

Sequencing the human microbiome in health and disease

Journal:	<i>Human Molecular Genetics</i>
Manuscript ID:	HMG-2013-D-00971.R1
Manuscript Type:	4 Invited Review Article
Date Submitted by the Author:	n/a
Complete List of Authors:	Cox, Michael; Imperial College London, National Heart and Lung Institute Cookson, William; Imperial College London, National Heart and Lung Institute Moffatt, Miriam; Imperial College, National Heart and Lung Institute
Key Words:	human microbiome, 16S rRNA, metagenomics, bacteria

SCHOLARONE™
Manuscripts

Review

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Sequencing the human microbiome in health and disease

Michael J Cox*, William OCM Cookson and Miriam F Moffatt

Molecular Genetics and Genomics Section, National Heart and Lung Institute,
Imperial College London, Dovehouse Street, London, SW3 6LY

*corresponding author
Tel: +44 207 351 8730
Fax: +44 207 351 8126
Email: michael.cox1@imperial.ac.uk

For Peer Review

Abstract

Molecular techniques have revolutionised the practise of standard microbiology. In particular, 16S rRNA sequencing, whole microbial genome sequencing and metagenomics are revealing the extraordinary diversity of microorganisms on Earth and their vast genetic and metabolic repertoire.

The increase in length, accuracy and number of reads generated by high-throughput sequencing has coincided with a surge in interest in the human microbiota, the totality of bacteria associated with the human body, in both health and disease. Traditional views of host/pathogen interactions are being challenged as the human microbiota are being revealed to be important in normal immune system function, to diseases not previously thought to have a microbial component and to infectious diseases with unknown aetiology.

In this review, we introduce the nature of the human microbiota and application of these three key sequencing techniques for its study, highlighting both advances and challenges in the field. We go on to discuss how further adoption of additional techniques, also originally developed in environmental microbiology, will allow the establishment of disease causality against a background of numerous, complex and interacting microorganisms within the human host.

Introduction

Microorganisms associated with the human body have been studied for many years in both health and disease. The first Human Microbiome Project perhaps began when Antonie van Leeuwenhoek scraped “gritty matter” from between his teeth and became the first to visualise bacteria, or “animalcules” in dental plaque in 1683 (1).

Since then, research on human-associated microorganisms has, for the pragmatic reason of combating infectious disease in human, veterinary and agricultural settings, tended to focus on pathogens. In addition, human associated-microbe research has been limited because many bacteria are difficult to grow in the laboratory. Now techniques pioneered in environmental microbiology are being applied to human diseases and are revealing complex interactions between microorganisms themselves and with their human hosts. This holds promise for a new understanding of infectious disease and for diseases not previously recognised to have a microbial component.

With the advent of high-throughput sequencing substantial numbers of samples can be processed rapidly and cost effectively. These technological advances have led to an interest in the human as a super-organism made up of interacting human and microbial components. Such interactions may be complex and occur at many levels that extend well beyond the traditional models of host-pathogen and immune-virulence. Five to 8% of the human genome, for example, consists of endogenous retroviruses (2); gut bacteria may increase the risk of cardiovascular disease by the metabolic degradation of L-carnitine (3); and the gut microbiota may confer good health in the elderly by as yet unknown mechanisms (4). Thus many aspects of human well-being may be influenced by our associated, integrated and ubiquitous microbiota.

In this review we will introduce three key techniques in the field and speculate on implications for future research in human disease.

The human microbiota

Microbiome literally means small biome, the ecosystem comprising all microorganisms in a particular environment together with their genes and environmental interactions. The assemblage of microorganisms themselves is referred to as the microbiota or microbial community and can include bacteria,

1
2
3 archaea, viruses, phage, fungi and other microbial eukarya. Microflora can also be
4 found as a general term in the literature, but *flora* refers specifically to plants, rather
5 than microbes and so the alternative terms are preferable. The human microbiota
6 consists of microorganisms that exist upon, within or in close proximity to the human
7 body. Bacteria are the most well studied group of microorganisms in this context but
8 archaea, viruses and eukarya such as fungi also account for a high proportion of the
9 human microflora (5, 6, 7). Parts of the body with significant bacterial populations
10 include the gut and the oral cavity, (8, 9), the skin (10), urogenital tracts (11, 12), the
11 nasopharynx (13) and body sites canonically considered sterile such as the lower
12 respiratory tract (14, 15).

13
14
15
16
17
18
19
20
21 Organisms in these locations may be present stably over long periods (16) and may be
22 considered endemic or may be transient (17). Composition of the microbiota tends to
23 be defined by the body site (18) indicating the presence of different selection
24 pressures. For example, moisture is an important driver of the community structure
25 found on skin (19). In addition, the human microbiota has been shown to be
26 individual, so twins share a similar but non-identical microbiota (20) and the use of
27 microbiota analyses in forensic analysis has been suggested (21).

28
29
30
31
32
33
34
35 In designing studies of the human microbiota, variability over time and the
36 individuality of composition necessitates longitudinal sampling, as an individual in
37 their healthy state is their own best disease control. Large, cross-sectional study
38 cohorts are nonetheless important when longitudinal sampling is not feasible.

39
40
41
42
43 Acquisition of the human microbiota is believed to start at birth (22). Bacteria may be
44 present in amniotic fluid (23) whilst the bacteria in the newborn gut have been shown
45 to be influenced by delivery mode with the microbiota sourced from the mother's
46 vagina during delivery or from skin with caesarean section (24). Accumulation of
47 further organisms continues during infancy and childhood as demonstrated by
48 longitudinal studies of the gut (8) and cross-sectional studies of the respiratory tract
49 (25).

