

Optimizing Microbiome Research: A Comparative Study of Laboratory Techniques and Bioinformatics Pipelines

PhD Thesis

Gábor Gulyás

Supervisor

Dr. Dóra Tombácz, PhD



**Doctoral School of Experimental and Preventive Medicine
Institute of Medical Biology
Albert Szent-Györgyi Medical School
University of Szeged**

Szeged

- 2025 -

Table of contents

Table of contents	1
List of publications.....	6
Scientific papers included in the thesis	6
Publications not related to the thesis	6
Abbreviations	9
1. Introduction	11
1.1. Microbiome	11
1.2. Gut microbiome.....	11
1.1. DNA purification techniques	12
1.2. Sequencing platform	13
1.2.1. Illumina sequencing	13
1.2.2. Single-molecule real-time (SMRT) sequencing.....	14
1.2.3. Nanopore sequencing	15
1.3. Comparison of Sequencing Platforms for Microbiome Research.....	16
1.4. Databases.....	17
2. Aims	19
3. Methods	20
3.1. Sample collection	20
3.2. DNA purification and library preparation	20
3.3. DNA quantification and quality assessment	21
3.4. Statistical analysis	21
3.5. Relative standard deviation (RSD).....	21
3.6. Sequencing	21
3.7. Bioinformatic analysis.....	21
3.7.1. Brief description of the minitax software	21
3.7.2. Analysis of the raw data	22

3.7.3.	Benchmarking datasets.....	23
3.7.4.	Databases.....	23
3.7.5.	Taxonomy version.....	24
3.7.6.	CAMISIM	24
3.7.7.	Statistics and downstream data analysis.....	24
3.7.8.	Performance metrics and statistics	25
3.7.9.	Downsampling	26
4.	Result.....	27
4.1.	Study design	27
4.2.	Comparison of DNA preparation techniques	32
4.3.	Assessment of library construction and sequencing methods.....	37
4.3.1.	Quality, yield, and reproducibility	37
4.3.2.	Microbial composition, diversity and dispersion	37
4.3.3.	Difference in ratio of Gram-positive and Gram-negative bacteria	43
4.4.	Comparison of DNA isolation methods and sequencing libraries on synthetic microbial community standards	43
4.4.1.	Quality and yield of DNA	46
4.4.2.	Microbial composition, diversity, and dispersion	46
4.4.2.1.	ZymoBIOMICS Microbial Community Standard (MCS)	46
4.4.2.2.	ZymoBIOMICS Gut Microbiome Standard (GMS)	47
4.4.3.	Difference in ratio of Gram(+) and Gram(-) bacteria	47
4.5.	Laboratory work: key findings	48
4.6.	Comparison of bioinformatics techniques.....	50
4.7.	Evaluation of minitax: benchmarking across various sequencing methods and data types	52
4.7.1.	Comparing minitax with Emu using ONT V1-V9 sequencing of MCS and GMS..	52

4.7.2.	Comparing minitax with Emu and DADA2 using Illumina V1-V2 sequencing of MCS	55
4.7.3.	Comparing minitax with sourmash using MCS data of PacBio HiFi WGS	55
4.7.4.	CAMISIM: simulated mouse gut datasets	55
5.	Discussion	56
4.	Summary	60
5.	Funding.....	61
6.	Acknowledgements	61
7.	References	62
8.	Supplementary figures.....	72
8.1.	Supplementary Figure 1	72
8.2.	Supplementary Figure 2	72
8.3.	Supplementary Figure 3	73
8.4.	Supplementary Figure 4	73
8.5.	Supplementary Figure 5	74
8.6.	Supplementary Figure 6	75
8.7.	Supplementary Figure 7	76
8.8.	Supplementary Figure 8	77
8.9.	Supplementary Figure 9	78
8.10.	Supplementary Figure 10	79
8.11.	Supplementary Figure 11	80
8.12.	Supplementary Figure 12	81
9.	Supplementary methods	82
9.1.	QIAGEN QIAamp Fast DNA Stool Mini Kit	82
9.2.	Invitrogen PureLink™ Microbiome DNA Purification Kit	82
9.3.	Macherey-Nagel NucleoSpin DNA Stool Mini kit	83
9.4.	Zymo Research Quick-DNA™ HMW MagBead Kit	84

9.5.	LIBRARY PREPARATION	84
9.5.1.	From partial regions of the 16S rRNA gene	84
9.5.1.1.	Zymo Research V1-V2.....	84
9.5.1.2.	Zymo Research V3-V4.....	85
9.5.1.3.	PerkinElmer NEXTFLEX® 16S V1-V3 Amplicon-Seq Kit for Illumina	85
9.5.2.	For the analysis of full-length 16S rRNA gene sequencing.....	86
9.5.2.1.	ONT Rapid Sequencing 16S Barcoding Kit (SQK-RAB204)	86
9.5.2.2.	PacBio Full-Length 16S Library Preparation Using SMRTbell Express Template Prep Kit 2.0 Sequel Ii System ICS v10.0 / Sequel II Chemistry 2.0 / SMRT Link v10.0	86
9.5.3.	Shotgun sequencing.....	87
9.5.3.1.	Illumina DNA Prep	87
9.6.	Data availability	87
9.7.	Code availability	88
10.	Supplementary datas.....	89
10.1.	Supplementary Data 1	89
10.2.	Supplementary Data 2	99
10.3.	Supplementary Data 3	99
10.4.	Supplementary Data 4	100
10.5.	Supplementary Data 5	101
10.6.	Supplementary Data 6	103
10.7.	Supplementary Data 7	103
10.8.	Supplementary Data 8	104
10.9.	Supplementary Data 9	104
10.10.	Supplementary Data 10	105
10.11.	Supplementary Data 11	106
10.12.	Supplementary Data 12	106
10.13.	Supplementary Data 13	107

10.14.	Supplementary Data 14	107
10.15.	Supplementary Data 15	108
10.16.	Supplementary Data 16	109
10.17.	Supplementary Data 17	109
10.18.	Supplementary Data 18	110
10.19.	Supplementary Data 19	110
10.20.	Supplementary Data 20	111
10.21.	Supplementary Data 21	111

List of publications

Scientific papers included in the thesis

1. **Gábor Gulyás** ; Balázs Kakuk* ; Ákos Dörmő* ; Tamás Járny ; István Prazsák ; Zsolt Csabai ; Miksa Máté Henkrich ; Zsolt Boldogkői ; Dóra Tombácz
Cross-comparison of gut metagenomic profiling strategies
COMMUNICATIONS BIOLOGY 7 : 1 Paper: 1445 , 22 p. (2024)
Folyóirat szakterülete: Scopus - Agricultural and Biological Sciences (miscellaneous)
SJR indikátor: D1
Nyilvános idéző összesen: 5, Független: 5, Függő: 0, Nem jelölt: 0

Publications not related to the thesis

1. Dóra Tombácz* ; Zoltán Maróti* ; Péter Oláh* ; Ákos Dörmő ; **Gábor Gulyás** ; Tibor Kalmár ; Zsolt Csabai ; Zsolt Boldogkői
Temporal transcriptional profiling of host cells infected by a veterinary alphaherpesvirus using nanopore sequencing
SCIENTIFIC REPORTS 15 : 1 Paper: 3247 , 12 p. (2025)
Scimago Journal Rank indicator: Q1
2. Dóra Tombácz ; Balázs Kakuk ; Gábor Torma ; Ádám Fülöp ; Ákos Dörmő ; **Gábor Gulyás** ; Zsolt Csabai ; Zsolt Boldogkői
Mapping the temporal transcriptomic signature of a viral pathogen through CAGE and nanopore sequencing
PLOS ONE 20 : 4 Paper: e0320439 , 28 p. (2025)
Scimago Journal Rank indicator: Q1
3. István Prazsák ; Dóra Tombácz* ; Ádám Fülöp* ; Gábor Torma ; **Gábor Gulyás** ; Ákos Dörmő ; Balázs Kakuk ; Lauren Spires McKenzie ; Zsolt Toth ; Zsolt Boldogkői
KSHV 3.0: A State-of-the-Art Annotation of the Kaposi's Sarcoma-Associated Herpesvirus Transcriptome Using Cross-Platform Sequencing
MSYSTEMS 9 : 2 Paper: e01007-23 , 19 p. (2024)
Scimago Journal Rank indicator: D1
Nyilvános idéző összesen: 10, Független: 8, Függő: 2, Nem jelölt: 0
4. Gergely Ármin Nagy ; Dóra Tombácz ; István Prazsák ; Zsolt Csabai ; Ákos Dörmő ; **Gábor Gulyás** ; Gábor Kemenesi ; Gábor E. Tóth ; Jiří Holoubek ; Daniel Růžek ; Balázs Kakuk ; Zsolt Boldogkői
Exploring the transcriptomic profile of human monkeypox virus via CAGE and native RNA sequencing approaches
MSPHERE 9 : 9 Paper: e00356-24 , 22 p. (2024)
Scimago Journal Rank indicator: Q1
5. Balázs Kakuk ; Ákos Dörmő ; Zsolt Csabai ; Gábor Kemenesi ; Jiří Holoubek ; Daniel Růžek ; István Prazsák ; Virág Éva Dani ; Béla Dénes ; Gábor Torma ; Ferenc Jakab ; Gábor E. Tóth ; Fanni V. Földes ; Brigitta Zana ; Zsófia Lanszki ; Ákos Harangozó ; Ádám Fülöp ; **Gábor Gulyás** ; Máté Mizik ; András Attila Kiss ; Dóra Tombácz ; Zsolt Boldogkői
In-depth Temporal Transcriptome Profiling of Monkeypox and Host Cells using Nanopore Sequencing
SCIENTIFIC DATA 10 : 1 Paper: 262 , 12 p. (2023)

- Scimago Journal Rank indicator: D1
Nyilvános idéző összesen: 11, Független: 10, Függő: 1, Nem jelölt: 0
6. Dóra Tombácz* ; Gábor Torma* ; **Gábor Gulyás*** ; Ádám Fülöp ; Ákos Dörmő ; István Prazsák ; Zsolt Csabai ; Máté Mizik ; Ákos Hornyák ; Zoltán Zádori ; Balázs Kakuk ; Zsolt Boldogkői
Hybrid sequencing discloses unique aspects of the transcriptomic architecture in equid alphaherpesvirus 1
HELIYON 9 : 7 Paper: e17716 , 16 p. (2023)
Scimago Journal Rank indicator: Q1
Nyilvános idéző összesen: 4, Független: 1, Függő: 3, Nem jelölt: 0
 7. Gábor Torma* ; Dóra Tombácz* ; Zsolt Csabai* ; Islam A. A. Almsarrhad ; Gergely Ármin Nagy ; Balázs Kakuk ; **Gábor Gulyás** ; Lauren Spires McKenzie ; Ishaan Gupta ; Ádám Fülöp ; Ákos Dörmő ; István Prazsák ; Máté Mizik ; Virág Éva Dani ; Viktor Csányi ; Ákos Harangozó ; Zoltán Zádori ; Zsolt Toth ; Zsolt Boldogkői
Identification of herpesvirus transcripts from genomic regions around the replication origins
SCIENTIFIC REPORTS 13 : 1 Paper: 16395 , 25 p. (2023)
Scimago Journal Rank indicator: D1
Nyilvános idéző összesen: 7, Független: 5, Függő: 2, Nem jelölt: 0
 8. See-Chi Lee ; Nenavath Gopal Naik ; Dóra Tombácz ; **Gábor Gulyás** ; Balázs Kakuk ; Zsolt Boldogkői ; Kevin Hall ; Bernadett Papp ; Steeve Boulant ; Zsolt Toth
Hypoxia and HIF-1 α promote lytic de novo KSHV infection
JOURNAL OF VIROLOGY 97 : 11 Paper: e0097223 , 18 p. (2023)
Scimago Journal Rank indicator: D1
Nyilvános idéző összesen: 4, Független: 3, Függő: 1, Nem jelölt: 0
 9. Dóra Tombácz ; Balázs Kakuk* ; Gábor Torma ; Zsolt Csabai ; **Gábor Gulyás** ; Vivien Tamás ; Zoltán Zádori ; Victoria A. Jefferson ; Florencia Meyer ; Zsolt Boldogkői
In-Depth Temporal Transcriptome Profiling of an Alphaherpesvirus Using Nanopore Sequencing
VIRUSES 14 : 6 Paper: 1289 , 25 p. (2022)
Scimago Journal Rank indicator: Q1
Nyilvános idéző összesen: 7, Független: 3, Függő: 4, Nem jelölt: 0
 10. István Prazsák ; Zsolt Csabai ; Gábor Torma ; Henrietta Papp ; Fanni Földes ; Gábor Kemenesi ; Ferenc Jakab ; **Gábor Gulyás** ; Ádám Fülöp ; Klára Megyeri ; Béla Dénes ; Zsolt Boldogkői ; Dóra Tombácz
Transcriptome dataset of six human pathogen RNA viruses generated by nanopore sequencing
DATA IN BRIEF 43 Paper: 108386 , 11 p. (2022)
Scimago Journal Rank indicator: Q2
Nyilvános idéző összesen: 2, Független: 2, Függő: 0, Nem jelölt: 0
 11. Dóra Tombácz ; Ákos Dörmő ; **Gábor Gulyás** ; Zsolt Csabai ; István Prazsák ; Balázs Kakuk ; Ákos Harangozó ; István Jankovics ; Béla Dénes ; Zsolt Boldogkői
High temporal resolution Nanopore sequencing dataset of SARS-CoV-2 and host cell RNAs
GIGASCIENCE 11 Paper: giac094 , 11 p. (2022)

