation within species

Traditionally, investigations below the species level have relied on studying cultured isolates. With the rise of metagenomics, the amount of high-resolution genetic data has increased. Generally, these data are analysed based on variation within specific genetic segments (for example, marker genes) or within genomes recovered through assembly (metagenome-assembled genomes; MAGs) (BOX 1). Although this enables unprecedented discoveries owing to the large scale of data produced, these new methods also have important limitations and introduce new complexity (see the table below). Although metagenomics provides important new benefits over studying isolates, the two methods remain complementary^{8,193}. To ensure future synergy between the two approaches, isolate genome and MAG data quality must be readily available and comparable, and a common vocabulary should be maintained.

Criteria	Culturing isolates	Metagenomic sequencing
Scope of microorganisms that can be studied below the species level	Must be culturable in isolation but can be of low abundance in original sample	Must be abundant or deeply sequenced
Ability to describe multiple species variants within one sample (of the same or of different species)	Requires multiple rounds of isolation Intractable for low abundance variants	Can be determined from sequencing data from one sample but sufficient sequencing depth is required to distinguish from sequencing error
Ability to determine whether genetic variants originate from the same organism (genetic linkage)	Possible (as long as variation within isolate colony is low, which is normally the case)	Very difficult or impossible in current typical approaches but improvements are possible (for example, long reads, time series data and Hi-C sequencing ¹⁹⁴)
Ability to put a within-species variant in context of its community	Limited and work intensive	Implicitly supported, though biases exist ^{147,161,195}
Ability to describe phenotypic differences between within-species variants	Heterogeneity can be assessed ¹³⁴ with clinical, environmental and industrial relevance	Limited to description of potential phenotypes
Support for follow-up study	Isolates can be directly experimented on (for example, response to drug exposures)	Extracted DNA can be further tested molecularly (for example, PCR)
Main method for genome recovery	Isolate shotgun genomic sequencing and assembly	Shotgun metagenomic DNA sequencing followed by assembly (MAG; BOX 1)
Quality of the recovered genomes	Often remain at draft level but usually are high quality with little contamination	May have higher error rates and be chimeric, contaminated and incomplete ^{192,196}
Quality assessment of the recovered genomes	Provided by central repositories, with various guidelines developed (see, for example, REF. ⁴²)	Routinely assessed but ad hoc by authors. Recommendations are emerging ¹⁹²
Determining the presence or absence of a gene in the recovered genome	Usually simple and correct	Difficult to be certain
Expected impact of long read sequencing	Longer contigs, less challenged by repetitive regions	Better genomes for the most abundant organisms, low abundance fraction still hard to access

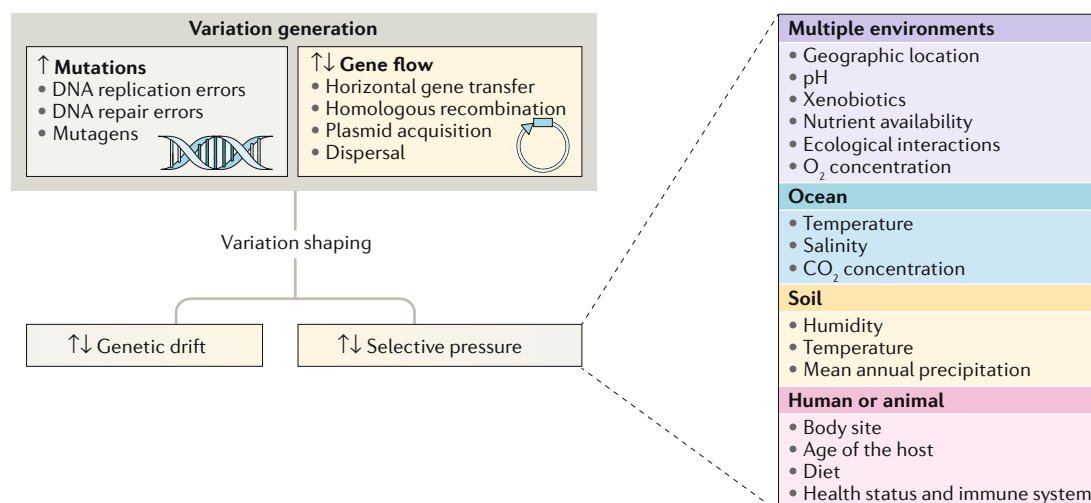


Fig. 1 | Drivers of variability within bacterial species. Within-species variability is introduced by mutations, which usually increase the amount of variation within a species (up arrow), and gene flow mechanisms, which can increase or decrease the amount of variation within a species. This variability is shaped by genetic drift and selective pressure, which can also increase or decrease the amount of variation. Selective pressures are shaped by many biotic and abiotic factors, some of which are known to drive adaptation in particular habitats more than in others.

Despite the overall consistency of genomic ANI data, defining bacterial and archaeal species remains controversial, with over 20 conceptual definitions of ‘species’^{57–60} and some researchers questioning the concept altogether⁴³. The biological and the phylogenetic concepts of species are the most applicable for bacteria and archaea⁶¹. The former defines species as a group of individuals that can interbreed, resulting in viable offspring, which translates to the possibility for homologous recombination in Prokaryotes, whereas the latter defines species as clades that are characterized by distinctive phenotypic properties. Both concepts predict a decline in the rates of homologous recombination^{36,62} and HGT⁶³ between different species. The multitude of potential species definitions are not necessarily well served by ANI-based genomic comparisons alone. Instead, other methods can be used to operationally define species in addition to or in place of ANI such as by phenotype, similarity in a universal single copy gene (for example, 16S rRNA) and gene content^{46,64}.

The genomic similarity within species is called ‘cohesion’. This is maintained predominantly through within-species recombination and selection against lower-fitness alleles^{55,65}. If an allele is more beneficial than all others in a population, it can spread completely through that population, resulting in a hard selective sweep^{66,67}. When recombination rates are low, it is likely that the whole genome will hitchhike to prevalence along with this adaptive allele, resulting in a genome-wide selective sweep⁶⁸. When hard, whole-genome selective sweeps occur, they can reduce diversity within a species and maintain dissimilarity between species^{65,69,70}.

Determinants of magnitude and structure of variation within a species. Diversity within species is generated, maintained and purged to different extents, such that some species are highly heterogenous whereas others are

tightly cohesive. These features of within-species variation depend on the populations observed (BOX 3) and can be described globally or locally. The balance between the forces that increase diversity and those that maintain cohesion shapes both the magnitude and the structure of variation within a species.

The amount of variation generated within a species depends on the mutation rate, generation time, tendency for inter-species HGT and population size, whereas the amount of variation that persists depends on the stringency of selective pressures in its habitats, population size⁷¹ and the frequency and severity of selective sweeps. The balance between divergence and cohesion is modulated by selection and drift, which are shaped by the biotic and abiotic factors of the ecological niche (FIG. 1). HGT can increase the genetic variation within a population if the material being transferred is novel to the receiving population, for example, if the donor cell was dispersed from a foreign population or is distantly related. Conversely, HGT can homogenize a population in terms of specific gene content or single-nucleotide variant (SNV) presence if it spreads this genetic material throughout the population, resulting in a gene-specific hard selective sweep⁷².

Within a species, a structured population can arise owing to a combination of soft selective sweeps — when multiple alternative adaptive alleles spread and coexist in a population⁷³ — along with drift and dispersal into new locations with similar or new ecological niches. For instance, when the rate of mutation generation is high and the rate of within-species recombination is low, strains may diverge into subgroups that are more internally cohesive relative to one another. Specifically, a reduction of the ratio of recombination to mutation events below 0.25 seems to enable subpopulations to diverge freely^{36,74}. This may result in the establishment of subspecies^{75,76}, which are groups of strains with a partially disrupted gene flow that might be in the process of speciation.

Selective sweep

A reduction of the genetic variation in a population owing to selection acting on novel mutations or existing alleles.

Hard selective sweep

One beneficial allele at a locus replaces most other alleles in the population.

Soft selective sweeps

Multiple beneficial alleles at a locus gain prevalence, replacing standing genetic variation in the population.

Intraspecific

Below species level, that is, at a higher resolution than species.

Metagenome-assembled genomes

(MAGs). Genome sequences recovered from metagenomic data, usually fragmented, and potentially incomplete or contaminated. Typically, shotgun metagenomic sequencing produces short DNA sequences that are then assembled and binned into 'genomes' using k-mer frequencies and abundance information.

Subspeciation can be caused or accelerated by physical or geographic barriers that block gene flow between subspeciating groups (allopatric), which leads to the divergence of subspecies either owing to natural selection or drift⁷⁷. However, subspeciation can also occur without spatial separation (sympatric). In this case, it is likely that there is a selective advantage to specialization, for example, to diminish competition for resources. Owing to the extreme dispersibility of bacteria and archaea, complete physical blocks to gene flow may be rare, and in-between scenarios may be possible. When occasional gene flow occurs and niches overlap, purifying the selection can maintain partial cohesion between subspecies, which can prevent divergence from establishing stable subspecies⁷⁸.

At one extreme, species can be monotypic; that is, they have a uniform or 'smeared' distribution of genetic similarities across their entire population. Monotypic species with low diversity are more likely to be specialists, with narrow geographic distributions or host ranges, or are the product of recent speciation^{79,80}. *Chlamydia trachomatis* is an example of a monotypic low diversity intracellular pathogenic species⁸¹. At the other extreme, species with subspecies (polytypic) and high diversity are more likely to be free-living generalists with multiple adaptations to distinct and fluctuating environments as well as broad geographic ranges or many partially overlapping niches^{78,80}. For example, *E. coli* has at least six phylogroups that tend to be more prevalent in different habitats⁸².

Much of the fundamental knowledge described above was obtained on a species-by-species basis through

culture-based and isolation-based experiments. The rise of microbiomic approaches enables the characterization of variation across many species on a large scale and offers promising new research avenues (BOX 2). To meaningfully place these new findings into context it is important to appropriately adapt the concepts and terminology from this body of knowledge for use in metagenomic studies.

Stratification of within-species variation

Within-species variation often needs to be stratified into meaningful groups to be studied and associated with categorical variables such as health status, geographic location or metabolic capability. The theory described above can support the conceptual definitions of such groups but, generally, these cannot be used directly in microbiological studies. Instead, operational definitions of variant groups must be devised based on criteria that can be measured. Typically, this is done on genetic or phenotypic scales. The appropriate metrics to use to operationally define variant groups, such as strains, depends on the biological questions being asked and the methodology being used (FIG. 2a).

Genetic stratification using metagenomic data. Within-species genetic variation can be measured in many ways, with some common metrics being overall genome similarity, the number of shared and unique genes, and/or the number and nature of SNVs. In this section, we discuss how these measures are taken and explore their strengths and limitations. When these analytical approaches are applied to the large amount of data produced by metagenomic sequencing, within-species profiling can be simultaneously performed in a high-throughput manner for many species (see, for example, REFS^{76,83–91} and examples below). However, this also raises various data quality issues, such as incomplete and partially erroneous data, as well as technical challenges such as large computational and storage requirements.

The overall similarity between conspecific genomes at intraspecific levels can be assessed from metagenomic data either directly from reads and reference genomes^{92–94} or through comparisons of metagenome-assembled genomes (MAGs)⁵⁴. Reference genome-based approaches can be limited by the low availability of appropriate reference genomes, especially in non-human microbiomes. Large sets of MAGs are now available, and methods to calculate ANI have improved in efficiency⁹⁵. However, calculating ANI for large genomic cohorts remains computationally challenging⁵⁴. Further, using MAGs in ANI comparisons can introduce inaccuracies owing to data quality limitations and incompleteness (BOX 3).