50
51
52
53
54
55
56 It is difficult to use culture methods to isolate more than a small percentage of the
57 microorganisms known to be present in any environment. Isolating a wider range of
58
59
60

1
2
3 organisms to be able to study them in detail takes considerable effort (26). The
4 molecular techniques described below offer an alternative means to gather substantial
5 detail about individual organisms and entire communities, whilst circumventing the
6 selection biases inherent in isolation by culture (Figure 1).
7
8
9

10 11 **16S rRNA gene sequencing.**

12 16S rRNA gene sequencing has been the first molecular tool to be generally applied
13 to the human microbiota. It gives a quantitative description of the bacteria present in a
14 complex biological mixture, allowing investigation of whole communities and the
15 identities of their constituent members.
16
17
18
19

20
21 The small subunit ribosomal or 16S rRNA gene (*rrnA*) is a highly conserved
22 component of the transcriptional machinery of all DNA-based life forms.
23 Phylogenetic mapping of *rrnA* variation was first used to establish the three domain
24 structure of life (27). The conserved nature of the gene was subsequently exploited to
25 develop more rapid methods for determining relationships between organisms directly
26 from environmental DNA and RNA extracts (28, 29).
27
28
29
30
31

32
33 The 16S rRNA gene consists of conserved and variable regions (Figure 2A). The
34 variable regions allow discrimination between different microorganisms. 16S rRNA
35 gene methods rely on the PCR (polymerase chain reaction) using “universal” primers
36 targeted at the conserved regions and designed to amplify as wide a range of different
37 microorganisms as possible (30). This is followed by assaying the amplified fragment
38 of the gene, originally with molecular fingerprinting approaches such as denaturing
39 gradient gel electrophoresis (31) or by cloning and sequencing of the PCR products
40 (32).
41
42
43
44
45
46
47

48 A number of different phylogenetic microarrays, consisting of multiple probes
49 designed to discriminate sub-groups of organisms have also been used (33), (34). The
50 probes and the organisms targeted must be pre-selected, limiting the utility of arrays
51 to detect novel organisms.
52
53
54
55

56 The coupling of 16S rRNA PCR with next-generation sequencing makes possible the
57 study of many samples at low cost (35). One of the most significant limitations of 16S
58
59
60

1
2
3 rRNA sequencing is the introduction of biases by PCR primer designs, which may
4 select for or against particular groups of organisms (36, 37). An under-recognised
5 problem with the use of the PCR is that bacterial contamination of reagents is
6 common and requires extensive controls at many levels of the process (38). The 16S
7 rRNA operon is also present in between one and fifteen copies in bacterial genomes,
8 which may additionally influence the apparent relative abundance of an organism
9 (39).
10
11
12
13
14
15

16 Analysis of 16S rRNA data relies on the clustering of related sequences at a particular
17 level of identity and counting the number of representatives of each cluster. Since
18 molecular methods for determining the identity of bacteria and archaea do not map
19 directly to the classic biochemically determined taxonomies, clusters of similar
20 sequences are referred to as operational taxonomic units (OTUs). OTU counts are
21 summarised in a table of relative abundances for each organism in each sample and
22 these tables are used for downstream analyses. A level of 97 % sequence identity is
23 frequently chosen as being representative of a species and 95 % for a genus when
24 using partial 16S rRNA gene sequences. Some bacteria and archaea will only be
25 identified at the genus or family, rather than species level (40).
26
27
28
29
30
31
32
33

34 Accuracy of identification is dependent on the reference database chosen. Five
35 percent of the 16S rRNA sequences in GenBank (the NIH sequence database) may be
36 erroneous (41). Curated databases such as The Ribosomal Database Project (42),
37 GreenGenes (43) and SILVA (44), where sequences undergo quality assessment and
38 alignments are manually optimised, are crucial for optimal phylogenetic placement of
39 test sequences. Two analysis pipelines are in common use for analysing 16S rRNA
40 gene sequence data: QIIME (45) and Mothur (46), though there is no standardised
41 way of applying these pipelines to datasets.
42
43
44
45
46
47
48

49 An important deliverable of 16S rRNA gene sequencing is the identification of
50 microorganisms that cause disease. Current microbial diagnostics provide information
51 about the presence or absence of known pathogens in patient samples, but the culture-
52 based techniques are much more targeted and selective when compared to 16S rRNA
53 gene sequencing (47). More than 50% of cases of pneumonia in children and adults
54
55
56
57
58
59
60

1
2
3 requiring hospitalisation have no diagnosis. DNA sequencing therefore has the
4 immediate potential to fill a major unmet clinical need.
5
6

7
8 Although identification and characterisation of disease-causing organisms is the
9 ultimate goal (see whole genome sequencing below), measures of the microbial
10 community structure, such as species richness, community evenness and diversity,
11 can reveal a great deal about dynamics and selection pressures experienced by the
12 system (Figure 2B and Table 1).
13
14
15
16

17
18 Association of these parameters with relevant environmental and clinical
19 measurements can give important insight into states of health and disease (48, 49).
20 Increased richness, evenness and diversity can be associated with stable, longer
21 established or less active ecosystems (50). Microbial community stability, resistance
22 to environmental pressures such as diet and antibiotic use, and resistance to invasion
23 with pathogens are also likely to be important in human disease states affecting the
24 bowel, mouth, lungs, skin and vagina (51).
25
26
27
28
29
30

31 **Whole genome sequencing (WGS)**

32 Complete genome sequencing is the foundation for the comprehensive understanding
33 of an organism's function. Bacteria were the first free-living organisms to undergo
34 complete genome sequencing, with *Haemophilus influenzae* being completed in 1995
35 (52). As of July 2013, the National Center for Biotechnology Information's microbial
36 genome site listed 2552 complete genomes, although the bacteria sequenced in their
37 entirety have been highly selected, with multiple genomes of commonly cultured
38 clinical strains and an absence of some entire phyla (53).
39
40
41
42
43
44
45

46 It is now feasible to map all the genes that characterise a particular group of
47 organisms, their 'pangenomes', by sequencing a broad range of isolates from different
48 sources (54). This reveals the genes that are core and that define the genomes of a
49 particular group as well as those that are accessory, perhaps possessed by a single
50 isolate or a subset with a particular lifestyle or pathology. It also allows inference of
51 how pathogenicity evolves within lineages of organisms that consist of both
52 pathogenic and non-pathogenic organisms.
53
54
55
56
57
58
59
60

1
2
3 Molecular epidemiology, where isolates associated with a particular disease outbreak
4 are typed and tracked to pinpoint their source also benefits from WGS. Multi-locus
5 sequence typing (MLST) relies on the sequencing of multiple house-keeping genes
6 for a particular species (55). This can reveal important details about transmission and
7 sources of an outbreak that can inform future responses, although it may suffer from a
8 lack of resolution (56). Epidemiological studies have almost exclusively been
9 conducted on stored isolates after an outbreak, limiting their usefulness in altering the
10 course of an outbreak as it occurs.