- Scimago Journal Rank indicator: D1
Nyilvános idéző összesen: 3, Független: 2, Függő: 1, Nem jelölt: 0
12. Dóra Tombácz ; Norbert Moldován ; Gábor Torma ; Tibor Nagy ; Ákos Hornyák ; Zsolt Csabai ; **Gábor Gulyás** ; Miklós Boldogkői ; Victoria A. Jefferson ; Zoltán Zádori ; Florencia Meyer ; Zsolt Boldogkői
Dynamic Transcriptome Sequencing of Bovine Alphaherpesvirus Type 1 and Host Cells Carried Out by a Multi-Technique Approach
FRONTIERS IN GENETICS 12 Paper: 619056 , 8 p. (2021)
Scimago Journal Rank indicator: Q2
Nyilvános idéző összesen: 5, Független: 1, Függő: 4, Nem jelölt: 0
13. Zoltán Maróti ; Dóra Tombácz ; Norbert Moldován ; Gábor Torma ; Victoria A. Jefferson ; Zsolt Csabai ; **Gábor Gulyás** ; Ákos Dörmő ; Miklós Boldogkői ; Tibor Kalmár ; Florencia Meyer ; Zsolt Boldogkői
Time course profiling of host cell response to herpesvirus infection using nanopore and synthetic long-read transcriptome sequencing
SCIENTIFIC REPORTS 11 : 1 Paper: 14219 , 11 p. (2021)
Scimago Journal Rank indicator: D1
Nyilvános idéző összesen: 6, Független: 4, Függő: 2, Nem jelölt: 0
14. Dóra Tombácz ; Gábor Torma ; **Gábor Gulyás** ; Norbert Moldován ; Michael Snyder ; Zsolt Boldogkői
Meta-analytic approach for transcriptome profiling of herpes simplex virus type 1
SCIENTIFIC DATA 7 : 1 Paper: 223 , 11 p. (2020)
Scimago Journal Rank indicator: D1
Nyilvános idéző összesen: 7, Független: 0, Függő: 7, Nem jelölt: 0
15. Norbert Moldován ; Gábor Torma ; **Gábor Gulyás** ; Ákos Hornyák ; Zoltán Zádori ; Victoria A. Jefferson ; Zsolt Csabai ; Miklós Boldogkői ; Dóra Tombácz ; Florencia Meyer ; Zsolt Boldogkői
Time-course profiling of bovine alphaherpesvirus 1.1 transcriptome using multiplatform sequencing
SCIENTIFIC REPORTS 10 : 1 Paper: 20496 , 14 p. (2020)
Scimago Journal Rank indicator: D1
Nyilvános idéző összesen: 11, Független: 0, Függő: 11, Nem jelölt: 0
16. Dóra Tombácz ; Norbert Moldován ; Balázs Zsolt ; **Gábor Gulyás** ; Zsolt Csabai ; Miklós Boldogkői ; Michael Snyder ; Zsolt Boldogkői
Multiple Long-read Sequencing Survey of Herpes Simplex Virus Dynamic Transcriptome
FRONTIERS IN GENETICS 10 Paper: 834 , 20 p. (2019)
Scimago Journal Rank indicator: Q1
Nyilvános idéző összesen: 39, Független: 20, Függő: 19, Nem jelölt: 0

Abbreviations

I:	Invitrogen
MN:	Macherey-Nagel
Q:	Qiagen
PE:	PerkinElmer
Z:	Zymo Research
I WGS:	Illumina DNA Prep mWGS library
ONT V1-V9:	Oxford Nanopore Technologies V1-V9 amplicon library
PacBio V1-V9:	Pacific Biosciences V1-V9 amplicon library
V1-V3:	PerkinElmer V1-V3 amplicon library
V1-V2:	Zymo Research V1-V2 amplicon library
V3-V4:	Zymo Research V3-V4 amplicon library
LRS:	Long-Read Sequencing
mWGS:	Metagenomic Whole Genome Sequencing
ONT:	Oxford Nanopore Technologies
PacBio:	Pacific Biosciences
SRS:	Short-Read Sequencing
WGS:	Whole-Genome Sequencing
GMS:	ZymoBIOMICS Gut Microbiome Standard D6331
MCS:	ZymoBIOMICS Microbial Community Standard D6300
ANOVA:	Analysis of Variance
HSD:	Tukey's Honestly Significant Difference
MAPQ:	Mapping Quality
NMDS:	Non-Metric Multi-Dimensional Scaling
PCoA:	Principal Coordinate Analyses

PERMANOVA:	Permutational Analysis of Variance
PERMIDISP:	Permutational Analysis of Dispersion
db:	database
Gram-positive:	Gram+
Gram-negative:	Gram-
HMW:	High Molecular Weight
LCA:	Lowest Common Ancestor

1. Introduction

1.1. Microbiome

The roots of microbiome research can be traced back to the early 20th century, when scientists first recognized that large numbers of microorganisms, including bacteria, fungi, and viruses, inhabit different parts of the human body such as the gut, oral cavity, and skin.¹ Early culture-based methods revealed only a small fraction of microbial diversity, but these initial findings laid the foundation for modern microbiome research. Today it is known that the human microbiota contains approximately 150 times more genetic information than the entire human genome, and is therefore often referred to as the “second genome,” contributing substantially to human physiology and health.²

Over the past two decades, explosive advances in metagenomics, as well as in molecular biology, genomics, and bioinformatics, have fundamentally transformed microbiome research. These new technologies have enabled comprehensive taxonomic and functional investigation of microbial communities without the need for cultivation.^{3,4} High-throughput sequencing methods, such as Illumina, Pacific Biosciences (PacBio), and Oxford Nanopore Technologies (ONT), have elevated the study of microbial ecosystems to a new level and provided deeper insights into their metabolic, immunological, and pathophysiological roles.^{5,6}

Current research is increasingly shifting toward multi-omics approaches: the integration of metatranscriptomics, metaproteomics, and metabolomics allows for the exploration of dynamic interactions between microorganisms and their host. This paradigm shift is moving microbiome science from descriptive studies toward functional and mechanistic understanding.^{7,8} The knowledge gained in this way is not only crucial for understanding human health but also opens broad applications in agriculture⁹, biotechnology¹⁰, and environmental sciences.¹¹

1.2. Gut microbiome

The composition of the gut microbiome plays the most important role in maintaining health. The bacteria residing in the gut perform a wide range of functions, including participation in digestion, protection against pathogens, vitamin production, and stimulation of the immune system.¹² The gut microbiome is typically composed of six major phyla: *Firmicutes*, *Bacteroidetes*, *Actinobacteria*, *Proteobacteria*, *Fusobacteria*, and *Verrucomicrobia*, among which *Firmicutes* and *Bacteroidetes* dominate.¹³

The human gastrointestinal tract is the largest microbial habitat in the body, hosting approximately 10^{18} microorganisms.¹⁴ The balance of the gut microbiome is closely linked to health and disease. Several studies have demonstrated that the gut microbiota plays a key role in nutrient extraction and metabolic processes, as it possesses a wide repertoire of metabolic genes that provide enzymes and biochemical pathways for energy and nutrient harvest.¹⁵ In addition, the microbiome is a crucial contributor to the biosynthesis of bioactive molecules such as vitamins, amino acids, and lipids.¹⁶

A healthy gut microbiome is stable, resilient, and in a symbiotic relationship with the host. A healthy community is characterized by diverse species composition, a rich gene pool, and a stable “core microbiome”.¹⁷ It is important to note, however, that the relative distribution of microorganisms is unique to each individual and can change over time within the same individual. The composition of the gut microbiota is influenced by age, environmental factors (e.g., medication use), and anatomical differences across different regions of the gastrointestinal tract.¹⁷

Previous studies have demonstrated a strong similarity between the human and canine gut microbiomes.¹⁸ Dogs represent a genetically more homogeneous population, and their lifestyle and diet can be more easily controlled, making them excellent model organisms for microbiome research.¹⁹

1.1. DNA purification techniques

DNA isolation from fecal samples represents a critical initial step in microbiome studies. The aim is to extract large amounts of high-quality DNA that are suitable for subsequent sequencing and bioinformatic analyses. Over the past decades, numerous DNA isolation techniques have been developed, which differ substantially in their methodology and efficiency.

Previous studies have demonstrated that different commercially available DNA isolation kits yield variable results in terms of DNA quantity, purity, and the representation of microbial community composition.²⁰ One of the major sources of variation is whether the kit includes a mechanical cell lysis step (such as bead-beating). This step is particularly important for disrupting the cell walls of Gram-positive bacteria, thereby ensuring efficient DNA recovery from these taxa.^{21,22}

These methodological differences greatly contribute to the low reproducibility of microbiome research and the limited comparability of results across studies. The low rate of reproducibility is one of the major challenges in biomedical research and is widely referred to as the

“reproducibility crisis”.²³ Although many laboratory and bioinformatic approaches exist for processing metagenomic data, there is currently no universally accepted standard workflow that would ensure complete comparability of results. The lack of standardization remains one of the greatest obstacles to the consistent and reproducible evaluation of metagenomic studies.^{20,24,25}

1.2. Sequencing platform

In metagenomic studies the selection of the appropriate sequencing platform is of fundamental importance. Short-read sequencing (SRS) has traditionally been the most frequently used method especially in the case of metagenomic whole genome sequencing (mWGS) which is also supported by numerous bioinformatic tools such as Kraken2 and sourmash.^{26–30} In recent years however long-read sequencing (LRS) has been gaining an increasingly important role represented primarily by the platforms of Oxford Nanopore Technologies (ONT) and Pacific Biosciences (PacBio).^{31–33}

1.2.1. Illumina sequencing

Illumina technology operates on the principle of sequencing by synthesis, the central element of which is the DNA-dependent DNA polymerase enzyme.^{34,35} In this method, DNA fragments are first provided with adapter sequences, which make it possible to attach the fragments to the surface of a special glass slide (flow cell). The inner surface of the flow cell is densely coated with oligonucleotides that are complementary to the adapter sequences, so the DNA fragments can specifically bind to them.³⁶

The fixed fragments are multiplied by bridge amplification. In this process the DNA molecules "bend over" to another oligonucleotide on the surface and serve as a template for the synthesis of a new strand. By repeating this cycle thousands of identical copies are created from each original fragment, and well separated DNA clusters are formed, which can also be detected visually. The cluster-based strategy has the advantage that individual sequencing errors can be filtered out by consensus analysis of the identical copies.^{37,38}

During the actual sequencing the polymerase enzyme incorporates fluorescently labeled reversible terminator nucleotides into the growing DNA strand. Each of the four nucleotides carries a different color fluorescent label. Since the nucleotides have a terminating group at the 3' end only one nucleotide can be incorporated at a time. At the moment of incorporation the detection system senses the emitted light upon laser excitation and records which base has been added to the DNA strand. Thus in every cycle only a single base is read, in complete synchrony across all clusters.^{34,39}