A decline in ANI and recombination rate can be an indicator of ongoing subdivision of a species⁵⁷. However, in contrast to species boundaries, within-species variants do not seem to display a universal threshold based on genome or marker genes that would categorize them into groups. Instead, the range and distribution of ANI values within species vary by taxon and population⁵⁴, which limits its utility for broad stratification. Further, genetic differences that are coded by a small number of

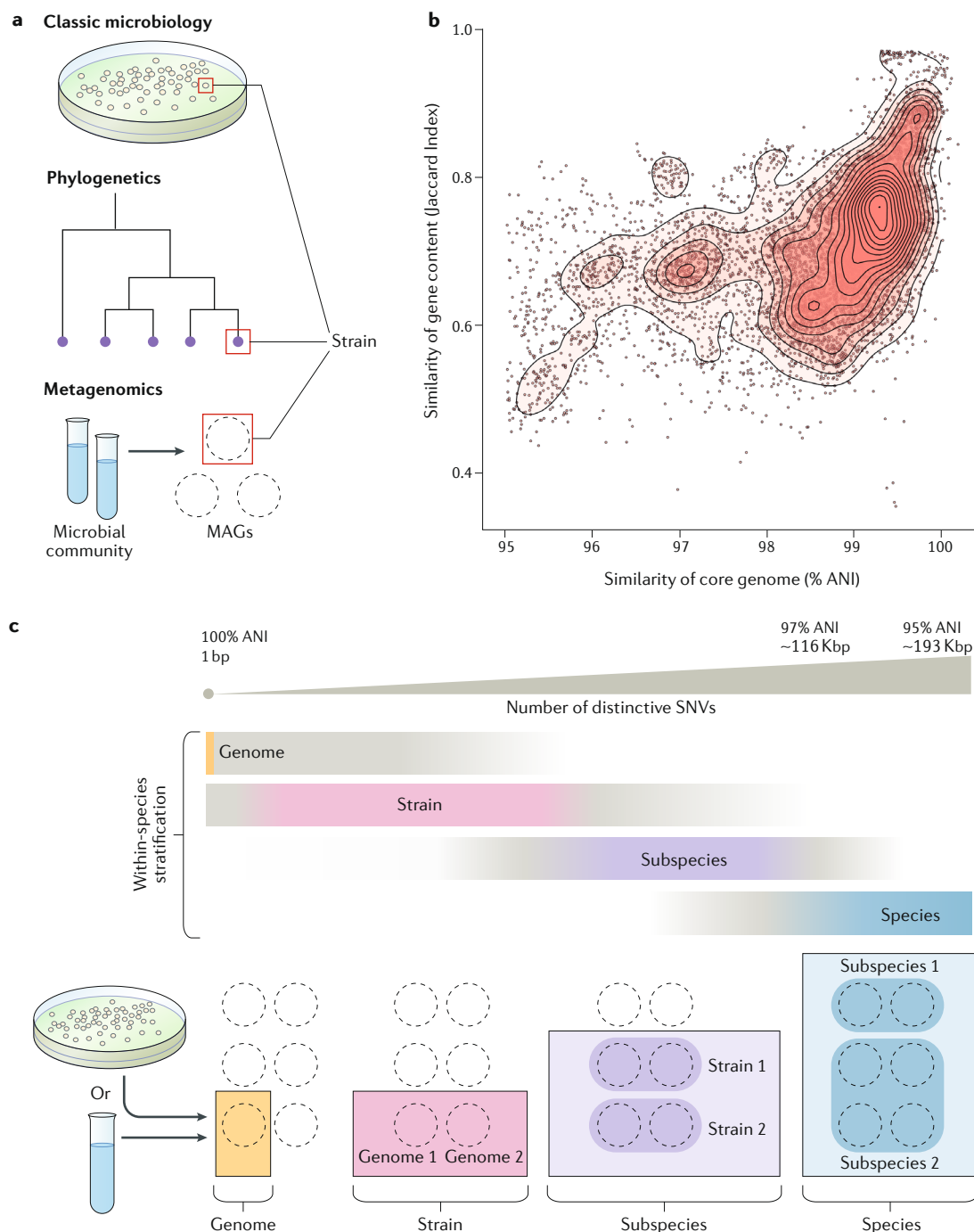
Box 3 | Challenges in studying variation within species in metagenomics

Investigations of variation within species in microbial communities are faced with study-design, technical and methodological challenges. A main study-design challenge is the 'unobserved variation' paradigm: you do not see what you do not sample. If low variability is seen within a species, it is difficult to prove that it is not due to under-sampling or sampling bias. This bias can be temporal (for example, due to strain turnover or extinctions) or spatial (for example, due to proximate sampling areas harbouring substantially different intraspecific profiles, such as in soil or skin). Shallow sequencing depth also biases against observing low abundance within-species variants. These biases are mitigated by the increasing number of deeply sequenced metagenomic samples. However, the integration of these samples across studies is still faced with technical challenges that are well known in metagenomics^{195,197–199}.

Although undoubtedly useful for investigating unknown and under-represented species, metagenome-assembled genomes (MAGs) have important limitations. MAGs are population consensus genomes; thus, loci may be polyallelic and unlinked^{164,196}. When compared to isolate genomes, MAGs often have low assembly quality, are less complete and are more likely chimeric^{104,164,192,196,200}. Therefore, and owing to the difficulty in detecting chimaeras below the species level, MAGs should not be considered equivalent to genomes sequenced from isolates¹⁹⁶. The use of the term 'complete MAG' should be adopted only for MAGs that are analogous to complete isolate genomes, which are usually a single circular contig with no gaps.

To avoid confusing isolate genomes and MAGs, the growing practice of uploading MAGs to public genome databases¹⁹⁶ should be discouraged, and the phrase 'genome-resolved metagenomics' should not be used for MAG studies that do not directly assess heterogeneity within MAGs. Single-cell sequencing approaches provide a promising alternative to MAGs for recovering genomes from metagenomes but are limited by high cost, low throughput, potential contamination and quality issues due to using a single molecule of DNA²⁰¹.

Continued technical advances, decreasing sequencing costs and an increasing integration of complementary methodologies will be necessary to counteract these challenges in data generation and integration.



nucleotides relative to the size of the genome, and thus have a small impact on ANI, can have a very large impact on phenotype. Therefore, at the small scale of ANI differences that occur within species, measures of gene content, SNVs and indels are more informative than ANI for defining biologically relevant within-species variants.

Gene content is the sum of all genes in a genome, including core genes (which are present in almost all conspecific variants) and accessory genes (which are only present in a subset). Differences in accessory gene content between variants can arise at the single-gene level¹⁹⁶ or at the genetic-segment level⁸³, which can include multiple genes (structural variation). Gene content differences can

be calculated based either on the presence or absence of a gene⁹⁷ or on the number of copies of that gene⁹⁸. Gene order (synteny) is considered within structural variation but has not yet been addressed directly by metagenomic methods. Metagenomic data can be used to study within-species gene content variation by looking for gene clusters⁹⁶ or by associating gene content with variants defined by SNV profiles^{76,99}. The relationship between gene content similarity and phylogeny is complicated by HGT. However, comparative studies of conspecific genomes have shown that pairwise similarity based on gene content is correlated with pairwise similarities based on core genome ANI^{100,101} (FIG. 2b), and that distinct SNV profiles can correspond to distinct gene profiles⁷⁶.

◀ Fig. 2 | **Within-species stratification.** **a** | Different operational definitions of 'strain', based on the field of investigation: a cultured isolate in classic microbiology, a leaf node in a phylogenetic tree and a metagenome-assembled genome (MAG) in metagenomics. **b** | Each point is a pairwise-comparison of one isolate genome versus all other conspecific isolate genomes. The data¹⁰⁰ are from 155 bacterial species, each with at least ten sequenced isolate genomes. The opacity of the red-coloured topographical overlay indicates the density of points. The plot shows the relationship between the similarity of the core genome, measured by average nucleotide identity (ANI), versus the similarity of gene content, measured by Jaccard Index. Genomes with higher similarity between their core gene sequences tend to have more genes in common (Spearman correlation $R=0.57$, $P<2.2\times10^{-16}$). However, a high ANI does not necessarily imply a highly similar gene content, with many genomes with an over 99% core genome ANI having less than 70% of genes in common. Most within-species ANI values are greater than 97%; the few data points below 95% ANI are not shown (83% and 4% of data points, respectively). The data are adapted from REF.¹⁰⁰. **c** | Spatial distribution of key terminology used to stratify variation within bacterial species, ranging from a single nucleotide variant (SNV) in the whole genome to the species-level threshold (97% ANI). The coloured portions of the bars reflect the recommended scope of use for each term, and the grey portions indicate the common, often unspecific, scope of use. Broadly speaking, conspecific genomes have identical nucleotides at homologous positions across 97% of their genome (97% ANI), which corresponds to differences in the order of 116,000 SNVs based on an average bacterial genome size (3.87 Mb (REF.¹⁸⁰)). The bottom panel illustrates the hierarchy of these terms, with a species potentially containing multiple subspecies, a subspecies containing multiple strains and a strain containing multiple (non-identical) genomes. These genomes can be sequenced from cultured isolates or through assembly of a metagenomic sample, creating a MAG that represents the consensus genome of a population of cells.

SNV differences can be used to compare conspecific variants at high resolution. These comparisons can consider the number of variant positions, their locations (for example, in core genes, accessory genes or intergenic regions), their spread across the genome (clustered or disperse) and their potential phenotypic impact (for example, synonymous or non-synonymous mutations). In metagenomes, the identification of SNVs can be *de novo*^{99,102}, based on MAGs¹⁰³, or based on pre-existing reference genes or genomes^{104–106}. The degree of similarity between the references and the actual community members can have major impact on the accuracy of the results¹⁰⁷. Identifying SNVs based on MAGs can reveal population dynamics, such as hard and soft selective sweeps in populations of lake bacteria¹⁰³, but can also introduce errors owing to the potential low quality of MAG references (BOX 3). Groups of conspecific genomes can be defined from metagenomic data based on the distinctive presence of SNVs (for example, SNP types¹⁰⁸) — from thousands of SNVs indicating population structure by defining subspecies⁷⁶ and subpopulations¹⁰⁹, to tens of SNVs delimiting strains¹⁰⁸. Isolate data have been used to show that single SNV differences can determine phenotype such as pathogenicity^{110,111} or antimicrobial drug resistance^{112,113}. The ability to detect low abundance SNVs in microbiomic data is limited when sequencing depths are shallow and population sizes are large. When SNVs are likely to have been vertically transferred, then they can be used to define haplotypes and lineages. Extending this approach, SNVs can be used to reconstruct phylogeny within a species¹¹⁴; however, care must be taken to use loci that are unlikely to have been in an HGT region such as housekeeping genes¹¹⁵.

When multiple genetic variants are in one chromosome they are 'linked'. Linked variants are inherited together, but this linkage can be disrupted by

recombination or mutation. Determining the linkage between alleles can be used to track lineages, reconstruct haplotypes (phasing variants) and detect HGT. However, metagenomic data are inherently limited in providing linkage data when the typical short-read, shotgun sequencing approach is used because this method breaks up DNA. The assembly of short reads may be able to recover linkage information; however, chimerism is common when there are multiple, highly similar genomes within one sample such as multiple conspecific strains (strain heterogeneity). Instead of providing exact profiles of linked alleles, shotgun metagenomics is usually limited to providing sets of multiallelic loci with allele frequency information; these can still be useful for many applications, as described in the final section of this Review. They can also be used to perform population genetic analyses for a species such as to calculate estimates of population diversity (for example, π , the diversity or average pairwise genetic difference between individuals), population structure (for example, fixation index (F_{st}) or allele similarity between populations) and selection pressure (for example, the ratio of non-synonymous and synonymous substitutions (dN/dS), the ratio of non-synonymous and synonymous polymorphisms (pN/pS), Tajima's D, or Fay and Wu's H)¹¹⁶.