11
12
13
14
15
16
17
18 High-throughput sequencing may bridge these gaps as it becomes feasible to generate
19 whole genome sequences faster than the outbreak source can be tracked. Sequencing
20 may deliver detailed information about virulence factors, antibiotic resistance and
21 strain source quickly enough to influence outbreak responses. The additional
22 genomic information allows much higher resolution of strains than MLST. In 2011 a
23 shigatoxigenic *Escherichia coli* O104:H4 outbreak occurred in Germany and isolates
24 were sequenced, data released and the annotation crowd-sourced within a week (57).
25 Application of WGS to a neonatal methicillin resistant *Staphylococcus aureus*
26 outbreak revealed a transmission event missed by other techniques and did so in
27 clinically relevant time-scales (58).

28
29
30
31
32
33
34
35
36 Annotation of sequence data and comparative genomics is onerous and not currently
37 within the realm of clinical diagnostic or epidemiological laboratories. Raw sequence
38 data must undergo quality control prior to assembly of the curated reads as a
39 representative genome. Coding regions and their putative functions are identified
40 during annotation and then the assembled sequence deposited in public databases.
41 Finally, the new genome may be considered in the context of previously sequenced
42 microorganisms (59, 60). Tools for facile annotation and analysis are in development,
43 and may allow translation of the methodology into routine clinical practice (61).
44 Importantly, whole genome sequencing still requires the isolation of the organisms
45 concerned which is not always feasible.

54 55 **Metagenomics**

56 Metagenomics describes the direct sequencing of the total DNA extracted from a
57 microbial community and combines elements of the above two approaches.
58
59
60

1
2
3
4
5 Identification of organisms present is improved relative to 16S rRNA gene
6 sequencing, and organisms such as phage and viruses that do not have a phylogenetic
7 marker gene can be assessed (62). This is tempered by an increased sequencing effort
8 to detect less common organisms, relative to 16S rRNA gene sequencing. Functional
9 annotation allows the comparison of the physiological capabilities between
10 communities and environmental conditions (63). With sufficient representation of the
11 organisms present it may be possible to reconstruct the metabolic and biogeochemical
12 pathways between microorganisms and to use this to direct isolation attempts (64, 65),
13 gaining insight into the fundamental functioning of the ecosystem.
14
15
16
17
18
19

20
21 Metagenomics is beginning to be applied to the human microbiota. The healthy
22 human gut has been postulated to consist of three metagenomically defined
23 enterotypes, typified by relative dominance of particular groups of organisms:
24 *Prevotella*, *Ruminococcus* and *Bacteroides* spp. (66). Start-up companies and
25 research projects are now offering to do for your gut microbiota what direct to
26 consumer (DTC) genetics companies do for the human genome, classifying your
27 faecal sample into an enterotype. It is unclear, however, what the gut enterotype might
28 mean to the consumer, as other studies have suggested continuous measures of
29 community structure are more representative of gut microbial diversity and that
30 enterotypes may be artefacts of the analytic methods employed (67).
31
32
33
34
35
36
37
38
39

40 More recently metagenome wide association studies (MWAS) have begun to emerge.
41 In a study of Type 2 diabetes (T2D), although microbial composition of the gut varied
42 between individuals, it did not vary greatly between cases and controls. Some
43 functional differences were observed in the gut microbiota, with a relative decline in
44 butyrate-producing bacteria found in T2D (68). Butyrate is a key compound in
45 host/microbiota metabolism in the gut, being a bacterial metabolite produced by key
46 strains such as *Roseburia* spp. and *Faecalibacterium prausnitzii* (69, 70) as well as a
47 preferred carbon source for gut epithelium cells (71). For atherosclerosis MWAS has
48 shown depletion of the same organisms in cases of disease, accompanied by genes
49 involved in peptidoglycan synthesis (72).
50
51
52
53
54
55
56
57

58 **Future directions**

59
60

1
2
3 Many of the existing molecular studies of the human microbiota are hypothesis
4 generating and associate particular OTUs or organisms with clinical measurements.
5 These may yield useful biomarkers of disease but fall short generally of establishing
6 true causality. Causality may be established by borrowing further techniques from
7 environmental microbial ecology. For example, 16S rRNA sequencing, WGS and
8 metagenomic studies may be combined with approaches with that allow activity of the
9 community to be measured, including metabolomics (73), metaproteomics (74) and
10 metatranscriptomics (75).
11
12
13
14
15
16
17

18 Little is known about the consistency of the microbiota at particular sites in different
19 populations and environments. It is likely that genetic and epidemiological
20 approaches that include systematic quantification of microbial communities will be
21 able to identify the host factors that support healthy microbial communities as well as
22 those that predispose to disease.
23
24
25
26
27

28 Fluorescent *in situ* hybridisation (FISH) is useful for direct visualisation and
29 identification of microorganisms in human tissue (76). This is important since DNA
30 and RNA extracts often use samples such as sputum and stool, which do not give
31 good information on the localisation and spatial heterogeneity of the organisms.
32 Techniques such as stable isotope probing (77) are also available that simultaneously
33 reveal function, identity and activity of microorganisms in close to *in situ* conditions.
34 Finally, single cell genomics shows promise as a technique that effectively allows
35 isolation of an organism, without culture (78).
36
37
38
39
40
41
42

43 Treatments targeted at the microbiota for the moment tend to focus on probiotics,
44 introduced bacteria such as *Lactobacillus* spp., or prebiotics, compounds that enrich
45 the growth of particular organisms deemed to be of benefit (79). Faecal transplants
46 are perhaps amongst the more surprising methods for amending microbial
47 communities. In the case of chronic *Clostridium difficile* infection transplanting a
48 faecal slurry from a healthy donor to the patient's bowel has demonstrated good
49 efficacy (80), although the importance of pre-screening the transplant community for
50 potential pathogens should be emphasised.
51
52
53
54
55
56
57
58
59
60

1
2
3 The pervasive role of the human microbiota in health and disease is still largely
4 unexplored. Three hundred years after van Leeuwenhoek first visualised his dental
5 calculus, only 50 % of the microorganisms present in the oral microbiota can be
6 cultivated. This is the highest proportion of any of the human host niches (81). The
7 challenge remains for human microbial ecologists to establish causal relationships
8 between the microbiota and the human host in the varied microbial niches of the body
9 and varying disease states, with the ultimate goal of translation into real clinical
10 benefit.
11
12
13
14
15
16
17