At the end of the sequencing cycle the terminating group and the fluorescent label are removed through a chemical reaction, allowing DNA synthesis to continue with the next nucleotide. This "reversible termination" ensures the high accuracy of the process because after the identification of each base the cycle restarts. The detected fluorescent signals are collected as high resolution images and are then converted into nucleotide sequences by software.^{35,39}

One of the greatest advantages of Illumina technology is the extremely high level of parallelization: in a single run billions of short DNA fragments can be read simultaneously. This makes the method particularly effective in the analysis of highly complex samples such as whole genomes, metagenomes or transcriptomes.^{38,40} The short read length (typically 50–300 base pairs) however limits the method in the investigation of genomic regions containing many repetitive sequences or complex structural variations.^{40,41}

The reduction of errors is ensured by the redundancy of the clusters and precise chemical control, which is why Illumina sequencing is today the most widely used second generation sequencing technology, especially when high accuracy and enormous amounts of data are required in a relatively short time.^{34,37}

1.2.2. Single-molecule real-time (SMRT) sequencing

The SMRT technology developed by Pacific Biosciences is based on the mechanism of DNA replication, whose central enzyme is the DNA-dependent DNA polymerase.^{42,43} During the operation of this enzyme, the reading of the DNA segment to be sequenced takes place with the help of fluorescently labeled nucleotides. When a labeled nucleotide is incorporated into the forming DNA strand, the fluorescent label molecule is released from it because during synthesis a phosphate group is cleaved from the nucleotide which at the same time provides the energy requirement of the reaction. Upon release of the label molecule a short millisecond-long light flash is generated which is detected by the detector.^{42,44}

The essence of the method is that the process takes place in tiny chambers called Zero Mode Waveguide (ZMW) cells.^{43,45} In these chambers the incorporation of a single nucleotide can be observed in real time while the background signal can be suppressed since at any given moment only a single nucleotide is present. Thanks to this sequencing is faster and more accurate than in the case of second-generation technologies. In the ZMW chamber a modified DNA polymerase enzyme is immobilized which attaches to the fragment to be sequenced with the help of an adapter sequence ligated to the ends of the cDNA strands. If the adapter sequence is missing from the end of the cDNA the polymerase cannot bind therefore no information is

generated about the nucleotide sequence of the given fragment.^{43,46} Each of the four types of nucleotides is labeled with a different fluorescent molecule which ensures their unambiguous distinction.^{42,44}

For the rapid and accurate operation of the polymerase a large number of free nucleotides is required. However, this causes significant background fluorescence. Earlier technologies solved this problem by inserting a washing step after nucleotide incorporation but before detection with which the background noise could be reduced.⁴⁷ In the case of SMRT the extremely small size of the ZMW chambers allows the background signal to be separated directly during measurement. Therefore there is no need for washing steps which significantly accelerates sequencing.^{43,47}

One of the greatest advantages of the PacBio SMRT technology is the long read length typically 10–20 kb but with newer developments reads longer than 30 kb can also be achieved.^{42,46} At the same time these reads initially showed a relatively high error rate which mainly resulted from insertions and deletions as well as electrophysiological noise.^{42,44} To reduce the error rate the Circular Consensus Sequencing (CCS) approach was developed in which a consensus sequence is created from multiple readings of the same molecule.^{43,46} With this method the newer generation HiFi reads already provide Illumina-like accuracy while preserving the advantages arising from the long read length.^{46,48}

1.2.3. Nanopore sequencing

The nanopore-based sequencing developed by Oxford Nanopore Technologies (ONT) is based on measuring the direct physical properties of DNA molecules.^{49,50} The central element of the method is a biological nanopore (for example α -hemolysin or MspA porin) which is embedded into a synthetic membrane. When an electric potential is applied across the two sides of the membrane an ion current forms through the pore. As the DNA strand passes through the bases inhibit the ion current in different ways thus characteristic current fluctuations are generated from which the nucleotide sequence can be determined.^{51,52}

To start sequencing special adapter sequences are ligated to the ends of the DNA fragments which also contain the motor protein. This motor protein regulates the speed at which the DNA passes through the pore so that the current signals corresponding to individual bases are sufficiently resolved.⁵³ One significant advantage of nanopore technology is that DNA can be read in real time and directly without the need for amplification or synthesis.^{50,54}

The method is capable of producing extremely long reads typically 10–100 kb but under favorable conditions DNA molecules of several megabases in length can also be sequenced.⁵⁵ This makes it possible to accurately map complex genome regions repetitive sequences and structural variations. At the same time the error rate of the technology was initially relatively high mainly in the form of insertions and deletions.^{51,52} Developments in recent years especially basecalling algorithms based on neural networks (for example Guppy) and consensus strategies have significantly improved accuracy bringing it increasingly closer to Illumina and PacBio HiFi level.⁵⁶

A unique property of nanopore sequencing is that it can not only determine the DNA nucleotide sequence but also directly detect post-translational modifications such as DNA methylation since these modifications also cause characteristic changes in the ion current.⁵⁷ This ability makes nanopore sequencing particularly valuable in epigenetic studies.

Another special advantage is that the platform is capable of sequencing both DNA and RNA natively. In the latter case the RNA molecule itself passes directly through the nanopore so no reverse transcription or PCR is required which provides a unique opportunity to study RNA structure and modifications (for example methylations). This function is currently not offered by any other sequencing platform.⁵⁸

Another advantage of the Oxford Nanopore platform is portability the MinION device is palm-sized and can be operated from a laptop via USB port thus enabling genomic studies even in field conditions.^{54,55}

1.3. Comparison of Sequencing Platforms for Microbiome Research

The Illumina technology is the most widely used second-generation platform which provides short but extremely accurate reads. Due to its high throughput and low cost it is ideal for studies with a large number of samples as well as for metagenomic and 16S amplicon-based microbiome profiling.^{26,40} Its drawback however is that the read length typically ranges between 50–300 base pairs which makes it difficult to study complex genome regions repetitive sequences and structural variants.⁴⁰

The PacBio SMRT technology belongs to the third-generation methods and provides long reads (10–20 kb and recently >30 kb). Although the raw error rate is higher the circular consensus sequencing (HiFi reads) allows a significant increase in sequence accuracy thus the method can achieve Illumina-level accuracy.^{46,48} The advantage of long reads is particularly evident in the

study of full-length 16S rRNA genes which enables taxonomic identification even at the species level in contrast to Illumina amplicon sequencing which is limited to short regions.⁵⁹

The Oxford Nanopore Technologies (ONT) platform is also a third-generation method which stands out with the production of ultra-long reads even in the megabase range.⁵⁵ A particular feature of ONT is that it is capable of native DNA and RNA sequencing without the need for PCR amplification or reverse transcription.⁵⁸ This provides a unique opportunity for the direct investigation of RNA modifications and DNA methylation.⁵⁷ Although the error rate was previously a significant limiting factor the continuously improving basecalling algorithms and consensus-based analyses have significantly increased accuracy.⁵⁶ In addition ONT devices such as the MinION Mk1B and Mk1D provide portability thus the technology can be applied in field and clinical settings as well.⁵⁵

The differences between the three platforms are particularly emphasized in microbiome research. With Illumina sequencing genus-level information can be obtained with high accuracy and cost efficiency however species-level resolution is limited due to the short read length.^{26,40} With PacBio HiFi reads full-length 16S rRNA genes can be sequenced which allows species-level identification although the method is more expensive and has lower throughput.⁵⁹ The ONT technology provides the possibility to sequence full 16S genes and even whole genomes natively as well as to perform epigenetic and transcriptomic studies which is a unique advantage in microbiome research.^{57,58} Comparative studies have shown that Illumina sequencing provides accurate but less detailed taxonomic profiles while PacBio and ONT technologies yield significantly better species-level resolution and detection of taxonomic diversity.^{59–62}

1.4. Databases

The choice of an appropriate reference database is one of the most critical factors in metagenomic analyses, as it fundamentally determines the quality of taxonomic and functional annotation. For short-read sequencing data, widely used resources include NCBI RefSeq and GenBank, which serve as general genomic repositories.⁶³ For microbiome-focused studies, dedicated reference databases such as SILVA, Greengenes, and RDP provide curated collections of 16S rRNA gene sequences, and are therefore essential for amplicon-based community profiling.^{64–66}

One of the most recent and widely adopted resources for taxonomic classification is the Genome Taxonomy Database (GTDB), which offers a standardized phylogeny-based framework for

bacteria and archaea, enabling more accurate and consistent species-level assignments compared with traditional NCBI taxonomy.⁶⁷

For the annotation of shotgun metagenomic data, commonly used resources include the extensive pathogen and microbiome databases developed for Kraken2²⁸, as well as reference sets compatible with sourmash.³⁰ In addition, the MetaPhlAn and HUMAnN databases are widely applied, supporting species-level taxonomic profiling as well as functional annotation of metagenomes and metatranscriptomes.^{68,69}

From a functional annotation perspective, KEGG⁷⁰, eggNOG⁷¹, and UniProtKB⁷² are of particular importance, as they allow the investigation of the metabolic potential and gene functions of microbial communities.

In recent years, databases focusing specifically on the human microbiome have gained increasing significance. For examples, the Human Microbiome Project (HMP)⁴, Integrative Human Microbiome Project⁷³ and the Unified Human Gastrointestinal Genome (UHGG) collection⁶ provide high-resolution reference genomes for studying the gut microbiome. These resources enable more accurate identification of species- and strain-level differences, which are critical in both clinical and ecological research.

2. Aims

Since the microbiota plays a crucial role in human health and actively contributes to a wide range of biological processes and disease development,⁷⁴ Such methods should enable the acquisition of data that most accurately reflect biological reality. In this context, the present study aims to address the following objectives:

1. Establishing a wet lab protocol that minimizes methodological bias in gut microbiome profiling
2. Assessing the efficiency and reproducibility of DNA extraction methods
3. Evaluating the performance of library preparation protocols
4. Investigating the accuracy and potential biases of data generated by different sequencing platforms (short read and long read)
5. Designing a bioinformatics pipeline that delivers the highest possible accuracy in microbiome analysis
6. Conducting a comparative evaluation of outputs from different annotation tools
7. Assessing the effectiveness of various reference databases in microbiome research
8. Quantifying the extent to which currently accepted standard methods deviate from biological reality
9. Developing a versatile annotation tool capable of processing both short read and long read sequencing datasets

3. Methods

We applied three different sample types to evaluate the performance of four different DNA isolation kits, using six library preparation methods across three sequencing platforms. The samples originated from dog feces and from two distinct microbial community mixtures, ZymoBIOMICS Microbial Community Standard (MCS) and ZymoBIOMICS Gut Microbiome Standard (GMS), which comprised eight and eighteen bacterial strains respectively.

3.1. Sample collection

A fecal sample was collected from a 13.5-year-old healthy male Pumi (a Hungarian purebred dog) within one minute after defecation, immediately frozen, and stored at -80°C . As controls, fecal samples were obtained from six other healthy Pumis, including four puppies, a 7 year old female, and a 6.5 year old neutered male (Table 3). For each DNA isolation, a single 2.5 gram fecal sample was used, which was divided according to the requirements of each isolation kit, resulting in a total of 16 isolations with four technical replicates per kit. To avoid potential bias, the sample was not partitioned based on its internal distribution but rather randomly, and the allocation to the kits was also carried out randomly. Due to the small sample size, homogenization was not performed, which has been considered justified by Liang et al. (2020)⁷⁵, since although homogenization is essential for metabolomic analyses, it is not required in microbiome studies. For control purposes, random 200 mg subsamples were taken from six different dogs. In addition, to validate the results, we used the ZymoBIOMICS Microbial Community Standard (MCS, Zymo Research, D6300) and the ZymoBIOMICS Gut Microbiome Standard (GMS, Zymo Research, D6331) mixtures.