Many software tools have been developed to measure and categorize diversity within species using metagenomic data. Generally, these have two broad aims: classification and discovery. Classification-oriented tools, for example, metaMLST¹¹⁷, PathoScope⁹⁴, MetaPhlAn2 (REF.¹¹⁸), StrainSifter¹¹⁹, Sigma⁹³, SPARSE⁹² and StrainEst¹²⁰, aim to detect if a known, characterized, within-species group (for example, a target genome, named strain, classic typed subspecies or multi-locus sequence type) is present in a sample. Discovery-oriented tools typically group within-species variation into clusters of similarity using one of three measures: gene content (for example, PanPhlAn⁹⁶), SNVs in whole or core genomes (for example, metaSNV¹⁰⁵) or SNVs in marker genes (for example, Lineages algorithm¹²¹, ConStrains¹⁰⁸, StrainPhlAn¹⁰⁶, DESMAN⁹⁹, StrainFinder¹²² and mOTUs2 (REF.⁵⁶)), which might be followed up with the detection of distinctive gene content (for example, DESMAN⁹⁹). Although many tools claim to provide strain level resolution, the term 'strain' is defined differently across software (see next section for discussion of definitions). The tools that can recover SNV linkage information *de novo* from SNV abundances across samples include ConStrains¹⁰⁸, DESMAN⁹⁹, StrainFinder¹²² and the Lineages algorithm¹²¹. When the assumption can be made that samples contain a single dominant within-species group, tools like StrainPhlAn¹⁰⁶ and metaSNV¹⁰⁵ can also be used to cluster SNVs into within-species groups (strains and subspecies, respectively).

Although these tools enable many applications of metagenomic data to study within-species variation (see below), they have some important limitations. For example, tools that rely on mapping reads to reference genomes or marker genes are inherently limited by the availability of appropriate reference genomes, which in some environments is very low (for example, freshwater

and soil). This limitation can be circumvented by building and using MAGs (for example, as in DESMAN), but MAG quality concerns must be considered, especially if time series data are not available (BOX 3). Other logistical limitations include requiring an extremely high depth of coverage (for example, the reported limitations for ConStrains^{88,99}) and not being able to handle large magnitudes of data (for example, the reported limitations for the Lineages algorithm^{99,105}). These selected examples demonstrate how limitations can arise in foundational software as the metagenomic field progresses towards larger and more complex datasets. These and other limitations result in some tools being difficult or impossible to run or not feasible to use with current reasonably sized datasets, preventing the results from being reproducible or extendible.

The software referenced in this Review are examples of tools that reportedly perform the methodological approaches described. These references are not endorsements or reports of accuracy or usability. The reported features of many tools have been compared in recent reviews^{19,20}, but a thorough comparison of accuracies has not yet been completed (although they are expected to be addressed in the Critical Assessment of Metagenome Interpretation (CAMI)¹²³ framework). Future work is expected to make comparisons for within-species analysis software; however, what exactly is meant by the specific terminology of each tool (for example, SNV type, strain populations and so on) and their mapping to common terms (for example, strain and subspecies) will have to be carefully handled.

Terms for genetic stratification. There are many terms that stratify variation within species (TABLE 1). From the terms that are most commonly used and recognized by the International Code of Nomenclature of Prokaryotes⁴⁴, we highlight three to cover the range of genetic variation within species: genome, strain and subspecies (FIG. 2c). In this section, we discuss conflicts in the usage of these terms in culture-based microbiology and metagenomics and suggest solutions.

For decades, the most common source of microbial genomes was the sequencing of isolates. Recently, the prevalence of isolate genomes has been overtaken owing to the rapid production of MAGs. A barrier to synergy between isolate-based and metagenomic research stems from the misinterpretation of MAGs as equivalent to isolate genomes (BOX 3). The former might represent a population containing considerable diversity, whereas the latter usually represents a cultured isolate with little diversity. Considering also the rise in single-cell sequencing, it will be useful to increasingly qualify the term ‘genome’ as cellular, isolate or metagenomic.

The term ‘strain’ is widely used across fields in microbiology and has many contrasting definitions (FIG. 2a). In bacteriology, “a strain is made up of the descendants of a single isolation in pure culture, and usually is made up of a succession of cultures ultimately derived from an initial single colony”^{98,124} founded by one or more cells⁴⁴. This is a strain in the taxonomic sense²¹ (taxonomic strain), used for type strains and culture collections. In this case, the origin of a strain is at isolation.

An alternative definition used, for example, in epidemiology, recognizes a strain as an entity existing in nature²¹. This ‘natural strain’ is defined as a set of conspecific isolates with distinctive genotypic and/or phenotypic characteristics¹²⁵. A ‘taxonomic strain’ can be thought of as an isolated, cultured sample of a natural strain²¹. Operationally, the boundaries of natural and taxonomic strains vary. For example, taxonomic strains can become phenotypically heterogeneous with as few as three mutations¹²⁶ but would still be called the same strain. By contrast, in some cases, isolates need to have less than three SNV differences¹²⁷ to be considered to come from the same natural strain. This demonstrates that the genetic thresholds for strain delineation have not been universally set in culture-centric microbiology.

These two definitions of strain, among others¹²⁸, continue to coexist in culture-centric microbiology, and adoption of the term in microbiomics has extended this complexity. The disambiguating prefixes ‘taxonomic’ versus ‘natural’ are rarely used; however, this duality can clarify the mixed usage of the term ‘strain’ in metagenomics. Strain-level metagenomics often poses two types of questions: classification and discovery. Classification questions ask if genetic segments (sequencing reads) belong to a particular taxonomic strain, such as detecting if the probiotic strain *Bifidobacterium bifidum* BB12 is present in a stool sample. Discovery questions ask if there are subgroups within a species that form natural strains, for example, by clustering the genetic variation of genomes or of genetic segments. Conflict can arise among metagenomic tools for strain discovery that use different definitions of a natural strain and will implicitly therefore give different results, for example, defining natural strains based on differential gene content⁹⁶ versus based on SNVs in shared genes¹⁰⁶.

A universally applicable, operational definition of strain with a strong biological basis has not been established and may not exist. In theory, genomes with as few as one SNV difference could be referred to as different strains. However, this practice is not recommended owing to the unmanageable number of strains it would produce from metagenomic data. There are no rules on how many SNVs define a separate strain and whether such SNVs need to be fixed in the population or need to effect phenotype. In practice, the choice of how to set this cut-off is implicit in the choice of the strain-level profiling tool (for example, more than 0.1% of the nucleotides on species-specific marker genes as set in StrainPhlAn) or is set by the analysis authors (for example, greater than 98% ANI¹²⁹). Given such variability in the operational definition of a strain, it becomes particularly valuable to use more specific terminology instead of the generic term ‘strain’ (see TABLE 1 and the section entitled ‘Applications of within-species variation’ for guidelines).

Subspecies are groups of conspecific strains, and many definitions of the term exist¹³⁰. In classic microbiology, subspecies are clusters of strains that are genetically or phenotypically distinct, have a type strain available⁴⁴ and are named (for example, *Bacillus subtilis* subsp. *subtilis*). Over time, the basis for classification of subspecies has shifted from qualitative phenotypic

Type strains

Living cultures that serve as a fixed reference point for the assignment of bacterial and archaeal names. They are descended from the original isolate used in a species’ description and share all of its relevant phenotypic and genotypic properties.

Microbiomics

The study of microbial communities (microbiomes) using one or more -omic approaches; for example, genomics, transcriptomics and proteomics.

Table 1 | Definitions of terms used to stratify or describe variation within species

Term	Definition	Notes
Genotype	The set of alleles of an organism	Variable throughout time owing to mutation and recombination
Haplotype	Set of alleles or SNVs that are inherited together from a single parent ¹⁶⁵	Genetic signature of a lineage or clonal line, which can be disrupted through recombination
Haplogroup	Group of similar haplotypes with a common ancestor that has one or more clade-specific SNVs ¹⁶⁶	In human context, used to describe a group of people who share a common ancestor
Lineage and sublineage	Unbranched sequence of ancestral and descendant entities. Each ancestor may have multiple descendants, but only one is included in the lineage. Each entity could be an organism, clade, population or subspecies, among others ¹⁶⁷ . A sublineage is a subsection of a lineage	Can be visualized as an unbranched path through an evolutionary tree
Clone	Population of bacterial cells derived from a single parent cell ⁴⁴	In evolutionary terms, it is assumed to include all the descendants of the parent cell (monophyletic) ²¹ . Cultured isolates are samples of clones
Isolate	A pure culture obtained from a single colony separated from others in vitro ¹⁶⁸	Presumed to be and usually is derived from a single organism
Clade	Group of taxonomic entities composed of one ancestor and all of its evolutionary descendants ¹⁶⁷	Synonym: monophyletic group
Strain	Set of genetically similar descendants of a single colony or cell ⁴⁴ . Depending on the field, it can be genetic or phenotypic based	Descriptive subdivision of a species. Used widely but often with loose and/or inconsistent definitions. Can be described as 'taxonomic' or 'natural' ²¹
Within-species variant	Any subclassification of a species	General term that does not imply a level of resolution or phylogeny
Classic or typed subspecies	Set of strains that are genetically or phenotypically distinct and have a type strain available in a culture collection ¹⁶⁹ ; for example, <i>Lactococcus lactis</i> ssp. <i>lactis</i> and <i>L. lactis</i> ssp. <i>cremoris</i>	The name of a classic subspecies cannot be validly published if the description is based on studies of a mixed culture ⁴⁴ Variety was used as synonym of subspecies (now deprecated) ⁴⁴
Population subspecies	Set of local populations of strains that live in a subdivision of a species' spatial range and differs from other populations of the same species by phenotypic or genotypic characteristics ^{75,130}	Species with subspecies are 'polytypic', without are 'monotypic'
Population	Group of organisms that live in a particular location or ecological niche at a given time	Can be used to refer to all members of a species or to a subset of the entire population
Subpopulation	Portion of a population that is partially isolated from others and in which allele frequencies evolve independently ¹⁷⁰	A 'metapopulation' is a group of subpopulations
Strain population	A set of strains living simultaneously in the same spatial location or niche	Distinct from population subspecies, which can include multiple populations or ecotypes
Ecotype	An ecologically homogeneous population ⁷² . A clade within a species that has adapted to a particular environment. The scale of genetic dissimilarity between ecotypes can vary greatly	Ecotypes must be ecologically distinct enough that they can coexist indefinitely ¹⁷¹ . A mutant within an ecotype can outcompete the other strains in its own ecotype, but not those from a different ecotype ⁶⁹
Phylotype	Clade in which all members contain a homologous sequence (one or more marker genes, genetic or inter-genic regions) that are distinctively similar	The threshold level of similarity may be arbitrarily chosen. Not limited to within species
SNV type or SNP type	Set of genomes that share a distinctive set of SNVs ¹⁰⁸	Also used to describe the type of a SNV (for example, the exact switch in nucleotides)
Structural variant	Set of genomes that share distinctive structural variations ¹⁷²	Structural variations can be defined as insertions, deletions and inversions greater than 50 bp in size ¹⁷²
Pathotype	Set of genomes that cause the same disease using the same set of virulence factors ¹⁷³	Based on observational data; phenotypic and genotypic. It is not necessarily a clade
Serotype and serovar	Cells or viruses classified together based on their cell surface antigens, allowing the epidemiologic classification of organisms to the subspecies level ^{174–176}	Different strains can belong to the same serotype ¹⁷⁷ . Certain serotypes are often associated with specific pathotypes ¹⁷⁸
Phagotype (or phage type)	Set of genomes susceptible to a particular bacteriophage and demonstrated by phage typing ¹⁷⁹	Also called 'lysotype' ¹⁷⁹

SNV, single-nucleotide variant.

measures to genomic similarities between isolates¹³¹. This change has resulted in classification switches, such as the demotion of species to subspecies (for example, in *Bifidobacterium longum*¹³²) and vice versa (for example, in *Polynucleobacter necessarius*¹³³). Thus, named classical

subspecies do not (yet) necessarily align with distinct genomic clusters. By contrast, in a population biology context, a subspecies is a set of local populations that live in a subdivision of a species' spatial range and that differ from other populations of the same species⁷⁵,

for example, by genotype or phenotype¹³⁰. Adapting the term ‘subspecies’ for microbiomics implies the same usage dichotomy as described for strains: classification of reads to an existing ‘classic subspecies’ and discovery of ‘population subspecies’ by clustering the within-species genetic variation observed across spatial scales.