18 **Acknowledgements**

19 We thank two anonymous reviewers for their constructive comments. The authors are
20 supported by a Wellcome Trust Senior Investigator Award to WOCMC and MFM.
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

References

1. Van Leuwenhoek,A. (1683) An Extract of a Letter from Mr. Anth. Van Leuwenhoek, concerning Animalcules Found on the Teeth; Of the Scaleyness of the Skin, &c. *Phil. Trans. (1683-1775)*, 646–649.
2. Belshaw,R., Pereira,V., Katzourakis,A., Talbot,G., Paces,J., Burt,A. and Tristem,M. (2004) Long-term reinfection of the human genome by endogenous retroviruses. *Proc. Natl. Acad. Sci. U.S.A.*, **101**, 4894–4899.
3. Koeth,R.A., Wang,Z., Levison,B.S., Buffa,J.A., Org,E., Sheehy,B.T., Britt,E.B., Fu,X., Wu,Y., Li,L., et al. (2013) Intestinal microbiota metabolism of L-carnitine, a nutrient in red meat, promotes atherosclerosis. *Nat. Med.*, **19**, 576–585.
4. Claesson,M.J., Cusack,S., O'Sullivan,O., Greene-Diniz,R., de Weerd,H., Flannery,E., Marchesi,J.R., Falush,D., Dinan,T., Fitzgerald,G., et al. (2011) Composition, variability, and temporal stability of the intestinal microbiota of the elderly. *Proc. Natl. Acad. Sci. U.S.A.*, **108 Suppl. 1**, 4586–4591.
5. Duncan,S.H., Louis,P. and Flint,H.J. (2007) Cultivable bacterial diversity from the human colon. *Lett. Appl. Microbiol.*, **44**, 343–350.
6. Delwart,E. (2013) A Roadmap to the Human Virome. *PLoS Path.*, **9**, e1003146.
7. Parfrey,L.W., Walters,W.A. and Knight,R. (2011) Microbial eukaryotes in the human microbiome: ecology, evolution, and future directions. *Front. Microbio.l*, **2**, 153, 1-6.
8. Yatsunenko,T., Rey,F.E., Manary,M.J., Trehan,I., Dominguez-Bello,M.G., Contreras,M., Magris,M., Hidalgo,G., Baldassano,R.N., Anokhin,A.P., et al. (2012) Human gut microbiome viewed across age and geography. *Nature*, **486**, 222–227.
9. Wade,W.G. (2013) Pharmacological Research. *Pharmacol. Res.*, **69**, 137–143.
10. Grice,E.A., Kong,H.H., Renaud,G., Young,A.C., NISC Comparative Sequencing Program, Bouffard,G.G., Blakesley,R.W., Wolfsberg,T.G., Turner,M.L. and Segre,J.A. (2008) A diversity profile of the human skin microbiota. *Genome Res.*, **18**, 1043–1050.
11. Price,L.B., Liu,C.M., Johnson,K.E., Aziz,M., Lau,M.K., Bowers,J., Ravel,J., Keim,P.S., Serwadda,D., Wawer,M.J., et al. (2010) The Effects of Circumcision on the Penis Microbiome. *PLoS ONE*, **5**, e8422.
12. Ravel,J., Gajer,P., Abdo,Z., Schneider,G.M., Koenig,S.S.K., McCulle,S.L., Karlebach,S., Gorle,R., Russell,J., Tacket,C.O., et al. (2011) Vaginal microbiome of reproductive-age women. *Proc. Natl. Acad. Sci. U.S.A.*, **108 Suppl 1**, 4680–4687.
13. Bogaert,D., Keijsers,B., Huse,S., Rossen,J., Veenhoven,R., van Gils,E., Bruin,J., Montijn,R., Bonten,M. and Sanders,E. (2011) Variability and Diversity of Nasopharyngeal Microbiota in Children: A Metagenomic Analysis. *PLoS ONE*, **6**,

- 1
2
3 e17035.
4
5 14. Hilty,M., Burke,C., Pedro,H., Cardenas,P., Bush,A., Bossley,C., Davies,J.,
6 Ervine,A., Poulter,L., Pachter,L., et al. (2010) Disordered microbial communities
7 in asthmatic airways. *PLoS ONE*, **5**, e8578.
8
9 15. Charlson,E.S., Bittinger,K., Haas,A.R., Fitzgerald,A.S., Frank,I., Yadav,A.,
10 Bushman,F.D. and Collman,R.G. (2011) Topographical Continuity of Bacterial
11 Populations in the Healthy Human Respiratory Tract. *Am. J. of Resp. and Crit.*
12 *Care Med.*, **184**, 957–963.
13
14 16. Faith,J.J., Guruge,J.L., Charbonneau,M., Subramanian,S., Seedorf,H.,
15 Goodman,A.L., Clemente,J.C., Knight,R., Heath,A.C., Leibel,R.L., et al. (2013)
16 The Long-Term Stability of the Human Gut Microbiota. *Science*, **341**, 1237439–
17 1237439.
18
19 17. Caporaso,J.G., Lauber,C.L., Costello,E.K., Berg-Lyons,D., Gonzalez,A.,
20 Stombaugh,J., Knights,D., Gajer,P., Ravel,J., Fierer,N., et al. (2011) Moving
21 pictures of the human microbiome. *Genome Biol.*, **12**, R50, 1-8.
22
23 18. Consortium,T.H.M.P. (2013) Structure, function and diversity of the healthy
24 human microbiome. *Nature*, **486**, 207–214.
25
26 19. Findley,K., Oh,J., Yang,J., Conlan,S., Deming,C., Meyer,J.A., Schoenfeld,D.,
27 Nomicos,E., Park,M., NIH Intramural Sequencing Center Comparative
28 Sequencing Program, et al. (2013) Topographic diversity of fungal and bacterial
29 communities in human skin. *Nature*, **498**, 367–370.
30
31 20. Turnbaugh,P.J., Hamady,M., Yatsunenکو,T., Cantarel,B.L., Duncan,A., Ley,R.E.,
32 Sogin,M.L., Jones,W.J., Roe,B.A., Affourtit,J.P., et al. (2009) A core gut
33 microbiome in obese and lean twins. *Nature*, **457**, 480–484.
34
35 21. Fierer,N., Lauber,C.L., Zhou,N., McDonald,D., Costello,E.K. and Knight,R.
36 (2010) From the Cover: Forensic identification using skin bacterial communities.
37 *Proc. Natl. Acad. Sci. U.S.A.*, **107**, 6477–6481.
38
39 22. Palmer,C., Bik,E.M., DiGiulio,D.B., Relman,D.A. and Brown,P.O. (2007)
40 Development of the human infant intestinal microbiota. *PLoS Biol.*, **5**, e177.
41
42 23. Han,Y.W., Redline,R.W., Li,M., Yin,L., Hill,G.B. and McCormick,T.S. (2004)
43 *Fusobacterium nucleatum* Induces Premature and Term Stillbirths in Pregnant
44 Mice: Implication of Oral Bacteria in Preterm Birth. *Inf. and Imm.*, **72**, 2272–
45 2279.
46
47 24. Dominguez-Bello,M.G., Costello,E.K., Contreras,M., Magris,M., Hidalgo,G.,
48 Fierer,N. and Knight,R. (2010) Delivery mode shapes the acquisition and
49 structure of the initial microbiota across multiple body habitats in newborns.
50 *Proc. Natl. Acad. Sci. U.S.A.*, **107**, 11971–11975.
51
52 25. Cardenas,P.A., Cooper,P.J., Cox,M.J., Chico,M., Arias,C., Moffatt,M.F. and
53 Cookson,W.O. (2012) Upper Airways Microbiota in Antibiotic-Naïve Wheezing
54 and Healthy Infants from the Tropics of Rural Ecuador. *PLoS ONE*, **7**, e46803.
55
56
57
58
59
60