3.2. DNA purification and library preparation

In our study we tested the following commercially available DNA isolation kits, each with four technical replicates: QIAGEN QIAamp Fast DNA Stool Mini Kit (Ref. #51604, Lot. #169025369), Invitrogen PureLink™ Microbiome DNA Purification Kit (Cat. #423482), Macherey Nagel NucleoSpin DNA Stool Mini Kit (Ref. #740472.50, Lot. #2302 001), and Zymo Research Quick DNA™ HMW MagBead Kit (Cat. #D6060). For library preparation we applied the following kits: Illumina DNA Prep (Doc. #1000000025416 v09), Oxford Nanopore Technologies Rapid Sequencing 16S Barcoding Kit (SQK RAB204), Pacific Biosciences Full Length 16S Library Preparation Using SMRTbell Express Template Prep Kit 2.0 (PN 101 916 900), PerkinElmer NEXTFLEX® 16S V1 V3 Amplicon Seq Kit for Illumina (Cat. #4202 02), and Zymo Research Quick 16S™ NGS Library Prep Kit (Cat. #D6400). DNA isolation and library preparation were carried out according to the manufacturer protocols, which are

described in detail in the Supplementary Methods to ensure that the exact versions applied in this study remain available even if future protocol modifications are introduced.

3.3. DNA quantification and quality assessment

DNA yield was quantified using the Qubit 4.0 Fluorometer, while DNA quality including assessment of fragment length distribution was evaluated with the TapeStation 4150 system.

3.4. Statistical analysis

The data were analyzed using one way analysis of variance (ANOVA) to compare the DNA yields obtained from the four different kits. ANOVA tests were performed separately for each sample type (dog feces, MCS, and GMS). To evaluate the significance of differences between groups we applied the Tukey Honest Significant Difference (HSD) ^{76,77} post hoc test which allowed us to identify which pairs of kits differed significantly from each other.

3.5. Relative standard deviation (RSD)

The RSD was calculated to assess the internal variability of kit performance using the following formula $RSD = (\text{standard deviation} / \text{mean}) \times 100 \%$.

3.6. Sequencing

For SRS a total of nine MiSeq Reagent Kit v2 and four MiSeq Reagent Kit v2 Nano kits were used (Table 4). Sequencing of the V1–V9 region was carried out on ONT MinION and PacBio Sequel IIe platforms. ONT V1–V9 barcoded libraries prepared from dog samples were loaded onto three MinION flow cells while an additional three flow cells were used for sequencing the MCS and GMS samples. PacBio V1–V9 libraries were also barcoded and sequenced on a single Sequel SMRT Cell 8M run.

3.7. Bioinformatic analysis

3.7.1. Brief description of the minitax software

We developed *minitax*, a versatile taxonomic assignment tool designed to address the challenges posed by different types of sequencing data. *Minitax* enables taxonomic profiling across multiple sequencing platforms (ONT, PacBio, Illumina) and library types, including metagenomic whole genome sequencing (mWGS) and 16S rRNA gene sequencing. The tool employs minimap2⁷⁸ with platform specific parameter settings for initial read alignment against the reference database. Alignment results are imported into the R environment via Rsamtools⁷⁹ and integrated with database information. The workflow relies on the data.table⁸⁰ package for efficient handling of large datasets and subsequently performs post alignment processing steps to determine the best alignment for each read.

Post alignment processing involves several key steps. First, *minitax* applies a general MAPQ based filtering (for example for MAPQ values between 1 and 59) and then selects the alignment with the highest MAPQ score for each read. Next, it retains the alignment with the highest CIGAR score based on platform specific scoring matrices (Supplementary Data 8). After filtering, *minitax* determines the lowest common taxonomic level (“tax.identity”) and its corresponding rank (“tax.identity.level”) for each read. Additional refinement options are provided including the *BestAln* method which evaluates the proportion of alignments supporting a given taxonomic level. If a taxon is supported by the majority of alignments (or above a predefined default threshold of 60 percent) and is more specific than the current assignment the identification is updated accordingly. A simpler approach is offered by the *RandAln* method which randomly selects an alignment from the filtered set whereas the *SpeciesEstimate* method uses all valid alignments normalizing read counts by the number of alignments to estimate species level abundance. The *LCA* (Lowest Common Ancestor) method applies a conservative approach assigning each read to the lowest common ancestor of all matched taxa thus ensuring that assignments remain at the most specific taxonomic level supported by all alignments.

Finally *minitax* aggregates read counts at the chosen taxonomic rank and exports the results of each processing step in .tsv format while providing the final output as a *phyloseq*⁸¹ object for downstream analysis. The tool is available on GitHub at <https://github.com/Balays/minitax>. CIGAR scoring schemes and other parameters can be user defined. In this study we applied the *RandAln* method although in many cases the *BestAln* method may yield more accurate results.

3.7.2. Analysis of the raw data

For Illumina V1–V2, V3–V, and V1–V3 regions the DADA2 pipeline⁸² was used for quality control, filtering, trimming of Illumina amplicon reads, followed by dereplication, chimera detection and removal, and generation of ASVs (amplicon sequence variants), concluding with taxonomic assignment. Taxonomic classification was performed using either the SILVA 16S database⁶⁴ (version 138.1) or the Emu database (version 3.4.4). Exact parameters are provided in SUPPTAB_X (dada2_config.tsv) and in the full workflow available on GitHub (dada2.WF.R). Illumina reads were additionally processed with Emu⁸³(v3.4.4) and the in house developed *minitax* tool (v1.0) using parameters “--type sr” and “--N 10” with default settings otherwise. In both cases the default Emu database was applied, while *minitax* also included an NCBI genome collection. Both programs employ minimap2⁷⁸ for initial read mapping.

For ONT V1–V9 regions raw signal data from the MinION platform were basecalled with Guppy⁸⁴ version 6.1.5 (MinKNOW 20.05.8) using the high accuracy model. Reads were demultiplexed based on SQK RAB204 barcodes. During basecalling a minimum quality threshold of 8 was applied, separating pass and fail reads, and only pass reads were retained for downstream analysis. These pass reads were processed with Emu⁸³ (v3.4.4) and *minitax* (v1.0) using parameters “--type map-ont” and “--N 10” with default settings. Both workflows used the default Emu database, while *minitax* additionally incorporated an NCBI genome collection. The ONT EPI2ME pipeline⁸⁵ (v3.6.1) was also applied.

For PacBio V1–V9 regions basecalling, demultiplexing, and HiFi read generation were performed with PacBio SMRT Link⁸⁶ version 10.2.0.133434. High quality CCS reads were filtered with “--min-qv 20” and the resulting reads were analyzed with Emu⁸³ (v3.4.4) and *minitax* (v1.0) using parameters “--type map-pb” and “--N 10” under default settings. Both used the default Emu database, while *minitax* also employed an additional NCBI genome collection.

For Illumina WGS data raw reads were trimmed with Trim Galore⁸⁷ and host derived reads were removed with BMTagger⁸⁸ using the *Canis lupus familiaris* reference genome (GCF_014441545.1). Quality filtered reads were then processed either with the fast and sensitive taxonomic classifier *sourmash*³⁰ (v4.8.2) using the GenBank genome database (March 2022), or with *minitax* (v1.0) using default parameters and an NCBI genome collection.

3.7.3. Benchmarking datasets

In addition to the Zymo D6300 Microbial Community Mixture sequenced with the previously described ONT V1–V9 and Illumina V1–V2 methods we evaluated the performance of *minitax* on two further publicly available datasets. The Zymo MCM D6331 Microbial Community PacBio HiFi WGS dataset was used by Portik et al.⁸⁹ for benchmarking long read mWGS software. Fastq files were downloaded from NCBI (accession SRX9569057) and used as input for the *sourmash* and *minitax* programs. For both the D6300 and D6331 workflows all reads mapping to any *Veillonella* genus members were reassigned to *Veillonella rogosae* and all reads classified as belonging to the *Roseburia* genus were assigned to *Roseburia hominis*. Furthermore ten samples each from the CAMISIM simulated mouse gut dataset⁹⁰ were randomly selected for PacBio and Illumina platforms and used as input for the programs.

3.7.4. Databases

We utilized approximately 18 000 genomes representing around 13 000 bacterial archaeal and eukaryotic species. These genomes were downloaded from NCBI in February 2022 and were

used as reference databases for both WGS and 16S gene sequencing data. For amplicon sequencing datasets the choice of reference database depended on the software applied with either the Emu or the SILVA 16S database (version 138)⁶⁴ being used.

3.7.5. Taxonomy version

For the comparison of DNA isolation and library preparation techniques we used the version of NCBI that corresponded to the genomes included in the collection. In sections 2 and 3 of the results however we switched to the most recent NCBI release from July 2024. This update was not applied in other sections because the databases used for comparison with *minitax* results were based on an earlier version of the taxonomy.

3.7.6. CAMISIM

For the comparison of different DNA isolation and library preparation techniques we used the version of NCBI that corresponded to the genomes included in the collection. In sections 2 and 3 of the results however we switched to the most recent NCBI release from July 2024. This update was not applied in other sections because the databases used for comparison with *minitax* results were based on an earlier version of the taxonomy.

3.7.7. Statistics and downstream data analysis

The outputs from the different programs were organized into *phyloseq* objects which were then merged into a single comprehensive *phyloseq* object containing all metagenomic read count data. Subsequent analyses were carried out in the R environment using this *phyloseq* object together with the *tidyverse* *FactoMineR* and *vegan* packages. All scripts used for downstream data processing and figure generation are available in the GitHub repository <https://github.com/Balays/Microbiome-Method-Comparison>.

PERMANOVA (Permutational Multivariate Analysis of Variance): Significant differences in microbial composition between groups were tested with a full model PERMANOVA using the formula `adonis2(otutab_t ~ DNA_isolation_method + library, data=sampdat, method="bray")`. This model evaluated the combined effects of DNA isolation method and library preparation on microbial community structure with 999 permutations.

PERMDISP2 (Permutational Analysis of Multivariate Dispersion): Differences in within group variability were assessed with the `betadisper` function applied to the Bray–Curtis dissimilarity matrix using `betadisper(distance_matrix, groups)`.

PCoA (Principal Coordinates Analysis): PCoA was carried out on the Bray–Curtis dissimilarity matrix separately for each DNA isolation method using the output of betadisper to evaluate the effect of library preparation within each method. Distances to centroids were calculated and the results were combined across isolation methods for visualization (Fig. 5b).

NMDS (Non-metric Multidimensional Scaling): To further explore relationships between samples NMDS was performed on normalized abundance data using Bray–Curtis dissimilarity with the command `ordinate(ps.prop, method="NMDS", distance="bray")`. NMDS was used to visualize similarities between samples with the relative closeness of points reflecting the similarity of microbial community compositions (Fig. 5a).

3.7.8. Performance metrics and statistics

To comprehensively assess the performance of *minitax* we calculated precision recall F1 and F0.5 scores at species detection thresholds of 1 percent 0.1 percent and 0.01 percent following the methodology of Portik et al.⁸⁹ In addition chi square tests were performed to determine whether significant differences existed between theoretical and observed distributions and Pearson correlations were computed between theoretical and observed community compositions with corresponding r^2 values at all taxonomic levels.

We also computed the Chi-squared statistic for each taxon using the formula:

$$\chi^2 = \frac{(\text{Observed} - \text{Expected})^2}{\text{Expected}}$$

which was then summed across all taxa to assess the overall goodness-of-fit between the theoretical and observed distributions. The Chi-squared test was used to determine whether the observed distribution significantly deviated from the theoretical distribution. A Bonferroni correction was applied to account for multiple comparisons, and the null hypothesis was rejected if the corrected p -value was below the significance threshold (0.05).

In addition, we calculated Pearson's correlation coefficients (r^2 values) to evaluate the linear relationship between the theoretical and observed microbial abundances.

We also identified significantly different taxa (between the expected and observed compositions) using DESeq2⁹¹ with the following commands:

```
de.seq <- phyloseq_to_deseq2(ps, ~DNA_isolation_method)
de.seq <- DESeq(de.seq, test = "Wald", fitType = "parametric", sfType = "poscounts")
```

Taxa with a p -value < 0.05 and a log2Fold difference of ≥ 2 or ≤ -2 were considered significantly different.

3.7.9. Downsampling

Raw Illumina mWGS reads were randomly downsampled to the following read counts 2 500 000, 2 000 000, 1 500 000, 1 000 000, 750 000, 500 000, 250 000, 225 000, 200 000, 175 000, 150 000, 125 000, 100 000, 75 000, 50 000, 20 000, 10 000, and 5 000. These subsets were subsequently analyzed with Kaiju⁹² (version 1.9.0) using the progenomes (v2) database for taxonomic classification. Similarly, raw reads from non-mWGS datasets were downsampled to the same set of read counts, and processed with Emu⁸³ (version 3.4.4) under default settings.