Although the strict definitions of these terms do not limit the relative amounts of variation they can each contain, in practice, it is useful to put them in context of each other and use them in the suggested ranges (FIG. 2c). As these ranges are guidelines, actual thresholds for group delineations should be included in reports when each term is used. Importantly, ‘strain’ is subordinate to ‘subspecies’ and thus should not be used to refer generally to any grouping subordinate to species (as it sometimes is). We also discourage use of the term ‘sub-species’ to mean ‘below species’ owing to its different definition but visual similarity to ‘subspecies’. Instead, we recommend using the terms ‘intraspecific’ or simply ‘within-species’. For example, inappropriate usage of ‘strain-level analysis’ or ‘subspecies analysis’ would be replaced with ‘intraspecific analysis’ or ‘within-species analysis’. Additionally, non-specific groupings within species can be referred to as ‘within-species variants’.

Phenotypic stratification in microbial communities.

Genetic variation within a species can manifest as phenotypic differences in complex ways. Different genetic variants can manifest as the same phenotype, whereas the same genetic variant can manifest as different phenotypes under different conditions¹³⁴. The scale of genetic differences and their phenotypic impact are also not necessarily correlated such as in the case of a dramatic increase in antibiotic resistance being conferred with as little as one SNV^{112,113}. Further, different phenotypes can be observed when bacteria are cultured in isolation or in coculture or are within their natural community. For example, *Pseudomonas aeruginosa* has distinct gene expression profiles in vitro versus during human infection, including genes involved in antibiotic resistance, cell–cell communication and metabolism, with implications for therapy development¹³⁵. Differences in phenotype can also be seen within species, for example, two strains of the halophilic bacterium *Salinibacter ruber* had similar expression patterns when cultured in isolation but had distinct patterns when grown in coculture¹³⁶. These examples highlight the importance of studying phenotypic variation within species directly in microbiomes, which can be done through several methods (BOX 1); for example, metatranscriptomics has been used to reveal functional diversity between conspecific symbionts in mussels¹³⁷, and metagenomically inferred replication rates have distinguished between intraspecific subpopulations of *Citrobacter koseri* in infants¹³⁸.

The complicated relationship between genotype and phenotype implies that phenotypic classification schemes can be at odds with genetic stratifications. In medicine and epidemiology, it has been useful to categorize bacteria into (possibly polyphyletic) groups

based on differential pathogenicity (pathotypes) or cell surface antigens (serotypes). For example, the enteric *E. coli* group includes both commensal and pathogenic strains, which are divided into at least seven pathotypes⁴⁵. In ecology, groups can also be defined based on behaviour and their functional role in a community, for example, based on the type of resources used and the way in which they are exploited^{139,140}. Species grouped in this way are called ‘guilds’, a concept and term that could similarly be used to describe groups of strains. This kind of grouping was designed to provide an appropriate resolution for the analysis of competition within ecosystems and the generalization of findings across communities. Although phenotype is the most relevant to many biological questions, it is hard to measure at large scale (though methods are progressing^{4,141}). With microbiome genetic sequencing, genotypes are much easier to measure in high throughput but linking them to phenotype is challenging as phenotype can change drastically with habitat and small genotypic differences.

Applications of within-species variation

The many scales and dimensions of variation within species reflect the wide range of biological questions that a within-species investigation can address. Isolate-based approaches have been used to investigate many biological questions that involve within-species variation^{142,143}. With the rise of metagenomic approaches, some of the same questions can now be investigated in high throughput and for many species in the community simultaneously (with important limitations; BOX 2; BOX 3). Below, we describe how many of the important biological applications that were pioneered using isolate-based methods can now be investigated using a metagenomics approach. We summarize common examples of such investigations into five major themes, built around key biological questions (FIG. 3). For each theme, we summarize the methodological approaches and appropriate terminology and provide examples of relevant studies or software.

Source tracking. Where did the cells in this sample originally come from? To determine patterns of transmission or dispersal of microbial cells, their exact source population must be identified. The probability that a cell was dispersed from or is a direct descendant of a particular source population can be calculated by comparing genetic material from the target cell or population against genetic material from its potential source population or ancestors (source tracking, transmission tracking or lineage tracking) (FIG. 3a). Strategies to determine source populations from metagenomic data include detecting the presence of shared SNVs^{86,88,89,91}, CRISPR signals¹⁴⁴, or strain-specific genes⁹⁰ and genome reconstruction¹⁴⁵. These approaches have been used to assess, for example, whether there is transmission of bacterial cells from the human oral cavity to the gut⁸⁸, from mother to infant^{86,89}, from probiotic treatment to the consumer⁹⁰, or from faecal microbiome transplant donor to recipient⁹¹. These strategies can be complicated by metagenomic disruption of allele linkage, multiple source populations and evolution of the target population

Polyphyletic

Describes a group of organisms that do not share an immediate common ancestor; not a clade.

Guilds

A guild is a group of species that use the same type of resources in a similar way; although originally defined as a group of species (Root, 1967), the concept could be applied to strains or subspecies.

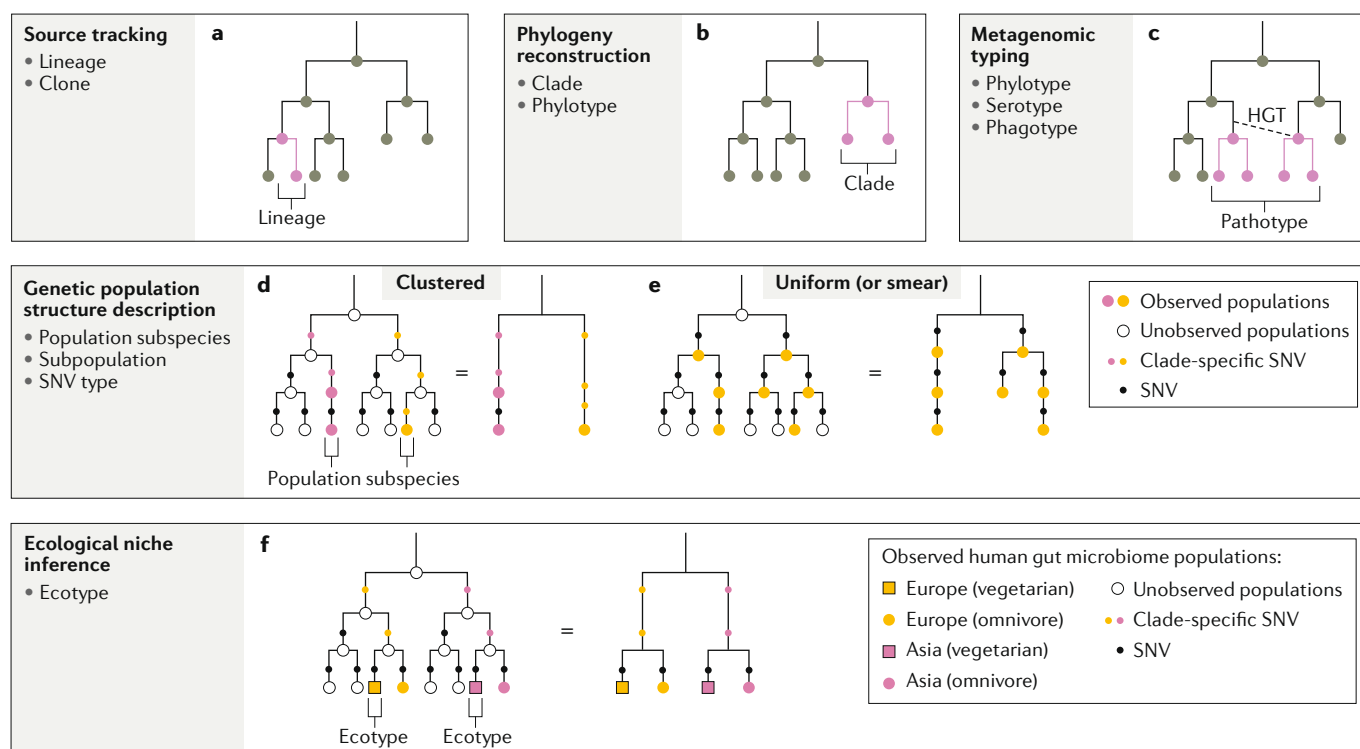


Fig. 3 | Applications of within-species variation. Five major areas of investigation for within-species-oriented metagenomic data analysis are illustrated, paired with the corresponding appropriate terminology. Trees depict the genetic similarity and ancestry of potentially coexisting populations, with nodes representing populations and edges representing genetic differences accumulating from top to bottom. **a** | Source tracking is concerned with identifying an unbranched path through a tree of ancestors and descendants (a 'lineage', pink edges and nodes). **b** | Phylogeny reconstruction aims to build a tree that reflects the history of within-species variants based on their genetic similarity. A phylogeny might be cut into complete sub-trees ('clades'), which may be called 'phylotypes'. **c** | Metagenomic typing detects the presence of a previously identified signature of interest within a species. For example, the presence of a gene associated with pathogenicity could be the criteria for detecting a 'pathotype'. This gene may have been transferred between clades via horizontal gene transfer (HGT), so may be at odds with the within-species phylogeny. **d,e** | The genetic population structure of a species can be described from the distribution of

the genetic similarities across observed variants. **d** | A 'clustered' structure occurs when there is a discontinuity across genetic similarities, enabling clades to be grouped into distinct clusters. Such a non-uniform structure is created by unobserved (extinct or unsampled) intermediate populations. A hypothetical within-species history with unobserved populations (white nodes) can be simplified (=), showing how unobserved populations can lead to a clustered genetic distribution, which may include distinct population subspecies. As single-nucleotide variants (SNVs) (black dots) accumulate through this history, some might be specific to a particular set of populations (coloured dots). **e** | When unobserved intermediate populations are rare or when they are spread widely through a species, the genetic distribution appears uniform or smeared, and distinct groups of populations are not seen. **f** | Ecological niche inference combines population observational data with phenotypic and/or habitat data to identify populations that have adapted to particular niches ('ecotypes'). Adaptive traits might be identified by comparing populations but potential geographic confounds must also be considered.

since dispersal from its source. Thus, although lineage tracking approaches can be useful for pathogen source detection¹⁴⁵, they can be insufficient for epidemiological outbreak analysis¹⁴⁶. In the context of source tracking, the general term 'strain' could be replaced by the more specific term 'lineage', which can be characterized by a haplotype. Determining genomic haplotypes from metagenomic data remains a challenge¹⁴⁷; however, long-read sequencing of single DNA molecules provides promise as error rates decline^{148,149}.