- 1
2
3 26. Tunney, M.M., Klem, E.R., Fodor, A.A., Gilpin, D.F., Moriarty, T.F., McGrath, S.J.,
4 Muhlebach, M.S., Boucher, R.C., Cardwell, C., Doering, G., et al. (2011) Use of
5 culture and molecular analysis to determine the effect of antibiotic treatment on
6 microbial community diversity and abundance during exacerbation in patients
7 with cystic fibrosis. *Thorax*, **66**, 579–584.
8
9
10 27. Woese, C.R. and Fox, G.E. (1977) Phylogenetic structure of the prokaryotic
11 domain: the primary kingdoms. *Proc. Natl. Acad. Sci. U.S.A.*, **74**, 5088–5090.
12
13 28. Lane, D.J., Pace, B., Olsen, G.J., Stahl, D.A., Sogin, M.L. and Pace, N.R. (1985)
14 Rapid determination of 16S ribosomal RNA sequences for phylogenetic analyses.
15 *Proc. Natl. Acad. Sci. U.S.A.*, **82**, 6955–6959.
16
17 29. Olsen, G.J., Lane, D.J., Giovannoni, S.J., Pace, N.R. and Stahl, D.A. (1986)
18 Microbial Ecology and Evolution: A Ribosomal RNA Approach. *Annu. Rev.*
19 *Microbiol.*, **40**, 337–365.
20
21 30. Lane, D.J. (1991) 16S/23S rRNA sequencing. In Stackebrandt, E., Goodfellow, M.
22 (eds), *Nucleic acid techniques in bacterial systematics*.
23
24 31. Muyzer, G., de Waal, E.C. and Uitterlinden, A.G. (1993) Profiling of complex
25 microbial populations by denaturing gradient gel electrophoresis analysis of
26 polymerase chain reaction-amplified genes coding for 16S rRNA. *App. and*
27 *Environ. Microbiol.*, **59**, 695–700.
28
29
30 32. Eckburg, P.B. (2005) Diversity of the Human Intestinal Microbial Flora. *Science*,
31 **308**, 1635–1638.
32
33 33. Desantis, T.Z., Brodie, E.L., Moberg, J.P., Zubietta, I.X., Piceno, Y.M. and
34 Andersen, G.L. (2007) High-Density Universal 16S rRNA Microarray Analysis
35 Reveals Broader Diversity than Typical Clone Library When Sampling the
36 Environment. *Microb. Ecol.*, **53**, 371–383.
37
38 34. Rajilić-Stojanović, M., Heilig, H.G.H.J., Molenaar, D., Kajander, K., Surakka, A.,
39 Smidt, H. and de Vos, W.M. (2009) Development and application of the human
40 intestinal tract chip, a phylogenetic microarray: analysis of universally conserved
41 phylotypes in the abundant microbiota of young and elderly adults. *Environ.*
42 *Microbiol.*, **11**, 1736–1751.
43
44
45 35. Metzker, M.L. (2009) Sequencing technologies—the next generation. *Nat. Rev.*
46 *Genetics*, **11**, 31–46.
47
48 36. Sim, K., Cox, M.J., Wopereis, H., Martin, R., Knol, J., Li, M.-S., Cookson, W.O.C.M.,
49 Moffatt, M.F. and Kroll, J.S. (2012) Improved Detection of Bifidobacteria with
50 Optimised 16S rRNA-Gene Based Pyrosequencing. *PLoS ONE*, **7**, e32543.
51
52
53 37. Klindworth, A., Pruesse, E., Schweer, T., Peplies, J., Quast, C., Horn, M. and
54 Glockner, F.O. (2012) Evaluation of general 16S ribosomal RNA gene PCR
55 primers for classical and next-generation sequencing-based diversity studies.
56 *Nucleic Acids Res.*, **41**, e1, 1–14.
57
58
59
60