Shannon index values were calculated from the resulting outputs using the *phyloseq*⁸¹ R package by analyzing each taxonomic abundance profile within *phyloseq*⁸¹, and the results were visualized with *ggplot2*⁹³ (<https://github.com/gabor-gulyas/Technical-article-downsample>).

4. Result

4.1. Study design

The goal of this study was to examine how different protocols influence the outcomes related to gut microbiome composition. These protocols include DNA extraction, library preparation, sequencing, and bioinformatics methods (Figs. 1, 2). We analyzed the amount, quality, and reproducibility of DNA extraction using four different isolation kits from Qiagen (Q), Macherey-Nagel (MN), Invitrogen (I), and Zymo Research (Z) (Fig. 3 and Supplementary Data 2a). We further investigated the microbial community by assessing dominant taxa and overall diversity. For mWGS, libraries were prepared with the Illumina DNA Prep Kit (I WGS), while amplicon libraries targeting the V1-V3 regions [PerkinElmer (PE V1-V3)], V1-V2, and V3-V4 [both from Zymo Research (Z V1-V2 and Z V3-V4, respectively)] of the 16S rRNA gene were also created (Supplementary Data 2b). In addition, we generated V1-V9 libraries covering the entire 16S rRNA gene for sequencing on two long-read sequencing (LRS) platforms: ONT MinION (ONT V1-V9) and PacBio Sequel IIe (PacBio V1-V9). All libraries were assessed for quality, yield, and reproducibility.

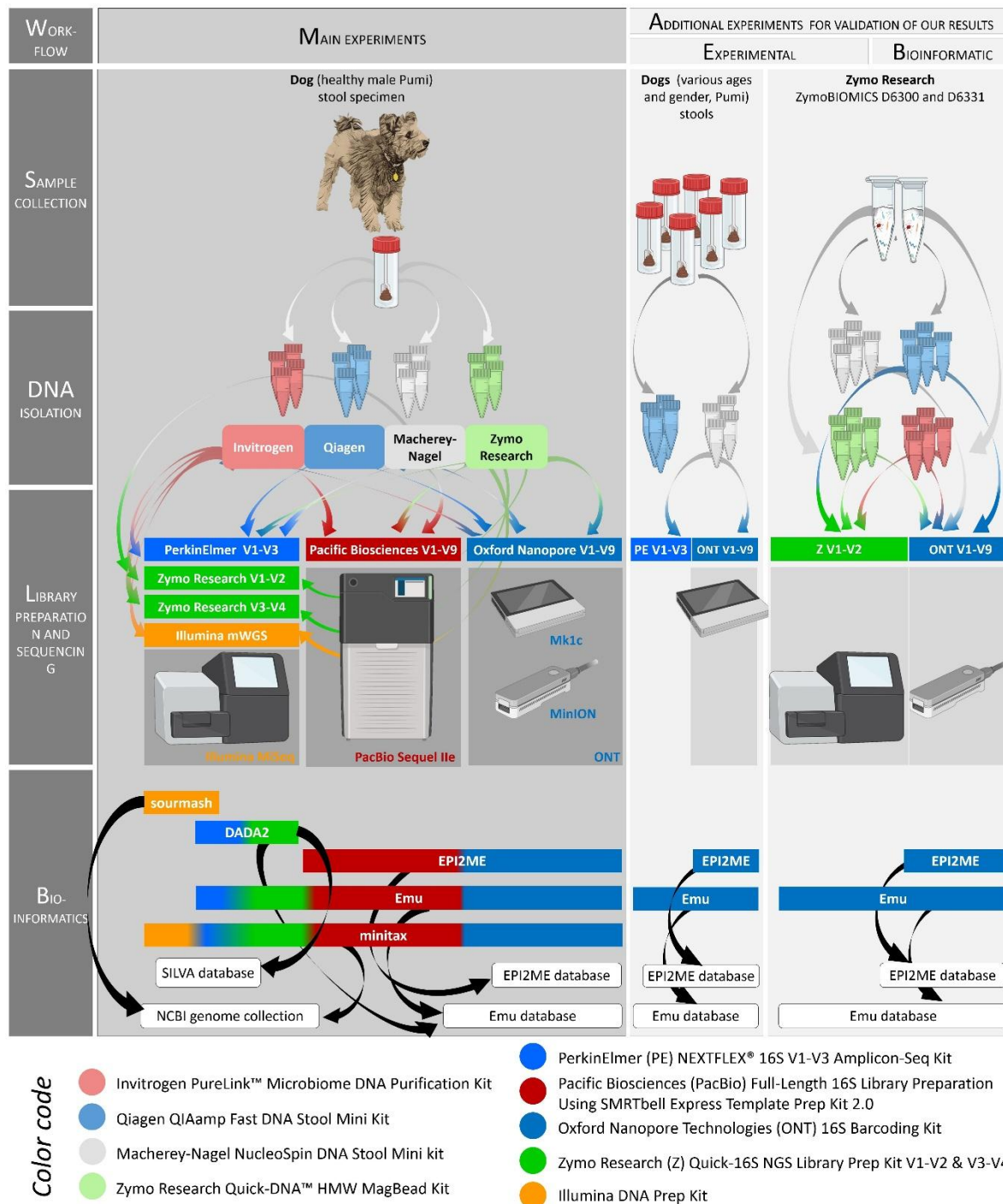


Figure 1: Workflow of the experimental part of the study.

The figure provides a detailed representation of the workflows conducted in this study. The figure illustrates the workflows conducted in this study, which included: 1) evaluating the efficacy of various DNA isolation kits in terms of quality, quantity, and microbial representation from canine stool samples; 2) comparing library preparation techniques on SRS and LRS platforms for reproducibility; 3) introducing “minitax”, a tool designed to ensure consistent analysis across multiple sequencing platforms; 4) assessing the influence of different

databases and tools on microbial profiling; and 5) comparing 16S V1-V9 sequencing on ONT and PacBio platforms to address literature gaps and emphasize bioinformatics workflows. Our goal was to identify reliable procedures for robust and reproducible gut microbiome profiling across both wet-lab and dry-lab methodologies. We performed additional experiments to validate the most extreme experimental and bioinformatics results, particularly focusing on methods that yielded the most inconsistent outcomes in comparison to other techniques. For this purpose, we utilized samples from six additional dogs of various genders and ages. The workflow involved: 1) DNA isolation using the Q kit, 2) Library preparation with the PerkinElmer V1-V3 kit, and 3) Analysis of V1-V9 libraries using the EPI2ME software. Furthermore, we carried out experiments employing a Microbial Community Standard (MCS; Zymo Research D6300) as well as a Gut Microbiome Standard (GMS; Zymo Research D6331) to validate the effectiveness of the four DNA isolation kits used. This included: 1) DNA isolation using the kits applied for the dog samples, 2) Preparation of V1-V2 and V1-V9 libraries, and 3) Sequencing on the corresponding Illumina and ONT platforms based on the library. Created in BioRender. BioRender.com/k32q619.

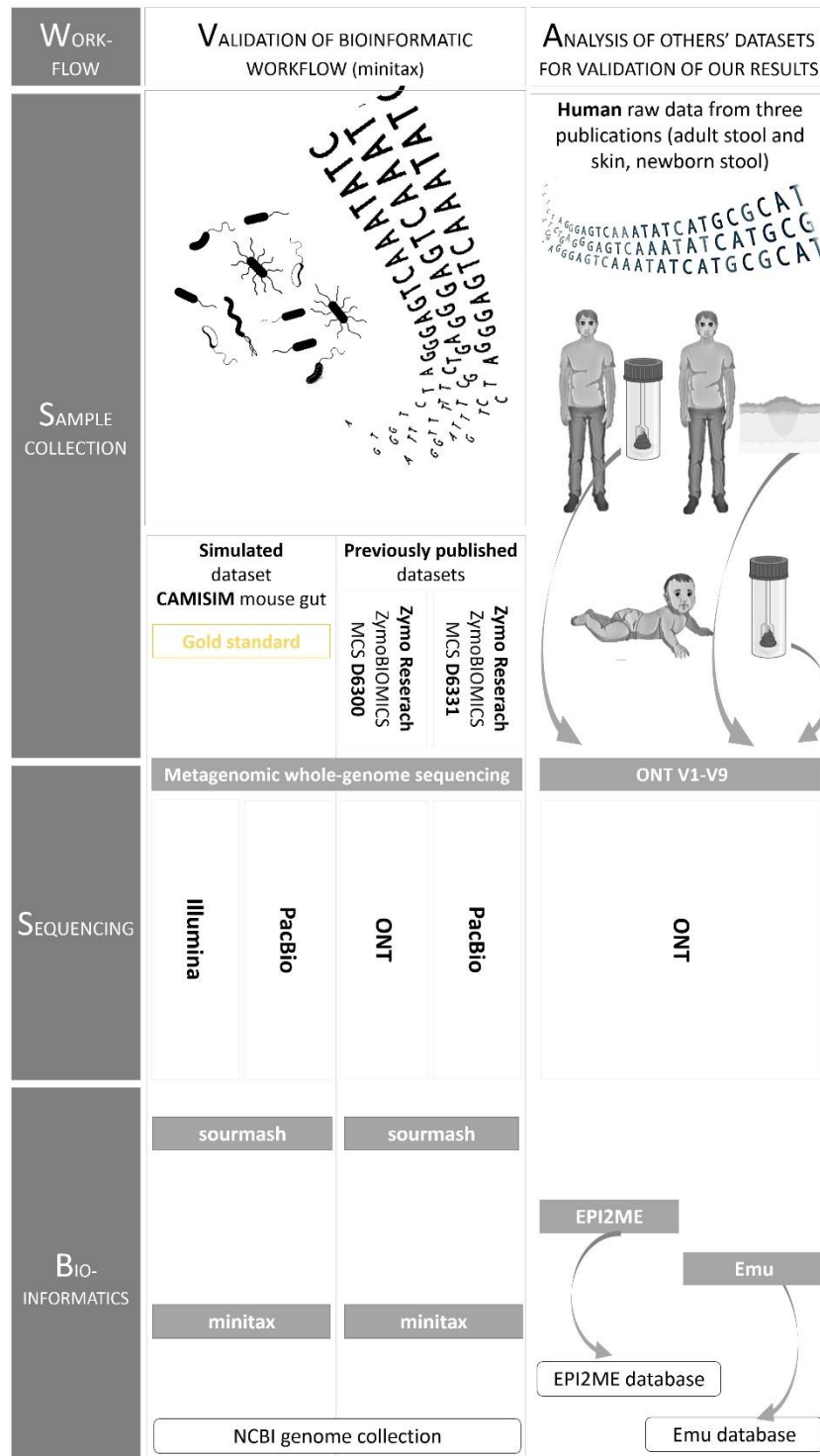


Figure 2: Workflow of the in silico analyses. Additionally, we conducted an in silico analysis to validate our bioinformatic tool, minitax. We compared the performance of minitax and sourmash by utilizing them for the analysis of: 1) simulated PacBio and Illumina data⁹⁰, and 2) previously published datasets⁸⁹. For further validation of our results, we also used previously published metagenomics data from human sources encompassing skin⁹⁴ as well as fecal samples from newborns⁹⁵ and adults⁹⁶. Created in BioRender. BioRender.com/e57v119.