Phylogeny reconstruction. What is the evolutionary history of variants within this species? In phylogeny reconstruction (FIG. 3b), the relative ancestry of multiple lineages within a species is inferred from genetic similarity. This similarity can be based on full genomes or genetic segments (for example, marker genes). Owing to HGT and homologous recombination, the phylogeny that would be reconstructed can vary based on the loci

chosen and the phylogeny of genetic segments may not reflect the overall genomic phylogeny¹⁵⁰. Alternatively, within-species phylogenetic studies might focus on reconstruction of the history of a particular gene or plasmid within a species. Phylogeography puts these histories in the context of observed geographic distributions^{151,152}. Phylogenetic analysis using isolate genomes is well established¹⁵³, and these methods can be applied to microbial communities if high quality genomes are recovered, for example, using MAGs or single amplified genomes⁸⁷. However, data quality issues must be considered before this application (BOX 3). Alternatively, a typical approach is to identify conspecific, homologous genetic segments in metagenomes (for example, through alignment to reference sequences), detect SNVs in them^{56,104–106} and then infer their most probable history¹⁰⁶. Groups within species can be defined based on phylogeny by cutting the resultant tree at an arbitrary level of similarity,

creating ‘phylotypes’. In this context, the general term ‘strain’ could be replaced by the more specific terms ‘clade’ or ‘phylotype’.

Genetic population structure description. Does this species have distinct subpopulations and/or subspecies? Describing the genetic population structure of a species can, for example, suggest its geographic history or explain heterogeneous associations with host disease states⁷⁶. A species’ population structure can be determined by overlaying genetic data with observational data to describe the distribution of genetic similarities between variants within and across populations¹⁵⁴. A uniform structure (smear) occurs when there is a smooth distribution in genetic similarity across the observed species variants; this occurs when populations of ancestral and sister clades exist, that is, there are few unobservable (extinct or undetectable) branches within a tree (FIG. 3e). By contrast, a clustered structure occurs when there is a discontinuity across genetic similarities, enabling clades to be grouped into distinct clusters. Such a non-uniform structure is created by extinct branches within a tree (FIG. 3d); this manifests as subpopulations, which are subsets of a whole population that have distinct frequencies of genetic variations (for example, alleles or SNVs).

Metagenomics can be used to study population genetics of species within microbiomes¹⁹ by looking for clustering of genetic similarities across potential subpopulations. Detecting subpopulations is sensitive to sampling effort, as discontinuities in genetic similarity can be due to failure to observe intermediates (BOX 3). Assessing such genetic similarities can be based on SNV allele frequencies in whole genomes^{76,105,109}, SNVs in marker genes^{56,106} or gene content differences¹⁵⁵. When MAGs or single amplified genomes are produced, genome-based ANI clustering can also be used¹⁵⁶. MAGs can also be used to track SNV and gene content differences such as changes in populations of lake bacteria over time¹⁰³. In this context, ‘strain’ is sometimes inappropriately used to refer to a subpopulation or subspecies. Subpopulations might be ecotypes if they have adapted to different niches, for example, through a genome-wide sweep instead of a gene-specific sweep^{72,157}.

Ecological niche inference. Have the variants within this species adapted to different conditions? Looking at within-species variants in conjunction with their habitats can provide information about their niche specificity (FIG. 3f). When genetic data are used to make inferences about uncharacterized habitats, this is sometimes referred to as ‘reverse ecology’⁷¹⁵⁸. These inquiries often aim to identify the genetic segments (for example, genes, operons or plasmids) that are key to adapting to particular environments. Acquisition of these segments might be from vertical or horizontal transmission and can thus be in contrast with the phylogenetic history of the species. For example, a gene can rapidly become ubiquitous across populations owing to frequent HGT under selective conditions (gene-specific sweep), such as in the presence of antibiotics⁷². A common approach to investigate these questions using metagenomic data

is to look at conspecific subpopulations of cells that are known to have adapted to different conditions (for example, different human host diets⁸⁵, soil versus plant host¹⁵⁹ or shifts in lake water habitats¹⁰³) and subsequently identify distinctive genes^{83,96,99,105}. Methods used in metagenome-wide association studies can also be applied here, though these are not often focused on the adaptive evolution of populations¹⁶⁰. In this context, the general term ‘strain’ could be replaced by the more specific term ‘ecotypes’⁷².

In the example shown in FIG. 3, genetic population structure investigations would focus on the allele frequency differences between European and Asian populations to decide whether these are distinct subpopulations or belong to one continuous population. Investigations on ecological niche inference would focus on the gene differences in the gut-associated microbiome species associated with different diets, regardless of whether the European and Asian populations are distinct subpopulations.

Typing. Does this species variant belong to a previously described subgroup of the species? Typing analyses assess the presence of genetic features (for example, SNVs, genes, operons or plasmids) of specific interest in conspecific species variants (FIG. 3c). In this context, within-species groups are not defined based on evolutionary history or habitat ranges but simply on the presence or absence of specific genetic features. Such features may confer habitat fitness, may be transient and may only be expressed under rare or artificial conditions such as antimicrobial resistance genes, pathogenicity genes (for example, enteropathogenic *E. coli*) or flagella. In this case, HGT is a major consideration because its result means that the presence of a genetic feature does not necessarily reflect phylogeny. An example of a ‘type’ in this context is serogroups, which are potentially polyphyletic groups within a species that are defined based on the presence of cell surface antigens, which allows their epidemiological classification.

Metagenomic approaches can be used to detect the genetic features that define a type. SNVs of known¹¹⁷ or novel¹⁰⁵ importance can be detected based on reference sequences. Detecting the presence of type-defining genes based on homology to reference sequences is well established in metagenomics^{147,161} but determining with certainty that these detected genes are present in a specific strain is more difficult owing to the possibility of HGT within the community. In metagenomic data, HGT can be studied directly, with¹⁶² or without¹⁶³ assembling genomes (reviewed in REF. 164).

Comparative analyses of within-species variants with the same phenotype can be used to discover the specific genetic features that are associated with (and may be causing) the phenotype (such as in metagenome-wide association studies¹⁶⁰). For example, conspecific cells could be grouped into a pathogenic ‘variant’ based only on their presence within hosts that are displaying similar symptoms, without knowing the evolutionary relationship of the cells or their typical habitats. In this context, the general term ‘strain’ could be replaced by the more specific term ‘pathotype’.

Genome-wide sweep

Alleles at the locus under selection cause other linked loci (for example, genome and plasmid) to gain or lose abundance across the population; also known as a broad sweep.

Gene-specific sweep

Only alleles at the locus under selection gain or lose abundance across the population; also known as a narrow or locus-specific sweep.

The themes described above have traditionally been investigated using isolate genomic approaches or low-resolution molecular methods (BOX 1). As metagenomic studies increasingly create large amounts of data, dozens of new methods have been established to investigate the same questions, often with their own novel vocabulary. Considering how these new methods map back to the fundamental biological questions they are addressing and the history of research in the area will help to control the explosion of terminology. Many studies will include a combination of the themes described above, but considering the fundamental units separately facilitates the break-down of complex questions and selection of the most appropriate methodology and terminology.

Conclusions

Despite often being the highest resolution taxonomic category considered in microbiome surveys, species can contain extreme phenotypic variability. Studying such variability used to be relatively limited in scope, with a few key isolate-based methods and a limited pool of culturable bacteria. With the development of

metagenomic sequencing, the number of species that can be studied and the number of methods that can be used have increased substantially. The possibility to stratify variation within species according to many criteria, and at many scales, has also led to a growing and frequently imprecise terminology. Understanding how the variability within a species arose and identifying the central biological question being asked can help to determine the correct terminology and methodology to use. In some cases, the most appropriate term may have an operational definition, and its details and cut-off thresholds might vary across studies. To facilitate communication and collaboration, and enable future comparative meta-studies, vocabulary that does not have strict and widely known definitions should be avoided when possible or explicitly described both in terms of the criteria and the thresholds being used. This Review aims to guide such descriptions and support a more informed development and application of within-species investigation techniques to metagenomic data.