- 1
2
3 38. Tanner, MA., Goebel, BM., Dojka, MA., and Pace, NR. (1998) Specific ribosomal
4 DNA sequences from diverse environmental settings correlate with experimental
5 contaminants. *Appl Environ. Microbiol.* 1998, **64**, 3110–3113
6
7
8 39. Klappenbach, JA., Saxman, P.R., Cole, JR and Schmidt, TM. (2001) rrndb: the
9 ribosomal RNA operon copy number database. *Nucleic Acids Res.*, **29**, 181-184
10
11
12 40. Stackebrandt, E. and Goebel, B.M. (1994) Taxonomic note: a place for DNA-DNA
13 reassociation and 16S rRNA sequence analysis in the present species definition in
14 bacteriology. *Int. J. Syst. Bacteriol.*, **44**, 846–849.
15
16 41. Ashelford, K.E., Chuzhanova, N.A., Fry, J.C., Jones, A.J. and Weightman, A.J.
17 (2005) At least 1 in 20 16S rRNA sequence records currently held in public
18 repositories is estimated to contain substantial anomalies. *App. and Environ.*
19 *Micro.*, **71**, 7724–7736.
20
21
22 42. Cole, JR., Wang, Q., Cardenas, E., Fish, J., Chai, B., Farris, R.J., Kulam-Syed-
23 Mohideen, AS., McGarrell, DM., Marsh, T., Garrity, GM. *et al.* (2009) The
24 Ribosomal Database Project: improved alignments and new tools for rRNA
25 analysis. *Nucleic Acids Res.*, **37** (suppl. 1) D141-D145.
26
27 43. McDonald, D., Price, M.N., Goodrich, J., Nawrocki, E.P., Desantis, T.Z., Probst, A.,
28 Andersen, G.L., Knight, R. and Hugenholtz, P. (2011) An improved Greengenes
29 taxonomy with explicit ranks for ecological and evolutionary analyses of bacteria
30 and archaea. *ISME J.*, **6**, 610–618.
31
32 44. Quast, C., Pruesse, E., Yilmaz, P., Gerken, J., Schweer, T., Yarza, P., Peplies, J. and
33 Glockner, F.O. (2012) The SILVA ribosomal RNA gene database project:
34 improved data processing and web-based tools. *Nucleic Acids Res.*, **41**, D590–
35 D596.
36
37 45. Caporaso, J.G., Kuczynski, J., Stombaugh, J., Bittinger, K., Bushman, F.D.,
38 Costello, E.K., Fierer, N., Peña, A.G., Goodrich, J.K., Gordon, J.I., et al. (2010) QIIME
39 allows analysis of high-throughput community sequencing data. *Nat. Meth.*, **7**,
40 335–336.
41
42 46. Schloss, P.D., Westcott, S.L., Ryabin, T., Hall, J.R., Hartmann, M., Hollister, E.B.,
43 Lesniewski, R.A., Oakley, B.B., Parks, D.H., Robinson, C.J., et al. (2009) Introducing
44 mothur: Open-Source, Platform-Independent, Community-Supported Software
45 for Describing and Comparing Microbial Communities. *Applied and Environ.*
46 *Microbiol.*, **75**, 7537–7541.
47
48 49
50 47. Duff, R.M., Simmonds, N.J., Davies, J.C., Wilson, R., Alton, E.W., Pantelidis, P.,
51 Cox, M.J., Cookson, W.O.C.M., Bilton, D. and Moffatt, M.F. (2013) A molecular
52 comparison of microbial communities in bronchiectasis and cystic fibrosis. *Eur.*
53 *Respir. J.*, **41**, 991–993.
54
55 48. Cox, M.J., Allgaier, M., Taylor, B., Baek, M.S., Huang, Y.J., Daly, R.A., Karaoz, U.,
56 Andersen, G.L., Brown, R., Fujimura, K.E., et al. (2010) Airway Microbiota and
57 Pathogen Abundance in Age-Stratified Cystic Fibrosis Patients. *PLoS ONE*, **5**,

- 1
2
3 e11044.
4
5 49. Ley, R.E., Bäckhed, F., Turnbaugh, P., Lozupone, C.A., Knight, R.D. and Gordon, J.I.
6 (2005) Obesity alters gut microbial ecology. *Proc. Natl. Acad. Sci. U.S.A.*, **102**,
7 11070–11075.
8
9 50. Legendre, P. and Legendre, L. (2012) Numerical Ecology 3rd ed. Elsevier.
10
11 51. Shade, A., Peter, H., Allison, S.D., Baho, D.L., Berga, M., Bürgmann, H.,
12 Huber, D.H., Langenheder, S., Lennon, J.T., Martiny, J.B.H., et al. (2012)
13 Fundamentals of microbial community resistance and resilience. *Front.*
14 *Microbiol.*, **3**, 417.
15
16 52. Fleischmann, R.D., Adams, M.D., White, O., Clayton, R.A., Kirkness, E.F.,
17 Kerlavage, A.R., Bult, C.J., Tomb, J.F., Dougherty, B.A. and Merrick, J.M. (1995)
18 Whole-genome random sequencing and assembly of *Haemophilus influenzae* Rd.
19 *Science*, **269**, 496–512.
20
21 53. Wu, D., Hugenholtz, P., Mavromatis, K., Pukall, R., Dalin, E., Ivanova, N.N.,
22 Kunin, V., Goodwin, L., Wu, M., Tindall, B.J., et al. (2009) A phylogeny-driven
23 genomic encyclopaedia of Bacteria and Archaea. *Nature*, **462**, 1056–1060.
24
25 54. Pallen, M.J. and Wren, B.W. (2007) Bacterial pathogenomics. *Nature*, **449**, 835–
26 842.
27
28 55. Maiden, M.C., Bygraves, J.A., Feil, E., Morelli, G., Russell, J.E., Urwin, R.,
29 Zhang, Q., Zhou, J., Zurth, K., Caugant, D.A., et al. (1998) Multilocus sequence
30 typing: a portable approach to the identification of clones within populations of
31 pathogenic microorganisms. *Proc. Natl. Acad. Sci. U.S.A.*, **95**, 3140–3145.
32
33 56. Roetzer, A., Diel, R., Kohl, T.A., Rückert, C., Nübel, U., Blom, J., Wirth, T.,
34 Jaenicke, S., Schuback, S., Rüsche-Gerdes, S., et al. (2013) Whole Genome
35 Sequencing versus Traditional Genotyping for Investigation of a *Mycobacterium*
36 *tuberculosis* Outbreak: A Longitudinal Molecular Epidemiological Study. *PLoS*
37 *Med.*, **10**, e1001387.
38
39 57. Rohde, H., Qin, J., Cui, Y., Li, D., Loman, N.J., Hentschke, M., Chen, W., Pu, F.,
40 Peng, Y., Li, J., et al. (2011) Open-source genomic analysis of Shiga-toxin-
41 producing *E. coli* O104:H4. *N. Engl. J. Med.*, **365**, 718–724.
42
43 58. Köser, C.U., Holden, M.T.G., Ellington, M.J., Cartwright, E.J.P., Brown, N.M.,
44 Ogilvy-Stuart, A.L., Hsu, L.Y., Chewapreecha, C., Croucher, N.J., Harris, S.R., et al.
45 (2012) Rapid whole-genome sequencing for investigation of a neonatal MRSA
46 outbreak. *N. Engl. J. Med.*, **366**, 2267–2275.
47
48 59. Edwards, D.J. and Holt, K.E. (2013) Beginner's guide to comparative bacterial
49 genome analysis using next-generation sequence data. *Microb. Inform. Exp.*, **3**, 2.
50
51 60. Loman, N.J., Constantinidou, C., Chan, J.Z.M., Halachev, M., Sergeant, M.,
52 Penn, C.W., Robinson, E.R. and Pallen, M.J. (2012) High-throughput bacterial
53 genome sequencing: an embarrassment of choice, a world of opportunity. *Nature*
54 *Rev. Microbiol.*, **10**, 599–606.
55
56
57
58
59
60