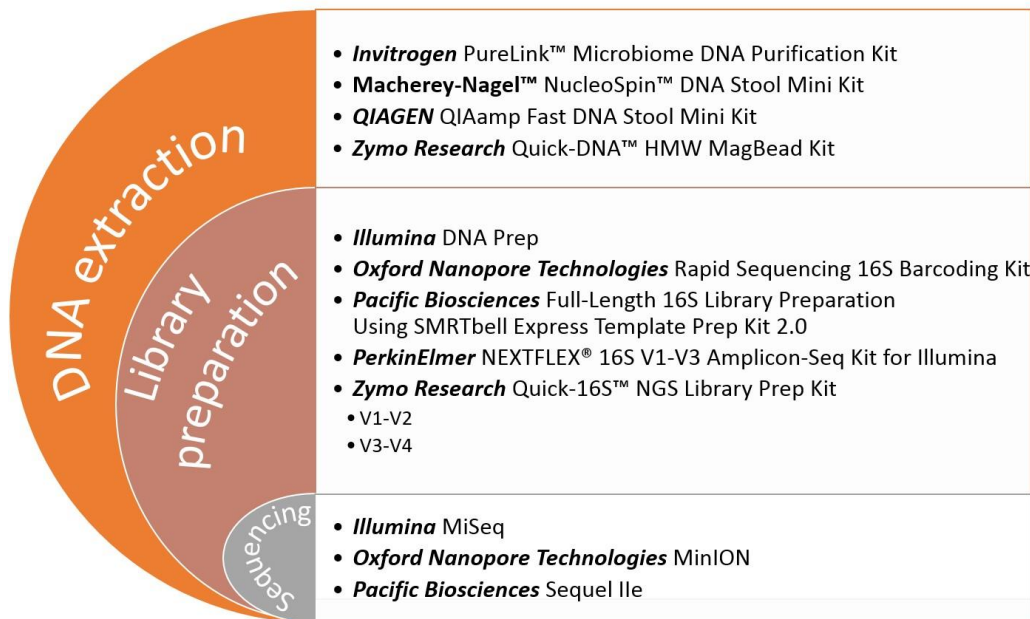


Figure 3: The kits and sequencers utilized in this project. This figure lists the DNA extraction and library preparation kits, along with sequencing devices which were used in this experiment.

The bioinformatics workflows included DADA2⁴⁴ for amplicon-based short-read sequencing (SRS) datasets, sourmash³⁰ for WGS samples, and Emu⁸³, a recently developed high-accuracy tool optimized for LRS 16S rRNA-Seq⁴⁶. For nanopore sequencing data, we also used the vendor-specific EPI2ME^{85,97} pipeline. Furthermore, we created a flexible and broadly applicable program named “minitax,” designed to process multiple data types and provide reliable taxonomic assignment across metagenomic datasets.

After sequencing reads were aligned to either a reference genome database (default option) or a 16S database using minimap2⁷⁸, minitax selected the best alignment and assigned the most likely taxonomy for each read based on mapping quality (MAPQ values) and CIGAR strings (Supplementary Data 2c).

To examine how different library preparation methods affect the accuracy of microbial composition, we standardized DNA extraction using the same isolation kit to minimize inconsistencies. We compared and validated the results by assessing various stages of the workflow with multiple sample types and datasets (Figs. 1, 2).

4.2. Comparison of DNA preparation techniques

Most DNA isolation kits rely on affinity-based purification, inhibitor removal solutions or columns, and lysis buffers, enzymes, or bead-beating for breaking cell walls. While bead-beating is frequently recommended in previous studies^{98,99}, many commercial kits do not include this step. We evaluated four commercial kits that differed in these features (Supplementary Data 2a). Among the four methods tested (Fig. 3), the Z approach required the most hands-on effort, setting it apart as the most labor-intensive option (Supplementary Fig. 1). The other three methods were more similar in this aspect. The kits used in this work showed notable differences in both yield and quality of extracted DNA (Supplementary Data 3). These differences were visible in the DNA amount (Fig. 4a, b, Supplementary Fig. 2), the ratio of microbial to host DNA (Fig. 4c, d), and reproducibility (Fig. 4e, f).

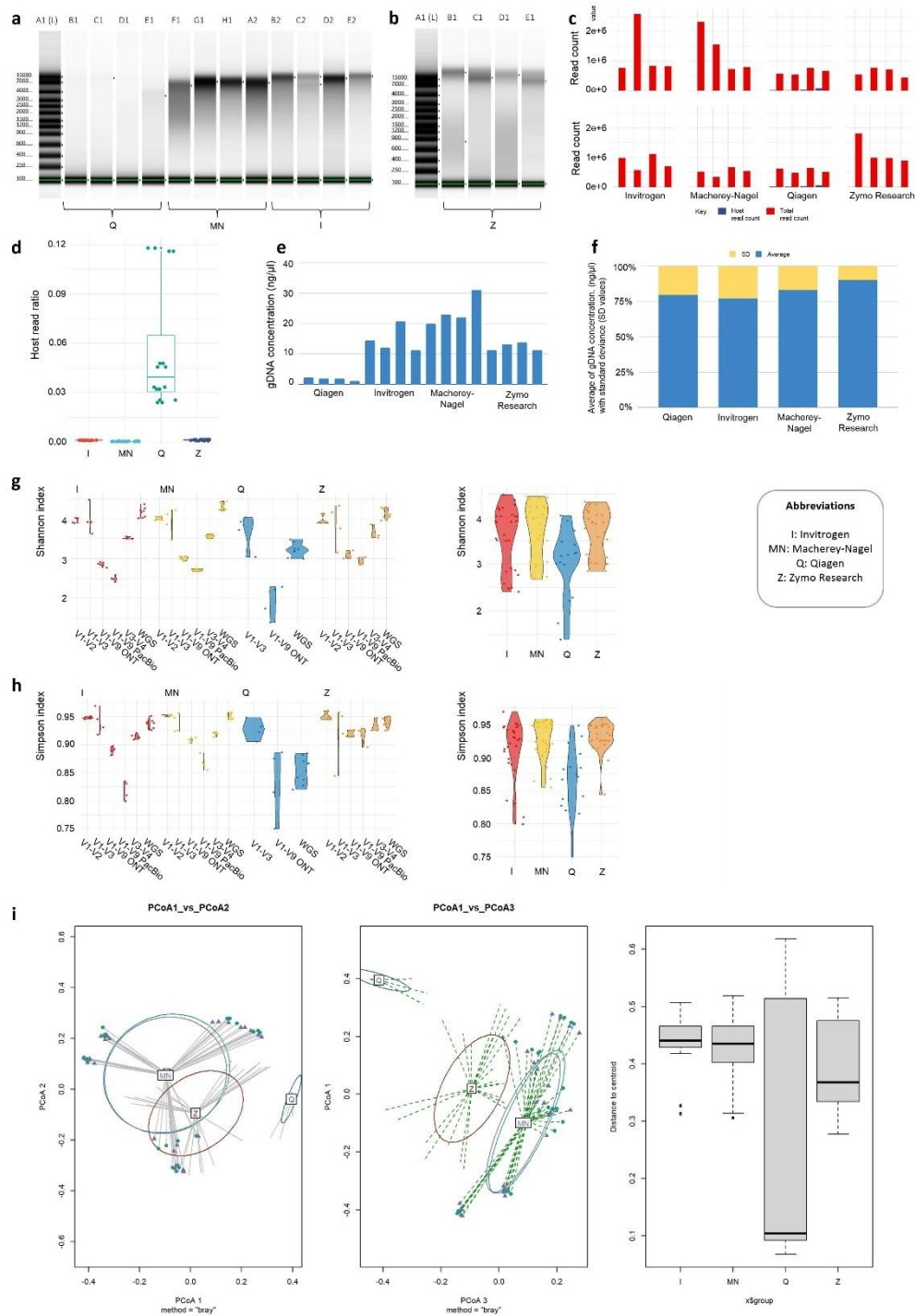


Figure 4: Impact of DNA isolation kit on yield, quality, and bacterial composition analysis.

Quality of the Isolated DNA; molecular weight of DNA **a**: lanes B1, E1, F1, A2 and B2, D2: four replicates of DNA samples were isolated using Q, MN and I kit, respectively. **b** lanes B1, E1: DNA samples extracted by Z kit. **c, d**: ratio of host DNA. **e, f**: Yield and reproducibility of the isolated DNA. **g-i**: Microbial composition, diversity and dispersion; Alpha-diversities in the different library preparation methods, according to each DNA isolation method: **g**: Shannon index; **h**: Simpson index; Beta-diversity analysis of each DNA isolation method: **i**: PCoA plots and distances to group centroids.

We found marked differences in DNA degradation across the kits, with the Q kit producing the most degraded samples (Fig. 4a, lanes B1, E1). To obtain DNA of sufficient integrity for full-length 16S rRNA sequencing with ONT, we applied the Quick-DNA High Molecular Weight (HMW) MagBead Kit from Zymo Research. Although longer DNA fragments indicate lower degradation, our findings suggest that HMW DNA is not strictly required for sequencing the complete V1-V9 region. DNA quality and yield from the MN (Fig. 4a, lanes F1, A2) and I (Fig. 4a, lanes B2, E2) kits were also appropriate for LRS, as the differences in average fragment length across these three kits had little impact on V1-V9 sequencing (Fig. 4a, lanes F1, E2 and Fig. 4b). However, for LRS-based WGS, maintaining HMW DNA was critical. Based on fragment length, the Z kit proved most suitable for this purpose (Supplementary Data 3). We further assessed each kit's ability to selectively extract bacterial DNA. Even when using the Q kit with its "Isolation of DNA from Stool for Pathogen Detection" protocol, host DNA contamination remained significantly higher than with the other kits (Fig. 4c, d).

For canine stool samples, all kits except the Q kit provided DNA of sufficient quality and quantity for sequencing (Fig. 4e, f and Supplementary Data 3a). Analysis of the Q kit on six different dog samples consistently produced low yields and poor quality (Supplementary Data 3b), whereas the other three kits showed reproducible performance across replicates. Among them, the Z kit displayed the greatest consistency, while the I kit had the highest variability (Fig. 4e, f). Substantial yield differences were observed across kits (Supplementary Data 3). For canine feces, the Q kit gave the lowest DNA yield, the I kit produced moderate amounts, and the MN kit delivered the highest yield. The Z kit achieved relatively high yields despite using only half the starting sample size (Table 1 and Fig. 4e, f). An ANOVA revealed strong differences among the kits with an F-value of 511.63 and a p-value below 0.0001. Pairwise comparisons indicated that the I kit yielded more DNA than Q, the MN kit outperformed both, and the Z kit also gave higher yields than Q, although MN still provided the best overall yield. Adjusting for input volume, the Z kit would likely match MN, suggesting the two are comparable when equal sample sizes are used.

The average of DNA concentration isolated from dog stool	Average gDNA cc. (ng/μl)
Qiagen	1.815 ± 0.464
Invitrogen	14.6 ± 4.35
Macherey-Nagel	24.0 ± 4.83
Zymo Research	12.4 ± 1.36 ng/μL
The average of DNA concentration isolated from the MCS sample	Average gDNA cc. (ng/μl)
Qiagen	0.364 ± 0.12
Invitrogen	0.287 ± 0.048
Macherey-Nagel	10.564 ± 0.91
Zymo Research	38.72 ± 3.33
The average of DNA concentration isolated from the GMS sample	Average gDNA cc. (ng/μl)
Qiagen	1.146 ± 0.157
Invitrogen	0.0683 ± 0.0076
Macherey-Nagel	12.25 ± 0.636
Zymo Research	4.087 ± 0.31
The average of the longest DNA fragments isolated from dog stool	Average length of DNA (bp)
Qiagen	9361
Invitrogen	26,341
Macherey-Nagel	21,325
Zymo Research	33,867
The average of the longest DNA fragments isolated from the MCS sample	Average length of DNA (bp)
Qiagen	>60,000
Invitrogen	>60,000
Macherey-Nagel	54,74
Zymo Research	>60000
The average of the longest DNA fragments isolated from the GMS sample	Average length of DNA (bp)
Qiagen	>60,000
Invitrogen	58,897
Macherey-Nagel	55,538
Zymo Research	>60,000

Table 1: Average concentrations and lengths of obtained genomic DNA samples

The first sections present the average concentrations of genomic DNA (gDNA) obtained from dog stool, MCS, and GMS samples, measured in ng/μL, respectively. Values are reported as means with standard deviations. The last blocks of tables report the average length of the longest DNA fragments (in base pairs, bp) obtained from dog stool, MCS, and GMS samples, respectively.

In canine fecal samples, the Z method generated the longest average fragment peaks, followed by I with moderately long fragments, MN with slightly shorter peaks, and Q with the shortest fragments (Supplementary Data 3 and Table 1). A Tukey HSD test confirmed significant differences in fragment length among several kit pairs (Supplementary Data 3c). Supplementary Fig. 3a illustrates these differences, while Supplementary Fig. 3b shows violin plots of fragment length distributions across kits with colors reflecting concentration.