Published online: 04 June 2020

- Wayne, L. G. et al. Report of the Ad Hoc Committee on Reconciliation of Approaches to Bacterial Systematics. *Int. J. Syst. Bacteriol.* **37**, 463–464 (1987).
- Leimbach, A., Hacker, J. & Dobrindt, U. E. coli as an all-rounder: the thin line between commensalism and pathogenicity. *Curr. Top. Microbiol. Immunol.* **358**, 3–32 (2013).
- Pierce, J. V. & Bernstein, H. D. Genomic diversity of enterotoxigenic strains of bacteroides fragilis. *PLoS One* **11**, e0158171 (2016).
- Maier, L. et al. Extensive impact of non-antibiotic drugs on human gut bacteria. *Nature* **555**, 623–628 (2018).
- Neuenschwander, S. M., Ghai, R., Pernthaler, J. & Salcher, M. M. Microdiversification in genome-streamlined ubiquitous freshwater Actinobacteria. *ISME J.* **12**, 185–198 (2018).
- Triplett, E. & Sadowsky, M. J. Genetics of competition for nodulation of legumes. *Annu. Rev. Microbiol.* **46**, 399–428 (1992).
- Nowrouzian, F. L., Adlerberth, I. & Wold, A. E. Enhanced persistence in the colonic microbiota of *Escherichia coli* strains belonging to phylogenetic group B2: role of virulence factors and adherence to colonic cells. *Microbes Infect.* **8**, 834–840 (2006).
- Whitman, W. B. & Bergey's Manual Trust. *Bergey's Manual of Systematics of Archaea and Bacteria* (Wiley, 2015).
- Zhao, S. et al. Adaptive evolution within gut microbiomes of healthy people. *Cell Host Microbe* **25**, 656–667 (2019).
- Lagier, J.-C. et al. Culturing the human microbiota and culturomics. *Nat. Rev. Microbiol.* **16**, 540–550 (2018).
- Tyson, G. W. et al. Community structure and metabolism through reconstruction of microbial genomes from the environment. *Nature* **428**, 37–43 (2004).
- Allen, E. E. et al. Genome dynamics in a natural archaeal population. *Proc. Natl Acad. Sci. USA* **104**, 1883–1888 (2007).
- Eppey, J. M., Tyson, G. W., Getz, W. M. & Banfield, J. F. Genetic exchange across a species boundary in the archaeal genus ferropasma. *Genetics* **177**, 407–16 (2007).
- Eppey, J. M., Tyson, G. W., Getz, W. M. & Banfield, J. F. Strainer: software for analysis of population variation in community genomic datasets. *BMC Bioinformatics* **8**, 398 (2007).
- Lo, I. et al. Strain-resolved community proteomics reveals recombining genomes of acidophilic bacteria. *Nature* **446**, 537–541 (2007).
- Denef, V. J. et al. Proteogenomic basis for ecological divergence of closely related bacteria in natural acidophilic microbial communities. *Proc. Natl Acad. Sci. USA* **107**, 2383–2390 (2010).
- Segata, N. On the road to strain-resolved comparative metagenomics. *mSystems* **3**, e00190-17 (2018).
- Suez, J., Zmora, N., Segal, E. & Elinav, E. The pros, cons, and many unknowns of probiotics. *Nat. Med.* **25**, 716–729 (2019).
- Denef, V. J. in *Population Genomics: Microorganisms* (eds Polz, M., Rajora, O.) 49–75 (Springer, 2018).
- Comprehensive review on the application of metagenomic approaches for microbial population genomics.
- Bobay, L.-M. & Raymann, K. Population genetics of host-associated microbiomes. *Curr. Mol. Biol. Rep.* **5**, 128–139 (2019).
- Dijkshoorn, L., Ursing, B. M. & Ursing, J. B. Strain, clone and species: comments on three basic concepts of bacteriology. *J. Med. Microbiol.* **49**, 397–401 (2000).
- Compares and summarises definitions of key terminology in a bacteriological (isolate-based) context.
- Brown, T. *Genomes* 2nd edn (Wiley-Liss, 2002).
- Alberts, B. et al. *Molecular Biology of the Cell* (Garland Science, 2002).
- Fijalkowska, I. J., Schaaper, R. M. & Jonczyk, P. DNA replication fidelity in *Escherichia coli*: a multi-DNA polymerase affair. *FEMS Microbiol. Rev.* **36**, 1105–21 (2012).
- Denamur, E. & Matic, I. Evolution of mutation rates in bacteria. *Mol. Microbiol.* **60**, 820–827 (2006).
- Dillon, M. M., Sung, W., Sebra, R., Lynch, M. & Cooper, V. S. Genome-wide biases in the rate and molecular spectrum of spontaneous mutations in *Vibrio cholerae* and *Vibrio fischeri*. *Mol. Biol. Evol.* **34**, 93–109 (2017).
- Strauss, C., Long, H., Patterson, C. E., Te, R. & Lynch, M. Genome-wide mutation rate response to pH change in the Coral Reef Pathogen *Vibrio shilonii* AK1. *mBio* **8**, e01021-17 (2017).
- Cooper, V. S., Vohr, S. H., Wrocklage, S. C. & Hatcher, P. J. Why genes evolve faster on secondary chromosomes in bacteria. *PLoS Comput. Biol.* **6**, e1000732 (2010).
- Bobay, L.-M., Traverse, C. C. & Ochman, H. Impermanence of bacterial clones. *Proc. Natl Acad. Sci. USA* **112**, 8893–8900 (2015).
- Andersson, J. O. & Andersson, S. G. E. Pseudogenes, junk DNA, and the dynamics of Rickettsia Genomes. *Mol. Biol. Evol.* **18**, 829–839 (2001).
- Mira, A., Ochman, H. & Moran, N. A. Deletional bias and the evolution of bacterial genomes. *Trends Genet.* **17**, 589–96 (2001).
- Lawrence, J. G. & Retchless, A. C. The interplay of homologous recombination and horizontal gene transfer in bacterial speciation. *Methods Mol. Biol.* **532**, 29–53 (2009).
- Lerner, A., Matthias, T. & Aminov, R. Potential effects of horizontal gene exchange in the human gut. *Front. Immunol.* **8**, 1630 (2017).
- Thomas, C. M. & Nielsen, K. M. Mechanisms of, and barriers to, horizontal gene transfer between bacteria. *Nat. Rev. Microbiol.* **3**, 711–721 (2005).
- Reviews the major concepts and mechanisms of HGT and their implications for genome flux across populations.
- Rocha, E. P., Cornet, E. & Michel, B. Comparative and evolutionary analysis of the bacterial homologous recombination systems. *PLoS Genet.* **1**, e15 (2005).
- Fraser, C., Hanage, W. P. & Spratt, B. G. Recombination and the nature of bacterial speciation. *Science* **315**, 476–480 (2007).
- Gasiunas, G., Sinkunas, T. & Siksnys, V. Molecular mechanisms of CRISPR-mediated microbial immunity. *Cell. Mol. Life Sci.* **71**, 449–465 (2014).
- Brouwer, M. S. M. et al. Horizontal gene transfer converts non-toxicogenic *Clostridium difficile* strains into toxin producers. *Nat. Commun.* **4**, 2601 (2013).
- Kaper, J. B. & O'Brien, A. D. Overview and historical perspectives. *Microbiol. Spectr.* **2** <https://doi.org/10.1128/microbiolspec.EHEC-0028-2014> (2014).
- Hallatschek, O., Hersen, P., Ramanathan, S. & Nelson, D. R. Genetic drift at expanding frontiers promotes gene segregation. *Proc. Natl Acad. Sci. USA* **104**, 19926–19930 (2007).
- Nemergut, D. R. et al. Patterns and processes of microbial community assembly. *Microbiol. Mol. Biol. Rev.* **77**, 342–356 (2013).
- Chun, J. et al. Proposed minimal standards for the use of genome data for the taxonomy of prokaryotes. *Int. J. Syst. Evol. Microbiol.* **68**, 461–466 (2018).
- Ford Doolittle, W. Population genomics: how bacterial species form and why they don't exist. *Curr. Biol.* **22**, R451–R453 (2012).
- International Committee on Systematics of Prokaryotes. International Code of Nomenclature of Prokaryotes: Prokaryotic Code (2008 Revision). *Int. J. Syst. Evol. Microbiol.* **69**, S1–S111 (2019).
- Croxen, M. A. et al. Recent advances in understanding enteric pathogenic *Escherichia coli*. *Clin. Microbiol. Rev.* **26**, 822–880 (2013).
- Konstantinidis, K. T. & Tiedje, J. M. Genomic insights that advance the species definition for prokaryotes. *Proc. Natl Acad. Sci. USA* **102**, 2567–2572 (2005).
- Richter, M. & Rosselló-Móra, R. Shifting the genomic gold standard for the prokaryotic species definition. *Proc. Natl Acad. Sci. USA* **106**, 19126–19131 (2009).