- 1
2
3 61. Richardson,E.J. and Watson,M. (2013) The automatic annotation of bacterial
4 genomes. *Brief. Bioinformatics*, **14**, 1–12.
5
6 62. Willner,D., Haynes,M.R., Furlan,M., Hanson,N., Kirby,B., Lim,Y. W.,
7 Rainey,P.B., Schmieler,R., Youle,M., Conrad,D., et al. (2012) Case Studies of
8 the Spatial Heterogeneity of DNA Viruses in the Cystic Fibrosis Lung. *Am. J. of*
9 *Resp. Cell and Mol. Biol.*, **46**, 127–131.
10
11 63. Turnbaugh,P.J., Ridaura,V.K., Faith,J.J., Rey,F.E., Knight,R. and Gordon,J.I.
12 (2009) The effect of diet on the human gut microbiome: a metagenomic analysis
13 in humanized gnotobiotic mice. *Science Trans. Med.*, **1**, 6ra14, 1-19.
14
15 64. Tyson,G.W., Chapman,J., Hugenholtz,P., Allen,E.E., Ram,R.J., Richardson,P.M.,
16 Solovyev,V.V., Rubin,E.M., Rokhsar,D.S. and Banfield,J.F. (2004) Community
17 structure and metabolism through reconstruction of microbial genomes from the
18 environment. *Nature*, **428**, 37–43.
19
20 65. Tyson,G.W., Lo,I., Baker,B.J., Allen,E.E., Hugenholtz,P. and Banfield,J.F. (2005)
21 Genome-Directed Isolation of the Key Nitrogen Fixer *Leptospirillum*
22 *ferrodiazotrophum* sp. nov. from an Acidophilic Microbial Community. *App. and*
23 *Environ. Microbiol.*, **71**, 6319–6324.
24
25 66. Arumugam,M., Raes,J., Pelletier,E., Le Paslier,D., Yamada,T., Mende,D.R.,
26 Fernandes,G.R., Tap,J., Bruls,T., Batto,J.-M., et al. (2011) Enterotypes of the
27 human gut microbiome. *Nature*, **473**, 174–180.
28
29 67. Koren,O., Knights,D., Gonzalez,A., Waldron,L., Segata,N., Knight,R.,
30 Huttenhower,C. and Ley,R.E. (2013) A Guide to Enterotypes across the Human
31 Body: Meta-Analysis of Microbial Community Structures in Human Microbiome
32 Datasets. *PLoS Comput. Biol.*, **9**, e1002863.
33
34 68. Qin,J., Li,Y., Cai,Z., Li,S., Zhu,J., Zhang,F., Liang,S., Zhang,W., Guan,Y.,
35 Shen,D., et al. (2012) A metagenome-wide association study of gut microbiota in
36 type 2 diabetes. *Nature*, **490**, 55–60.
37
38 69. Duncan,S.H. (2006) Proposal of *Roseburia faecis* sp. nov., *Roseburia hominis* sp.
39 nov. and *Roseburia inulinivorans* sp. nov., based on isolates from human faeces.
40 *Int. J. Syst. Evol. Microbiol.*, **56**, 2437–2441.
41
42 70. Lopez-Siles, M., Khan, TM., Duncan, SH., Harmsen, HJM., Hermie, JM., Garcia-
43 Gil, J., and Flint, HJ. (2012) Cultured Representatives of Two Major Phylogroups
44 of Human Colonic *Faecalibacterium prausnitzii* Can Utilize Pectin, Uronic
45 Acids, and Host-Derived Substrates for Growth. *Applied and Environ.*
46 *Microbiol.*, **78**, 420-428
47
48 71. Cummings,J.H., Pomare,E.W., Branch,W.J., Naylor,C.P. and Macfarlane,G.T.
49 (1987) Short chain fatty acids in human large intestine, portal, hepatic and venous
50 blood. *Gut*, **28**, 1221–1227.
51
52 72. Karlsson,F.H., Fålk,F., Nookaew,I., Tremaroli,V., Fagerberg,B., Petranovic,D.,
53 Bäckhed,F. and Nielsen,J. (2012) Symptomatic atherosclerosis is associated with
54 an altered gut metagenome. *Nature Commun.*, **3**, 1245.
55
56
57
58
59
60

- 1
2
3 73. Nicholson,J.K., Holmes,E. and Wilson,I.D. (2005) Gut microorganisms,
4 mammalian metabolism and personalized health care. *Nat. Rev. Microbiol.*, **3**,
5 431–438.
6
- 7 74. Verberkmoes,N.C., Russell,A.L., Shah,M., Godzik,A., Rosenquist,M.,
8 Halfvarson,J., Lefsrud,M.G., Apajalahti,J., Tysk,C., Hettich,R.L., et al. (2008)
9 Shotgun metaproteomics of the human distal gut microbiota. *ISME J.*, **3**, 179–
10 189.
11
- 12 75. Poretsky,R.S., Hewson,I., Sun,S., Allen,A.E., Zehr,J.P. and Moran,M.A. (2009)
13 Comparative day/night metatranscriptomic analysis of microbial communities in
14 the North Pacific subtropical gyre. *Environ. Microbiol.*, **11**, 1358–1375.
15
- 16 76. Geißdörfer,W., Moos,V., Moter,A., Loddenkemper,C., Jansen,A., Tandler,R.,
17 Morguet,A.J., Fenoller,F., Raoult,D., Bogdan,C., and Schneider,T. (2012) High
18 frequency of *Tropheryma whipplei* in culture negative endocarditis. *J. Clin.*
19 *Micro.* **50**, 216-222.
20
- 21 77. Radajewski,S., Ineson,P., Parekh,N.R. and Murrell,J.C. (2000) Stable-isotope
22 probing as a tool in microbial ecology. *Nature*, **403**, 646–649.
23
- 24 78. Rinke,C., Schwientek,P., Sczyrba,A., Ivanova,NN., Anderson,IJ., Cheng,J-F.,
25 Darling,A., Malfatti,S., Swan,BK., Gies,EA., et al. (2013) Insights into the
26 phylogeny and coding potential of microbial dark matter. *Nature*, **499**, 431–437.
27
- 28 79. Gareau,M.G., Sherman,P.M. and Walker,W.A. (2010) Probiotics and the gut
29 microbiota in intestinal health and disease. *Nat. Rev. Gastroenterol. Hepatol.*, **7**,
30 503–514.
31
- 32 80. van Nood,E., Vrieze,A., Nieuwdorp,M., Fuentes,S., Zoetendal,E.G., de
33 Vos,W.M., Visser,C.E., Kuijper,E.J., Bartelsman,J.F.W.M., Tijssen,J.G.P., et al.
34 (2013) Duodenal Infusion of Donor Feces for Recurrent *Clostridium difficile*. *N.*
35 *Engl. J. Med.*, **368**, 407–415.
36
- 37 81. Wilson,M.J., Weightman,A.J. and Wade,W.G. (1997) Applications of molecular
38 ecology in the characterization of uncultured microorganisms associated with
39 human disease. *Rev. in Med. Microbiol.*, **8**, 91.
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Figure Legends