To determine how much the DNA extraction kit influenced microbial composition, we analyzed data using the same database and bioinformatics workflow. Because commonly used pipelines are optimized for specific sequencing methods, we applied the minitax tool with an NCBI genome collection as the reference.

Earlier studies suggested that α -diversity indices can reflect DNA extraction efficiency²⁰. In our study, the Q kit consistently led to a strong reduction in both richness (Shannon index) and evenness (Simpson index) compared to the other kits (ANOVA, p -values < 0.000 for all comparisons) (Fig. 4g, h). Other pairwise richness differences were not significant, although the I kit showed reduced evenness relative to Z (Supplementary Data 4). We also identified read number thresholds for each sequencing approach at which α -diversity values remained stable without significant drop (Supplementary Fig. 4).

For β -diversity, we calculated a Bray-Curtis distance matrix at the sample level. From these distances, NMDS, PERMDISP, and PERMANOVA analyses were performed. These revealed clear compositional differences depending on the DNA extraction method. The full-model PERMANOVA indicated that DNA extraction accounted for 28.2% of the variance in microbial community composition across sequencing methods and platforms ($p < 0.001$).

We further examined the dispersion of samples for each extraction method, calculated as the average distance to the centroid in multivariate space (PERMDISP2). These values showed notable variability that can strongly influence conclusions about microbial composition. The results are visualized in PCoA plots for each extraction method (Fig. 4i). These plots revealed strong overlap between I and MN, partial overlap of Z with these two, and Q clustering apart from all others.

When analyzed separately for each library, the I kit consistently produced the lowest dispersion values with small standard deviations, showing a mean distance of 0.0626 and a standard deviation of 0.0586 (Supplementary Fig. 5a). The MN kit also gave low dispersion across nearly all libraries except PE V1-V3. The Z kit produced higher dispersions in LRS datasets but results

comparable to I and MN in SRS libraries. The Q kit showed variability comparable to the others in the two libraries suitable for deeper analysis (PE V1-V3 and I WGS). PERMANOVA tests for each library type confirmed that DNA extraction methods influenced microbial composition across both V-regions and WGS. The MN and I kits consistently yielded similar community profiles across protocols, while Z and Q diverged significantly from both MN and I, and often from each other. These trends are represented in the library-specific PCoA plots based on PERMDISP values (Supplementary Fig. 5b). Pairwise tests confirmed these differences across most V-regions, with the I/M pair being the only exception.

4.3. Assessment of library construction and sequencing methods

In this part of the study, our aim was to assess the consistency of variability among six different library preparation strategies. To do this, we compared the outcomes of each DNA isolation method (I, MN, Q, and Z) across the available sequencing libraries. Because the Q kit presented quality and quantity issues, only the WGS, PerkinElmer (PE) V1-V3, and ONT V1-V9 libraries were prepared from those samples. As most commonly applied programs are designed specifically for either amplicon sequencing or WGS, we used the same tool (minitax) across all sequencing methods, with an NCBI genome collection serving as the reference database.

4.3.1. Quality, yield, and reproducibility

Among the tested methods, the ONT 16S rRNA amplicon library required the least experimental input, while the Illumina WGS and PacBio 16S rRNA amplicon workflows required the most manual effort (Supplementary Fig. 6). With the exception of the PE method, which frequently produced two non-specific DNA fragments (400 and 800 bp), the other approaches showed excellent performance in quality, yield, and reproducibility for canine samples (Supplementary Fig. 7). However, read quality with Nanopore sequencing was somewhat lower (Supplementary Data 5).

4.3.2. Microbial composition, diversity and dispersion

Analysis of α -diversity revealed that in the short-read sequencing (SRS) libraries, richness and evenness values were similar across the three kits, except for Q. In the long-read sequencing (LRS) libraries, the I method produced somewhat lower diversity, the Z method the highest, and MN also showed comparably high values. For β -diversity, we calculated Bray-Curtis distances at the sample level. The NMDS plot in Fig. 5a shows that samples tended to cluster according to library preparation rather than DNA extraction method. Samples prepared with the I and MN kits displayed broader spread, suggesting that community composition varied with library preparation. The V3-V4 group formed a cluster clearly separated from V1-V3, which

was more closely related to V1-V2. The Z method produced somewhat tighter groupings, while the Q kit's WGS and V1-V3 libraries clustered closely together. Taken together, both extraction and sequencing approaches influenced community composition, with sequencing protocols being the stronger factor. A full-model PERMANOVA confirmed this, showing that library type was the dominant variable, explaining 58.8% of the total variance. We also performed PERMDISP for each group (Fig. 5b). PCoA visualization showed that none of the extraction methods produced identical microbial profiles across all libraries. Noticeable overlaps were only seen in V1-V2 and V1-V3 libraries prepared using I and MN (Fig. 5b-d). LRS and WGS libraries created from Z DNA yielded results that were broadly similar to those from I and MN. Variability within each library type differed, with the tightest clustering observed in the V3-V4 libraries of I samples, pointing to higher consistency.

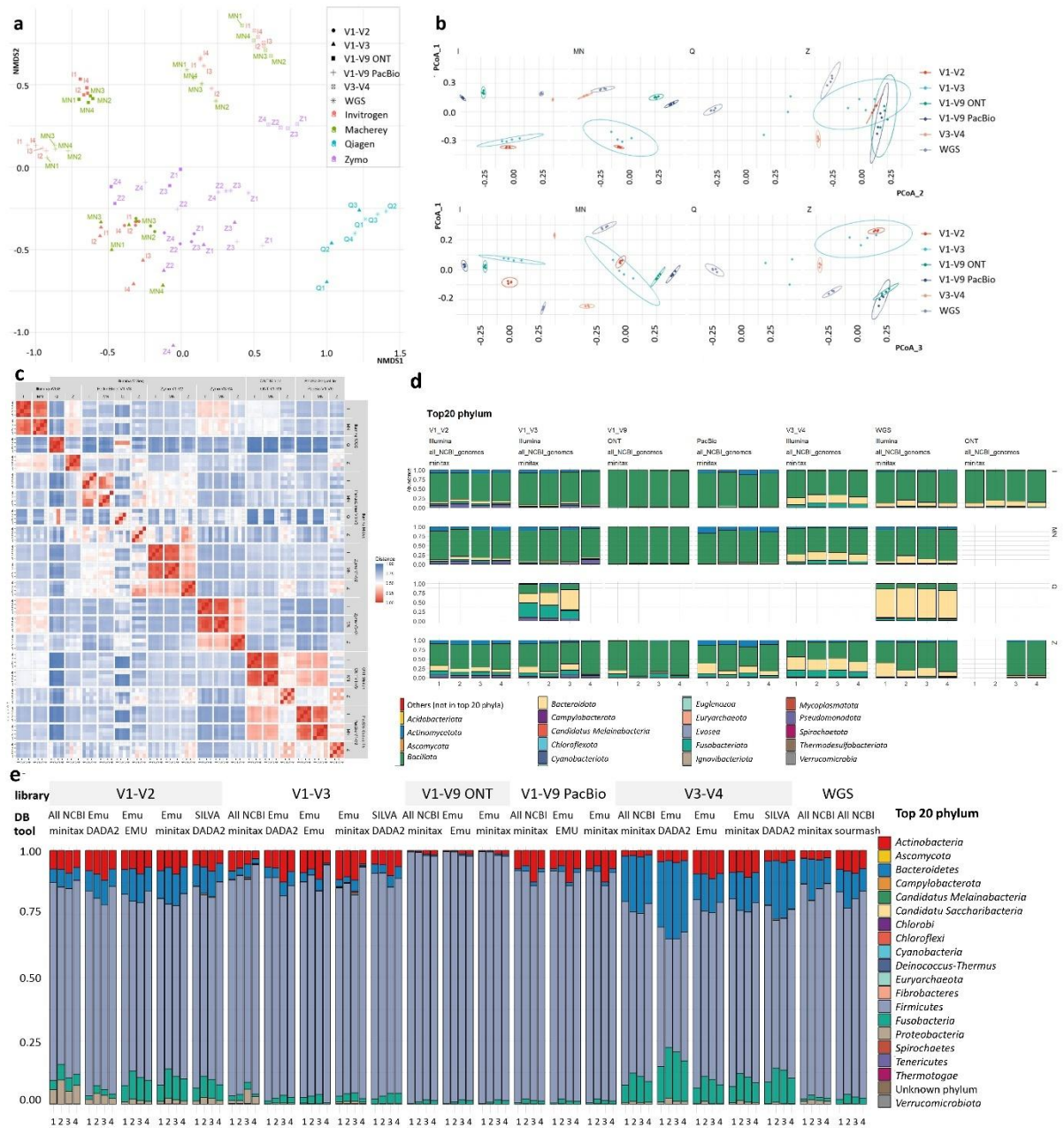


Figure 5: Comparison of library preparation protocols. a: β -diversity analysis – NMDS. The sample-wise Bray-Curtis distances between the samples were calculated and plotted using NMDS. Colors show the DNA isolation methods, while shapes indicate the library preparation protocols. **b:** β -diversity analysis—PERMDISP. PCoA visualization of the dispersion values from the PERMDISP results. The plot is faceted according to the DNA isolation methods and colors indicate the library preparation protocols. **c:** β -diversity analysis—Heatmap. The sample-wise Bray-Curtis distances between the samples were calculated and visualized on a heatmap. The WGS data composition for the I and MN isolation kits is most similar to the bacterial composition from the Z V3-V4 library, followed by the ONT V1-V9 library. Notably, for both I and MN DNA isolation kits, the WGS and ONT V1-V9 libraries exhibit fairly similar

compositional profiles. The Z DNA isolates produce relatively consistent microbial community profiles across different library preparation techniques. The PE V1-V3 library for Q DNA shows some similarity to WGS. The I and MN DNA samples exhibit significant overlap with the V1-V2 libraries prepared from I, MN, and Z DNA, with Z DNA showing a particularly high degree of overlap with itself. Additionally, Z DNA shows good overlap with the V1-V3 library compared to other libraries. The V3-V4 library consistently shows the least similarity to other libraries across all DNA isolation methods. The PacBio results are most similar to ONT, regardless of the DNA isolation kit used. When using Z DNA isolation method, PacBio shows significant similarity in microbial composition with other SRS amplicon-based sequencing results, particularly with V1-V2, V1-V3, and then V3-V4. WGS is significantly different. The Z-isolated DNA ONT V1-V9 libraries also show similar compositional similarities with these, with PacBio possibly showing a closer match. **d:** Relative abundance of the top 20 phyla in each sample with. Rows indicate the DNA isolation methods, while columns indicate the library preparation protocols. **e:** Barplot showing the top 20 phyla in the MCS samples, sequenced using Invitrogen DNA isolation kit, according to library preparation protocols, analyzed with minitax, DADA2, and Emu programs using several different databases.

The β -diversity heatmap in Fig. 5c illustrates pairwise comparisons across DNA extraction and library preparation methods. Regardless of the isolation kit, PacBio results showed similarity between V1–V2 and V3–V4, with the latter two being most alike. In these same regions, I and MN results were closer to each other than to Z. When libraries were compared internally, V1–V3 consistently performed the weakest.

Because overlaps in Fig. 5b were incomplete and variability across libraries was evident, we measured the distance of each sample to the centroid of its group. Considerable variability was observed in V1–V3 libraries regardless of extraction method. With the exception of V1–V3 and LRS libraries made from Z DNA, most library types showed relatively low variability (Supplementary Fig. 5c).

Pairwise PERMANOVA with library preparation as the only factor confirmed that all libraries differed significantly from each other, even after multiple-testing correction. The partial overlaps visible in the PCoA plots (Supplementary Fig. 5b) were not significant by PERMANOVA, which tests centroid separation. In V1–V3 libraries, however, differences were influenced not only by centroid shifts but also by differences in dispersions.