48. Mende, D. R., Sunagawa, S., Zeller, G. & Bork, P. Accurate and universal delineation of prokaryotic species. *Nat. Methods* **10**, 881–884 (2015).
49. Goris, J. et al. DNA-DNA hybridization values and their relationship to whole-genome sequence similarities. *Int. J. Syst. Evol. Microbiol.* **57**, 81–91 (2007).
50. Džink, J. L., Sheenan, M. T. & Socransky, S. S. Proposal of three subspecies of *Fusobacterium nucleatum* Knorr 1922: *Fusobacterium nucleatum* subsp. *nucleatum* subsp. nov., comb. nov.; *Fusobacterium nucleatum* subsp. *polymorphum* subsp. nov., norm. rev., comb. nov.; and *Fusobacterium nucleatum* subsp. *vincentii* subsp. nov., norm. rev., comb. nov. *Int. J. Syst. Bacteriol.* **40**, 74–78 (1990).
51. Kook, J. K. et al. Genome-based reclassification of *Fusobacterium nucleatum* subspecies at the species level. *Curr. Microbiol.* **74**, 1137–1147 (2017).
52. Konstantinidis, K. T. & Delong, E. F. Genomic patterns of recombination clonal divergence and environment in marine microbial populations. *ISME J.* **2**, 1052–1065 (2008).
53. Caro-Quintero, A. & Konstantinidis, K. T. Bacterial species may exist, metagenomics reveal. *Environ. Microbiol.* **14**, 347–355 (2012).
54. Jain, C., Rodríguez-R. L. M., Phillippy, A. M., Konstantinidis, K. T. & Aluru, S. High throughput ANI analysis of 90K prokaryotic genomes reveals clear species boundaries. *Nat. Commun.* **9**, 5114 (2018).
55. Olm, M. R. et al. Consistent metagenome-derived metrics verify and delineate bacterial species boundaries. *mSystems* **5**, e00731-19 (2020).
56. Milanese, A. et al. Microbial abundance, activity and population genomic profiling with mOTUs2. *Nat. Commun.* **10**, 1014 (2019).
57. Mayden, R. L. in *Species. The Units of Biodiversity* (eds Claridge, M. F., Dawah, H. A. & Wilson, M. R.) 381–423 (Chapman & Hall, 1997).
58. Wilkins, J. S. How to be a chaste species pluralist-realist: the origins of species modes and the synapomorphic species concept. *Biol. Philos.* **18**, 621–638 (2003).
59. Hey, J. The mind of the species problem. *Trends Ecol. Evol.* **16**, 326–329 (2001).
60. Baptiste, E. et al. Prokaryotic evolution and the tree of life are two different things. *Biol. Direct.* **4**, 34 (2009).
61. Konstantinidis, K. T., Ramette, A. & Tiedje, J. M. The bacterial species definition in the genomic era. *Philos. Trans. R. Soc. B Biol. Sci.* **361**, 1929–1940 (2006).
62. Bobay, L.-M. & Ochman, H. Biological species are universal across life's domains. *Genome Biol. Evol.* **9**, 491–501 (2017).
63. Moldovan, M. A. & Gelfand, M. S. Pangenomic definition of prokaryotic species and the phylogenetic structure of *Prochlorococcus* spp. *Front. Microbiol.* **9**, 428 (2018).
64. Snel, B., Bork, P. & Huynen, M. A. Genome phylogeny based on gene content. *Nat. Genet.* **21**, 108–110 (1999).
65. Achtman, M. & Wagner, M. Microbial diversity and the genetic nature of microbial species. *Nat. Rev. Microbiol.* **6**, 431–440 (2008).
66. Barton, N. H. The effect of hitch-hiking on neutral genealogies. *Genet. Res.* **72**, 123–133 (1998).
67. Hermisson, J. & Pennings, P. S. Soft sweeps and beyond: understanding the patterns and probabilities of selection footprints under rapid adaptation. *Methods Ecol. Evol.* **8**, 700–716 (2017).
68. Shapiro, B. J. et al. Population genomics of early events in the ecological differentiation of bacteria. *Science* **336**, 48–51 (2012).
69. Cohan, F. M. Bacterial species and speciation. *Syst. Biol.* **50**, 513–524 (2001).
70. Cohan, F. M. in *Selective Sweep* (ed. Nurminsky, D.) 78–93 (Springer, 2007).
71. Charlesworth, B. Effective population size and patterns of molecular evolution and variation. *Nat. Rev. Genet.* **10**, 195–205 (2009).
72. Cohan, F. M. Bacterial speciation: genetic sweeps in bacterial species. *Curr. Biol.* **26**, R112–R115 (2016).
73. Hermisson, J. & Pennings, P. S. Soft sweeps: molecular population genetics of adaptation from standing genetic variation. *Genetics* **169**, 2335–2352 (2005).
74. González-Torres, P., Rodríguez-Mateos, F., Antón, J., Gabaldón, T. & Heitman, J. Impact of homologous recombination on the evolution of prokaryotic core genomes. *mBio* **10**, e02494-18 (2019).
75. Monroe, B. A modern concept of the subspecies. *Auk* **99**, 608–609 (1982).
76. Costea, P. et al. Subspecies in the global human gut microbiome. *Mol. Syst. Biol.* **13**, 960–960 (2017).
77. Retchless, A. C. & Lawrence, J. G. Temporal fragmentation of speciation in bacteria. *Science* **317**, 1093–1096 (2007).
78. Shapiro, B. J. in *Population Genomics: Microorganisms* (eds Polz, M. F. & Rajora, O. P.) 31–47 (Springer Nature, 2018).
79. Sheppard, S. K., Guttman, D. S. & Fitzgerald, J. R. Population genomics of bacterial host adaptation. *Nat. Rev. Genet.* **19**, 549–565 (2018).
80. Bobay, L.-M. & Ochman, H. Factors driving effective population size and pan-genome evolution in bacteria. *BMC Evol. Biol.* **18**, 153 (2018).
81. Smelov, V. et al. Chlamydia trachomatis strain types have diversified regionally and globally with evidence for recombination across geographic divides. *Front. Microbiol.* **8**, 2195 (2017).
82. Tenailon, O., Skurnik, D., Picard, B. & Denamur, E. The population genetics of commensal *Escherichia coli*. *Nat. Rev. Microbiol.* **8**, 207–217 (2010).
83. Zeevi, D. et al. Structural variation in the gut microbiome associates with host health. *Nature* **568**, 43–48 (2019).
84. Lloyd-Price, J. et al. Strains, functions and dynamics in the expanded human microbiome project. *Nature* **550**, 61–66 (2017).
85. De Filippis, F. et al. Distinct genetic and functional traits of human intestinal *Prevotella copri* strains are associated with different habitual diets. *Cell Host Microbe* **25**, 444–453 (2019).
86. Ferretti, P. et al. Mother-to-infant microbial transmission from different body sites shapes the developing infant gut microbiome. *Cell Host Microbe* **24**, 133–145 (2018).
87. Stewart, R. D. et al. Compendium of 4,941 rumen metagenome-assembled genomes for rumen microbiome biology and enzyme discovery. *Nat. Biotechnol.* **37**, 953–961 (2019).
88. Schmidt, T. S. et al. Extensive transmission of microbes along the gastrointestinal tract. *eLife* **8**, e42693 (2019).
89. Asnicar, F. et al. Studying vertical microbiome transmission from mothers to infants by strain-level metagenomic profiling. *mSystems* **2**, e00164-16 (2017).
90. Zmora, N. et al. Personalized gut mucosal colonization resistance to empiric probiotics is associated with unique host and microbiome features. *Cell* **174**, 1388–1405 (2018).
91. Smillie, C. S. et al. Strain tracking reveals the determinants of bacterial engraftment in the human gut following fecal microbiota transplantation article strain tracking reveals the determinants of bacterial engraftment in the human gut following fecal microbiota transplants. *Cell Host Microbe* **23**, 229–240 (2018).
92. Zhou, Z., Luhmann, N., Alikhan, N. F., Quince, C. & Achtman, M. in *Research in Computational Molecular Biology. RECOMB 2018. Lecture Notes in Computer Science*, vol 10812 (ed. Raphael, B.) 225–240 (Springer, 2018).
93. Ahn, T.-H., Chai, J. & Pan, C. Sigma: strain-level inference of genomes from metagenomic analysis for biosurveillance. *Bioinformatics* **31**, 170–177 (2015).
94. Hong, C. et al. PathoScope 2.0: a complete computational framework for strain identification in environmental or clinical sequencing samples. *Microbiome* **2**, 33 (2014).
95. Ondov, B. D. et al. Mash: fast genome and metagenome distance estimation using MinHash. *Genome Biol.* **17**, 132 (2016).
96. Scholz, M. et al. Strain-level microbial epidemiology and population genomics from shotgun metagenomics. *Nat. Methods* **13**, 435–438 (2016).
97. Zhu, A., Sunagawa, S., Mende, D. R. & Bork, P. Inter-individual differences in the gene content of human gut bacterial species. *Genome Biol.* **16**, 82 (2015).
98. Greenblum, S., Carr, R. & Borenstein, E. Extensive strain-level copy-number variation across human gut microbiome species. *Cell* **160**, 583–594 (2015).
99. Quince, C. et al. DESMAN: a new tool for de novo extraction of strains from metagenomes. *Genome Biol.* **18**, 181 (2017).
100. Maistrenko, O. M. et al. Disentangling the impact of environmental and phylogenetic constraints on prokaryotic within-species diversity. *ISME J.* **14**, 1247–1259 (2020).
101. Andreani, N. A., Hesse, E. & Vos, M. Prokaryote genome fluidity is dependent on effective population size. *ISME J.* **11**, 1719–1721 (2017).
102. Nayfach, S., Shi, Z. J., Seshadri, R., Pollard, K. S. & Kyrpides, N. C. New insights from uncultivated genomes of the global human gut microbiome. *Nature* **568**, 505–510 (2019).
103. Bendall, M. L. et al. Genome-wide selective sweeps and gene-specific sweeps in natural bacterial populations. *ISME J.* **10**, 1589–1601 (2016).
104. Nayfach, S., Rodríguez-Mueller, B., Garud, N. & Pollard, K. S. An integrated metagenomics pipeline for strain profiling reveals novel patterns of bacterial transmission and biogeography. *Genome Res.* **26**, 1612–1625 (2016).
105. Costea, P. I. et al. metaSNV: a tool for metagenomic strain level analysis. *PLoS One* **12**, e0182392 (2017).
106. Truong, D. T., Tett, A., Pasolli, E., Huttenhower, C. & Segata, N. Microbial strain-level population structure and genetic diversity from metagenomes. *Genome Res.* **27**, 626–638 (2017).
107. Bush, S. J. et al. Genomic diversity affects the accuracy of bacterial single-nucleotide polymorphism-calling pipelines. *Gigascience* **9**, gaa007 (2020).
108. Luo, C. et al. ConStrains identifies microbial strains in metagenomic datasets. *Nat. Biotechnol.* **33**, 1045–1052 (2015).
109. Delmont, T. O. et al. Single-amino acid variants reveal evolutionary processes that shape the biogeography of a global SAR11 subclade. *eLife* **8**, e46497 (2019).
110. Jackson, R. W. et al. Identification of a pathogenicity island, which contains genes for virulence and avirulence, on a large native plasmid in the bean pathogen *Pseudomonas syringae* pathovar phaseolicola. *Proc. Natl Acad. Sci. USA* **96**, 10875–10880 (1999).
111. Scholz, B. K., Jakobek, J. L. & Lindgren, P. B. Restriction fragment length polymorphism evidence for genetic homology within a pathovar of *Pseudomonas syringae*. *Appl. Environ. Microbiol.* **60**, 1093–1100 (1994).
112. Pan, X. S., Yague, G. & Fisher, L. M. Quinolone resistance mutations in *Streptococcus pneumoniae* gyrA and parC proteins: mechanistic insights into quinolone action from enzymatic analysis, intracellular levels, and phenotypes of wild-type and mutant proteins. *Antimicrob. Agents Chemother.* **45**, 3140–3147 (2001).
113. Forslund, K., Sunagawa, S., Coelho, L. P. & Bork, P. Metagenomic insights into the human gut resistome and the forces that shape it. *BioEssays* **36**, 316–329 (2014).
114. Petkau, A. et al. SNVPhyl: a single nucleotide variant phylogenomics pipeline for microbial genomic epidemiology. *Microb. Genom.* **3**, e000116 (2017).
115. Jain, R., Rivera, M. C., Lake, J. A. & Lake, J. Horizontal gene transfer among genomes: the complexity hypothesis. *Proc. Natl Acad. Sci. USA* **96**, 3801–3806 (1999).
116. Polz, M. F. & Rajora, O. P. (eds) *Population Genomics: Microorganisms*. (Springer, 2019).
117. Zolfo, M., Tett, A., Jousset, O., Donati, C. & Segata, N. MetaMLST: multi-locus strain-level bacterial typing from metagenomic samples. *Nucleic Acids Res.* **45**, e7 (2017).
118. Truong, D. T. et al. MetaPhlAn2 for enhanced metagenomic taxonomic profiling. *Nat. Methods* **12**, 902–903 (2015).
119. Tamburini, F. B. et al. Precision identification of diverse bloodstream pathogens in the gut microbiome. *Nat. Med.* **24**, 1809–1814 (2018).
120. Albanese, D. & Donati, C. Strain profiling and epidemiology of bacterial species from metagenomic sequencing. *Nat. Commun.* **8**, 2260 (2017).
121. O'Brien, J. D. et al. A Bayesian approach to inferring the phylogenetic structure of communities from metagenomic data. *Genetics* **197**, 925–937 (2014).
122. Smillie, C. S. et al. Strain tracking reveals the determinants of bacterial engraftment in the human gut following fecal microbiota transplantation. *Cell Host Microbe* **23**, 229–240 (2018).
123. Sczyrba, A. et al. Critical assessment of metagenome interpretation - a benchmark of metagenomics software. *Nat. Methods* **14**, 1063–1071 (2017).