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1: A schematic demonstrating the processes for 16S rRNA gene sequencing, whole genome sequencing and metagenomics. Sample collection, DNA extraction, sequencing and sequence analysis are required in all three techniques. 16S rRNA gene sequencing and whole genome sequencing involve additional steps.

2A: The approximately 1.5 kb 16S rRNA gene of *Escherichia coli* showing the 9 variable regions that make it an ideal target as a phylogenetic marker gene.

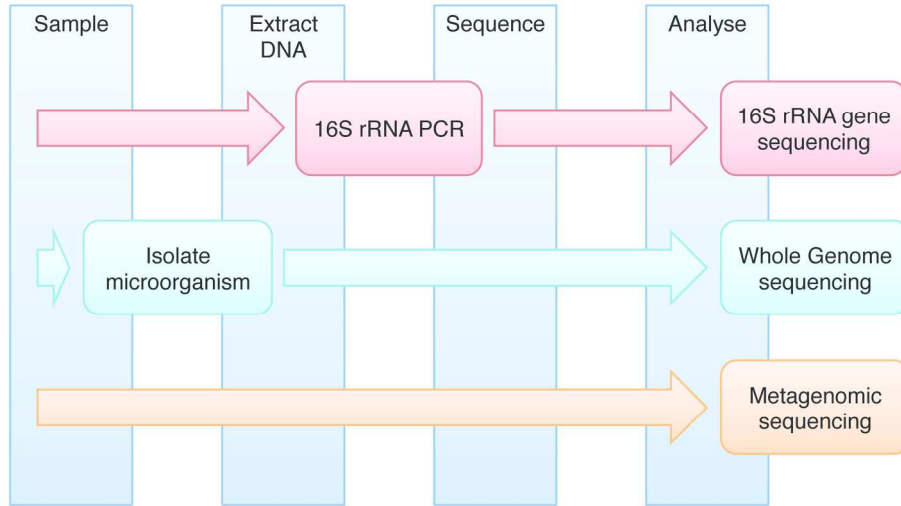
2B: A diagram demonstrating species richness and evenness and how they describe the composition of a community. Each shape represents an individual and the colour and nature of the shape represents a different type of organism. Increased numbers of different types of organism is described as increased species richness. When no one organism is dominant, the community is described as even.

Table 1: Explanation of commonly used ecological terms in the field of microbiota research.

Term	Explanation
Evenness	A measure of the skew in abundance of community members. Is there one dominant organism or are all evenly represented?
Richness	The number of different types of organism present.
Diversity	A combination of richness and evenness; can be considered to be a summary statistic for community structure as membership, abundance and evenness are taken into account.
Simpson index	A common diversity index indicating the probability that two individuals taken at random from a population are the same. Often presented as the inverse so that increasing diversity is mirrored by an increasing index value.
Shannon index	Alternatively, Shannon entropy – another common diversity index that quantifies the uncertainty of predicting the next individual taken from a sample.
Alpha diversity	Within sample diversity.
Beta diversity	Between sample diversity.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Figure 1



A schematic demonstrating the processes for 16S rRNA gene sequencing, whole genome sequencing and metagenomics. Sample collection, DNA extraction, sequencing and sequence analysis are required in all three techniques, though 16S rRNA gene sequencing and whole genome sequencing involve additional steps.

180x127mm (300 x 300 DPI)

Review

Figure 2A

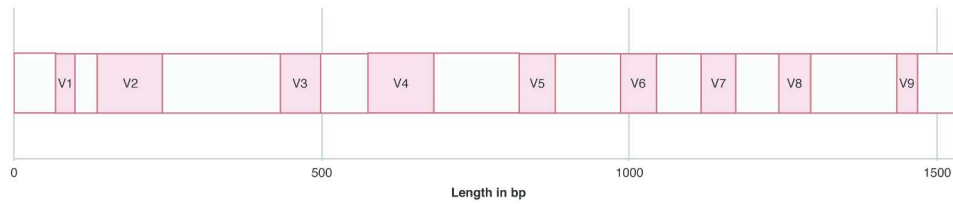
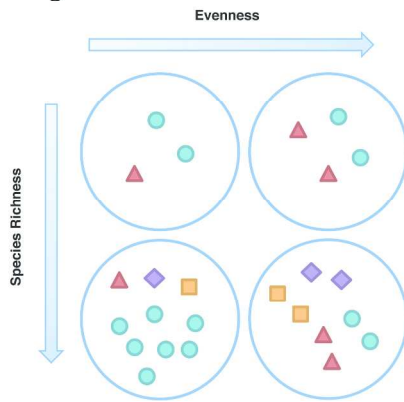


Figure 2B



2A The approximately 1.5 Kb 16S rRNA gene of *Escherichia coli*, showing the 9 variable regions that make it an ideal target as a phylogenetic marker gene.

2B A diagram demonstrating species richness and evenness and how they describe the composition of a community. Each shape represents an individual and the colour and nature of the shape represents a different type of organism. Increased numbers of different types of organism is described as increased species richness. When no one organism is dominant, the community is described as even.

254x177mm (300 x 300 DPI)