124. Brenner, D. J., Staley, J. T. & Krieg, N. R. Classification of Prokaryotic Organisms and the Concept of Bacterial Speciation. in *Bergey's Manual of Systematics of Archaea and Bacteria*. 1–9 (John Wiley & Sons, Ltd, 2015).
125. Struelens, M. J. et al. Consensus guidelines for appropriate use and evaluation of microbial epidemiologic typing systems. *Clin. Microbiol. Infect.* **2**, 2–11 (1996).
126. Spira, B., de Almeida Toledo, R., Maharjan, R. P. & Ferenci, T. The uncertain consequences of transferring bacterial strains between laboratories - *rpoS* instability as an example. *BMC Microbiol.* **11**, 248 (2011).
127. Kong, L. Y. et al. *Clostridium difficile*: investigating transmission patterns between infected and colonized patients using whole genome sequencing. *Clin. Infect. Dis.* **68**, 204–209 (2019).
128. Saak, C. C. & Gibbs, K. A. The self-identity protein IdsD is communicated between cells in swarming proteus mirabilis colonies. *J. Bacteriol.* **198**, 3278–3286 (2016).
129. Brooks, B. et al. Strain-resolved analysis of hospital rooms and infants reveals overlap between the human and room microbiome. *Nat. Commun.* **8**, 1814 (2017).
130. Patten, M. A. Subspecies and the philosophy of science. *Auk* **132**, 481–485 (2015).
131. Meier-Kolthoff, J. P. et al. Complete genome sequence of DSM 30083(T), the type strain (U5/41(T)) of *Escherichia coli*, and a proposal for delineating subspecies in microbial taxonomy. *Stand. Genomic Sci.* **9**, 2 (2014).
132. Fukuyama, M. et al. Unification of *Bifidobacterium infantis* and *Bifidobacterium suis* as *Bifidobacterium longum*. *Int. J. Syst. Evol. Microbiol.* **52**, 1945–1951 (2002).
133. Hahn, M. W., Schmidt, J., Pitt, A., Taipale, S. J. & Lang, E. Reclassification of four *Polynucleobacter necessarius* strains as representatives of *Polynucleobacter asymbioticus* comb. nov., *Polynucleobacter duraquae* sp. nov., *Polynucleobacter yangtzensis* sp. nov. and *Polynucleobacter sinensis* sp. nov., and emended description of *Polynucleobacter necessarius*. *Int. J. Syst. Evol. Microbiol.* **66**, 2883–2892 (2016).
134. Ackermann, M. A functional perspective on phenotypic heterogeneity in microorganisms. *Nat. Rev. Microbiol.* **13**, 497–508 (2015).
135. Cornforth, D. M. et al. Pseudomonas aeruginosa transcriptome during human infection. *Proc. Natl Acad. Sci. USA* **115**, E5125–E5134 (2018).
136. González-Torres, P. et al. Interactions between closely related bacterial strains are revealed by deep transcriptome sequencing. *Appl. Environ. Microbiol.* **81**, 8445–8456 (2015).
137. Ansoorge, R. et al. Functional diversity enables multiple symbiont strains to coexist in deep-sea mussels. *Nat. Microbiol.* **4**, 2487–2497 (2019).
138. Olm, M. R. et al. Identical bacterial populations colonize premature infant gut, skin, and oral microbiomes and exhibit different in situ growth rates. *Genome Res.* **27**, 601–612 (2017).
139. Pedrós-Alió, C. in *Plankton Ecology* (ed. Sommer, U.) 297–336 (Springer, 1989).
140. Root, R. B. The niche exploitation pattern of the blue-gray gnatcatcher. *Ecol. Monogr.* **37**, 317–350 (1967).
141. Mateus, A. et al. Thermal proteome profiling in bacteria: probing protein state in vivo. *Mol. Syst. Biol.* **14**, e8242 (2018).
142. Land, M. et al. Insights from 20 years of bacterial genome sequencing. *Funct. Integr. Genomics* **15**, 141–161 (2015).
143. Gutleben, J. et al. The multi-omics promise in context: from sequence to microbial isolate. *Crit. Rev. Microbiol.* **44**, 212–229 (2018).
144. Lam, T. J. & Ye, Y. CRISPRs for strain tracking and their application to microbiota transplantation data analysis. *Cris. J.* **2**, 41–50 (2019).
145. Mu, A. et al. Reconstruction of the genomes of drug-resistant pathogens for outbreak investigation through metagenomic sequencing. *mSphere* **4**, e00529-18 (2019).
146. Didelot, X., Walker, A. S., Peto, T. E., Crook, D. W. & Wilson, D. J. Within-host evolution of bacterial pathogens. *Nat. Rev. Microbiol.* **14**, 150–162 (2016).
147. Quince, C., Walker, A. W., Simpson, J. T., Loman, N. J. & Segata, N. Shotgun metagenomics, from sampling to analysis. *Nat. Biotechnol.* **35**, 833–844 (2017).
- Reviews how microbial communities can be studied using metagenomic sequencing, with comments**
- on sources of bias and comparisons of analytical methods.**
148. Koren, S. & Phillippy, A. M. One chromosome, one contig: complete microbial genomes from long-read sequencing and assembly. *Curr. Opin. Microbiol.* **23**, 110–120 (2015).
149. Somerville, V. et al. Long-read based de novo assembly of low-complexity metagenome samples results in finished genomes and reveals insights into strain diversity and an active phage system. *BMC Microbiol.* **19**, 143 (2019).
150. Jiang, X. et al. Dissemination of antibiotic resistance genes from antibiotic producers to pathogens. *Nat. Commun.* **8**, 15784 (2017).
151. Linz, B. et al. An African origin for the intimate association between humans and *Helicobacter pylori*. *Nature* **445**, 915–918 (2007).
152. Thorell, K. et al. Rapid evolution of distinct *Helicobacter pylori* subpopulations in the Americas. *PLoS Genet.* **13**, e1006546 (2017).
153. Gardy, J. L. et al. Whole-genome sequencing and social-network analysis of a tuberculosis outbreak. *N. Engl. J. Med.* **364**, 730–739 (2011).
154. Gregory, A. C. et al. Marine DNA viral macro- and microdiversity from pole to pole. *Cell* **177**, 1109–1123 (2019).
155. Arevalo, P., VanInsberghe, D., Elsherbini, J., Gore, J. & Polz, M. F. A reverse ecology approach based on a biological definition of microbial populations. *Cell* **178**, 820–834 (2019).
156. Garcia, S. L. et al. Contrasting patterns of genome-level diversity across distinct co-occurring bacterial populations. *ISME J.* **12**, 742–755 (2018).
157. Kopac, S. et al. Genomic heterogeneity and ecological speciation within one subspecies of *Bacillus subtilis*. *Appl. Environ. Microbiol.* **80**, 4842–4853 (2014).
158. Levy, R. & Borenstein, E. in *Evolutionary Systems Biology* Vol. 751 (ed. Soyer, O. S.) 329–345 (Springer, 2012).
159. Burghardt, L. T. et al. Select and resequence reveals relative fitness of bacteria in symbiotic and free-living environments. *Proc. Natl Acad. Sci. USA* **115**, 2425–2430 (2018).
160. Wang, J. & Jia, H. Metagenome-wide association studies: fine-mining the microbiome. *Nat. Rev. Microbiol.* **14**, 508–522 (2016).
161. Knight, R. et al. Best practices for analysing microbiomes. *Nat. Rev. Microbiol.* **16**, 410–422 (2018).
162. Song, W., Wemheuer, B., Zhang, S., Steensen, K. & Thomas, T. MetaChIP: community-level horizontal gene transfer identification through the combination of best-match and phylogenetic approaches. *Microbiome* **7**, 36 (2019).
163. Seiler, E., Trappe, K. & Renard, B. Y. Where did you come from, where did you go: refining metagenomic analysis tools for horizontal gene transfer characterisation. *PLoS Comput. Biol.* **15**, e1007208 (2019).
164. Douglas, G. M. & Langille, M. G. I. Current and promising approaches to identify horizontal gene transfer events in metagenomes. *Genome Biol. Evol.* **11**, 2750–2766 (2019).
165. Cox, C. B., Moore, P. D. & Ladle, R. J. (eds) *Biogeography: An Ecological and Evolutionary Approach*. (Wiley-Blackwell, 2016).
166. Arora, D., Singh, A., Sharma, V., Bhaduria, H. S. & Patel, R. B. HgsDb: haplogroups database to understand migration and molecular risk assessment. *Bioinformatics* **11**, 272–275 (2015).
167. Cantino, P. & de Queiroz, K. PhyloCode: a phylogenetic code of biological nomenclature. *PhyloCode*. www.ohio.edu/phyloCode (2010).
168. Tenover, F. C. et al. Interpreting chromosomal DNA restriction patterns produced by pulsed-field gel electrophoresis: criteria for bacterial strain typing. *J. Clin. Microbiol.* **33**, 2233–2239 (1995).
169. Schlöter, M., Leubhn, M., Heulin, T. & Hartmann, A. Ecology and evolution of bacterial microdiversity. *FEMS Microbiol. Rev.* **24**, 647–660 (2000).
170. Hamilton, M. *Population Genetics* (Wiley-Blackwell, 2009).
171. Cohan, F. M. Transmission in the origins of bacterial diversity, from ecotypes to Phyla. *Microbiol. Spectr.* **5**, <https://doi.org/10.1128/microbiolspec.MTBP-0014-2016> (2017).
172. Alkan, C., Coe, B. P. & Eichler, E. E. Genome structural variation discovery and genotyping. *Nat. Rev. Genet.* **12**, 363–376 (2011).
173. Kaper, J. B., Nataro, J. P. & Mobley, H. L. T. Pathogenic *Escherichia coli*. *Nat. Rev. Microbiol.* **2**, 123–140 (2004).
174. Samuel, B. *Medical Microbiology* (Univ. of Texas Medical Branch, 1996).
175. Kenneth, R., George, R. & Sherris, J. C. (eds) *Medical Microbiology: An Introduction to Infectious Diseases* (McGraw-Hill Medical, 2004).
176. Houghton Mifflin Company. *The American Heritage Medical Dictionary - Serovar*. (Houghton Mifflin, 2007).
177. Silva, N. A. et al. Genomic diversity between strains of the same serotype and multilocus sequence type among pneumococcal clinical isolates. *Infect. Immun.* **74**, 3513–3518 (2006).
178. Fratomico, P. M. et al. Advances in molecular serotyping and subtyping of *Escherichia coli*. *Front. Microbiol.* **7**, 644 (2016).
179. Miller-Keane & Marie, O. *Miller-Keane Encyclopedia and Dictionary of Medicine, Nursing, and Allied Health* 7th edn. (W. B. Saunders, 2003).
180. diCenzo, G. C. & Finan, T. M. The divided bacterial genome: structure, function, and evolution. *Microbiol. Mol. Biol. Rev.* **81**, e00019-17 (2017).
181. Hamady, M. & Knight, R. Microbial community profiling for human microbiome projects: Tools, techniques, and challenges. *Genome Res.* **19**, 1141–1152 (2009).
182. Nocker, A., Burr, M. & Camper, A. K. Genotypic microbial community profiling: a critical technical review. *Microb. Ecol.* **54**, 276–289 (2007).
- Reviews foundational methods that enabled microbial diversity to be assessed directly within a microbial community, sometimes at within-species resolution.**
183. Eren, A. M., Boris, G. G., Huse, S. M. & Mark Welch, J. L. Oligotyping analysis of the human oral microbiome. *Proc. Natl Acad. Sci. USA* **111**, E2875–E2884 (2014).
184. Eren, A. M. et al. Minimum entropy decomposition: unsupervised oligotyping for sensitive partitioning of high-throughput marker gene sequences. *ISME J.* **9**, 968–979 (2015).
185. Callahan, B. J. et al. DADA2: high-resolution sample inference from illumina amplicon data. *Nat. Methods* **13**, 581–583 (2016).
186. Amir, A. et al. Deblur rapidly resolves single-nucleotide community sequence patterns. *mSystems* **2**, e00191-16 (2017).
187. Tikhonov, M., Leach, R. W. & Wingreen, N. S. Interpreting 16S metagenomic data without clustering to achieve sub-OTU resolution. *ISME J.* **9**, 68–80 (2015).
188. Johnson, J. S. et al. Evaluation of 16S rRNA gene sequencing for species and strain-level microbiome analysis. *Nat. Commun.* **10**, 5029 (2019).
189. Nielsen, H. B. et al. Identification and assembly of genomes and genetic elements in complex metagenomic samples without using reference genomes. *Nat. Biotechnol.* **32**, 822–828 (2014).
190. Yu, F. B. et al. Microfluidic-based mini-metagenomics enables discovery of novel microbial lineages from complex environmental samples. *eLife* **6**, e26580 (2017).
191. Shi, X. et al. Microfluidics-based enrichment and whole-genome amplification enable strain-level resolution for airway metagenomics. *mSystems* **4**, e00198-19 (2019).
192. Bowers, R. M. et al. Minimum information about a single amplified genome (MISAG) and a metagenome-assembled genome (MIMAG) of bacteria and archaea. *Nat. Biotechnol.* **35**, 725–731 (2017).
- Establishes minimal quality reporting requirements for MAGs.**
193. Almeida, A. et al. A new genomic blueprint of the human gut microbiota. *Nature* **568**, 499–504 (2019).
194. Beitel, C. W. et al. Strain- and plasmid-level deconvolution of a synthetic metagenome by sequencing proximity ligation products. *PeerJ* **2**, e415 (2014).
195. Costea, P. I. et al. Towards standards for human fecal sample processing in metagenomic studies. *Nat. Biotechnol.* **35**, 1069–1076 (2017).
196. Shaiber, A. & Eren, A. M. Composite metagenome-assembled genomes reduce the quality of public genome repositories. *mBio* <https://doi.org/10.1128/mBio.00725-19> (2019).
- Provides an example of how assembling genomes from metagenomes (creating MAGs) can lead to poor quality genomic data and why these genomes should not be considered the same as genomes from isolates.**

197. Schmidt, T. S. B., Raes, J. & Bork, P. The human gut microbiome: from association to modulation. *Cell* **172**, 1198–1215 (2018).
Reviews the known connections between human gut microbiome and health, including discussion of strain-level variation.
198. Salter, S. J. et al. Reagent and laboratory contamination can critically impact sequence-based microbiome analyses. *BMC Biol.* **12**, 87 (2014).
199. Goldstein, S., Beka, L., Graf, J. & Klassen, J. L. Evaluation of strategies for the assembly of diverse bacterial genomes using MinION long-read sequencing. *BMC Genomics* **20**, 23 (2019).
200. Alneberg, J. et al. Genomes from uncultivated prokaryotes: a comparison of metagenome-assembled

- and single-amplified genomes. *Microbiome* **6**, 173 (2018).
201. Gawad, C., Koh, W. & Quake, S. R. Single-cell genome sequencing: current state of the science. *Nat. Rev. Genet.* **17**, 175–188 (2016).

Acknowledgements

Funding for research in the authors' laboratories was provided by the European Research Council (ERC) (grant ERC-AdG-669830 MicrobioS), the European Union's Horizon 2020 Research and Innovation Programme (grant 825694 MICROB-PREDICT), the German Federal Ministry of Education and Research (Bundesministerium für Bildung und Forschung, BMBF) (grant 01GL1746B PRIMAL) and the European Molecular Biology Laboratory (EMBL).

Author contributions

The authors contributed equally to all aspects of the article.

Competing interests

The authors declare no competing interests.

Peer review information

Nature Reviews Microbiology thanks C. Quince and the other, anonymous, reviewer(s) for their contribution to the peer review of this work.

Publisher's note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

© Springer Nature Limited 2020