Community composition analysis for I samples revealed that the V1–V3 libraries contained the highest proportion of *Bacillota* (formerly *Firmicutes*), while *Bacteroidota* (formerly *Bacteroidetes*) and *Fusobacteriota* (formerly *Fusobacteria*) were detected in much smaller amounts. In contrast, the V1–V2 and V3–V4 libraries showed a more balanced representation of *Bacteroidota* and *Fusobacteriota* (Fig. 5d). We compared these results with published datasets^{18,100–104} (Supplementary Data 6 and Fig. 6), all of which focused on specific aspects of the canine gut microbiome. A common feature across these studies, as well as our own, was the dominance of five phyla. Nonetheless, their relative abundances varied markedly across methods (Fig. 6), highlighting the strong influence of methodological choices on outcomes.

a

Invitrogen						
	V1-V2	V1-V3	V3-V4	WGS	V1-V9 (ONT)	V1-V9 (PacBio)
<i>Firmicutes</i>	60,147%	67,040%	63,333%	78,122%	69,710%	48,198%
<i>Bacteroidetes</i>	9,474%	2,741%	8,237%	1,156%	0,234%	0,099%
<i>Fusobacteria</i>	9,868%	2,080%	6,555%	2,082%	0,625%	1,100%
<i>Proteobacteria</i>	0,595%	0,849%	2,151%	1,029%	0,120%	0,156%
<i>Actinobacteria</i>	9,621%	4,539%	4,753%	17,149%	0,415%	0,908%

Macherey-Nagel						
	V1-V2	V1-V3	V3-V4	WGS	V1-V9 (ONT)	V1-V9 (PacBio)
<i>Firmicutes</i>	59,972%	66,778%	62,846%	74,882%	40,797%	29,404%
<i>Bacteroidetes</i>	8,561%	2,955%	7,517%	0,834%	0,287%	0,098%
<i>Fusobacteria</i>	9,428%	1,453%	6,409%	1,313%	0,390%	0,494%
<i>Proteobacteria</i>	0,528%	0,774%	2,154%	0,703%	0,083%	0,082%
<i>Actinobacteria</i>	11,527%	3,880%	5,752%	21,857%	0,217%	0,654%

Zymo Research						
	V1-V2	V1-V3	V3-V4	WGS	V1-V9 (ONT)	V1-V9 (PacBio)
<i>Firmicutes</i>	43,191%	50,512%	44,391%	78,736%	75,603%	63,706%
<i>Bacteroidetes</i>	15,003%	11,389%	14,198%	1,762%	2,402%	0,533%
<i>Fusobacteria</i>	17,569%	8,854%	12,108%	1,683%	6,342%	10,450%
<i>Proteobacteria</i>	0,971%	1,676%	4,658%	1,461%	1,320%	1,701%
<i>Actinobacteria</i>	11,834%	3,413%	5,576%	15,979%	1,178%	2,350%

Qiagen						
	V1-V2	V1-V3	V3-V4	WGS	V1-V9 (ONT)	V1-V9 (PacBio)
<i>Firmicutes</i>		17,261%		51,174%		
<i>Bacteroidetes</i>		26,861%		18,961%		
<i>Fusobacteria</i>		27,488%		20,869%		
<i>Proteobacteria</i>		5,033%		6,620%		
<i>Actinobacteria</i>		1,584%		1,342%		

b

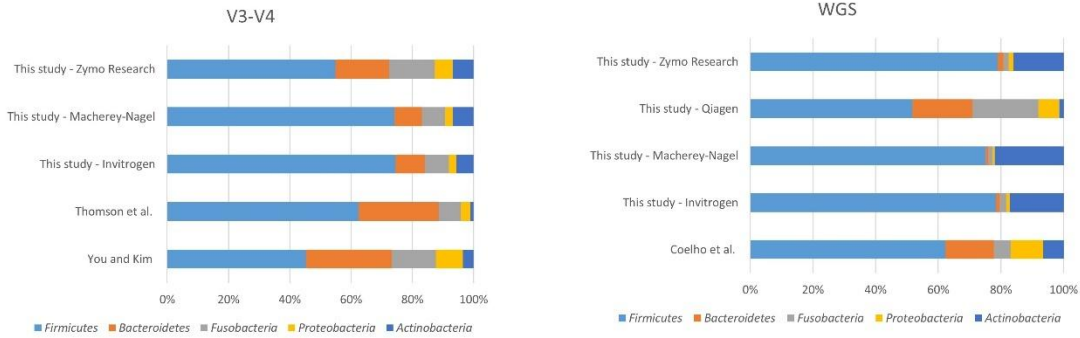


Figure 6: Comparison of the bacterial composition of various samples at the phylum level.

This figure presents a comparison of bacterial composition at the Phylum level in various samples. **a:** Tables show the proportions of the most abundant Phyla in samples obtained using different DNA extraction and library preparation methods. **b:** Bar charts highlight the discrepancies between the proportions of the most abundant bacterial Phyla from other groups' data and our datasets. Taxa not included in the defined composition are grouped together under the label "other".

4.3.3. Difference in ratio of Gram-positive and Gram-negative bacteria

To investigate whether certain DNA extraction methods preferentially capture Gram-negative (Gram⁻) bacteria compared to Gram-positive (Gram⁺) ones, we grouped species abundances according to their cell-wall staining properties and performed a similar analysis. We observed that the differences between extraction methods could largely be explained by the varying resistance of bacterial cell walls to lysis. Only the I and MN kits produced microbial community profiles with nearly identical Gram⁺ to Gram⁻ ratios. By contrast, all other pairwise comparisons revealed clear discrepancies in these ratios. Regarding library-level comparisons, we found several significant differences, with the V3-V4 region being particularly prominent in this respect (Supplementary Fig. 8).

4.4. Comparison of DNA isolation methods and sequencing libraries on synthetic microbial community standards

To assess how accurate our DNA extraction and library preparation approaches were, we used two synthetic microbial standards ZymoBIOMICS Microbial Community Standard D6300 (MCS) and ZymoBIOMICS Gut Microbiome Standard D6331 (GMS), both from Zymo Research] that have defined compositions (Fig. 7, Supplementary Fig. 9a, and Supplementary Data 7). From the MCS community, we prepared Illumina V1-2 and ONT V1-9 libraries, and from the GMS community, ONT V1-9 libraries, using DNA extracted by each of the four kits. The MCS contains eight bacterial species, of which five are Firmicutes and three are Proteobacteria. The GMS includes 18 bacterial species and one archaeal species, representing a more diverse mixture than MCS and offering a closer model of the human gut microbiome. Although Matsuo et al. reported that the V1-V9 primers are biased against *Bifidobacterium*¹⁰⁵, we chose to calculate our statistics without excluding this taxon, as neither the standard's manufacturer (Zymo) nor the primer manufacturer (ONT) recommends its removal. Unlike the fecal DNA analyses, where we focused on the consistency of methods within groups and the degree of similarity between them, our goal here was to determine how well each method reproduced the expected composition of the synthetic communities. We carried out this comparison using both the minitax and Emu software with the NCBI genome collection and the Emu database as references.

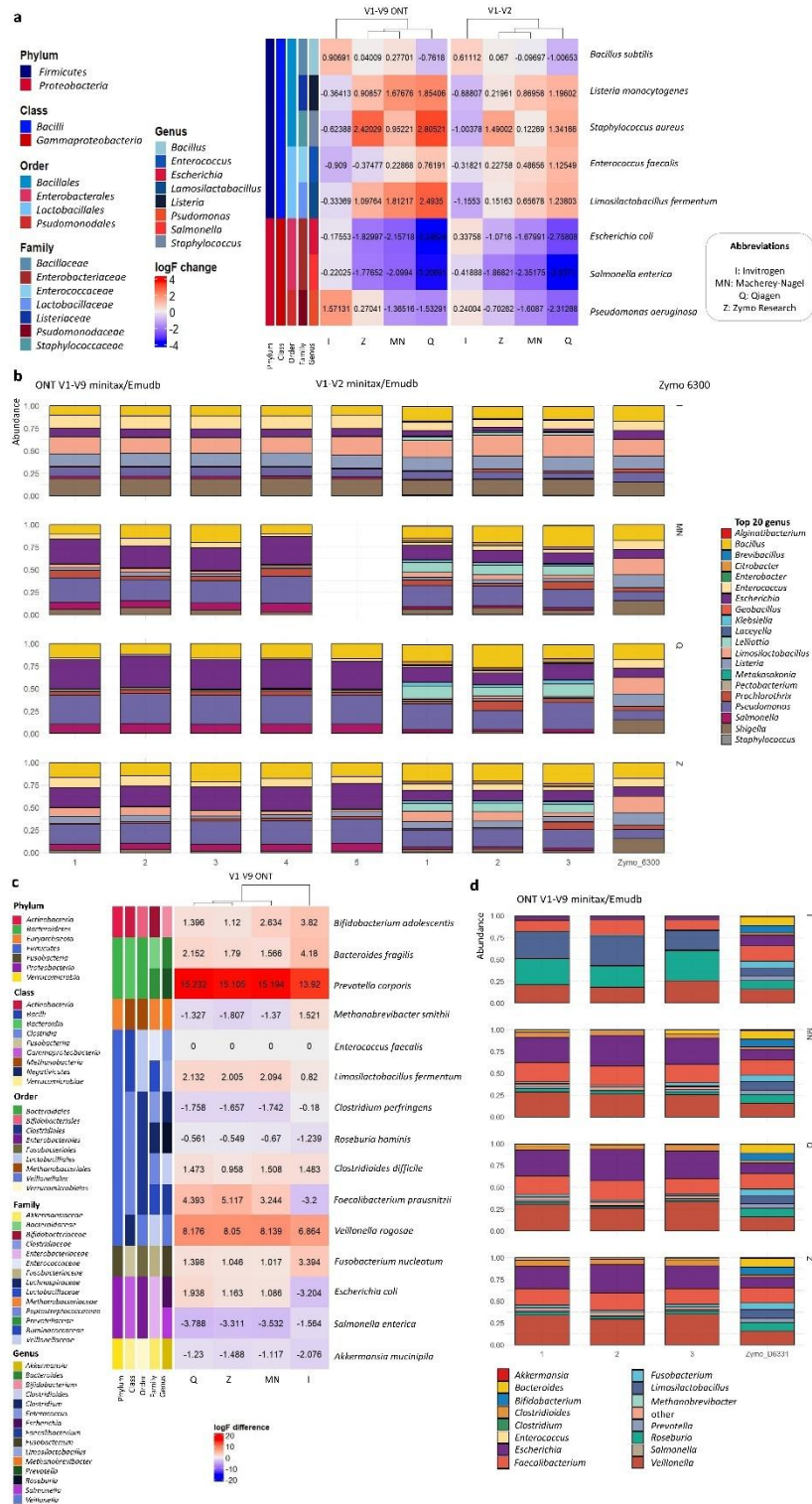


Figure 7: Comparative analysis of DNA isolation and library preparation protocols in MCS and GMS. a: Heatmap of sample-wise differences in the MCS samples. The abundance values identified by the Emu application were compared to the theoretical values provided by Zymo and log2 fold changes were estimated and are shown within the boxes. Deeper blue colors indicate lower experimental values compared to the theoretical, while more red colors indicate

higher experimental values. **b:** Barplots showing the top 20 phyla in MCS samples using the Illumina V1–V2 and ONT V1–V9 methods, according to the DNA isolation methods, analyzed with minitax using the Emu genome collection database. **c:** Heatmap of sample-wise differences in the GMS. The abundance values identified by the Emu tool were compared to the theoretical values provided by Zymo Research, with log2 fold changes calculated and displayed within the boxes. Darker blue colors represent lower experimental values compared to the theoretical, while dark red colors indicate higher experimental values. **d:** Barplots show the top 20 phyla in GMS samples using the ONT V1-V9 methods, according to the DNA isolation methods, analyzed with minitax using the Emu database.

4.4.1. Quality and yield of DNA

In our experiments, DNA yield differed significantly across methods for both MCS and GMS samples (Table 1). For MCS, the Z kit achieved the highest average yield, exceeding that of all others, including MN, which itself performed better than both Q and I. ANOVA confirmed the differences, with an F-value of 124.75 and a p-value below 0.0001. Post-hoc analysis showed that Z generated significantly more DNA than Q and I, and also outperformed MN. For GMS samples, the MN kit provided the highest yield, while Z produced a slightly lower but still substantial amount. Among the remaining kits, Q yielded more DNA than I, which produced the least. ANOVA for GMS also showed strong differences, with an F-value of 147.92 and a p-value below 0.0001. The HSD analysis indicated that MN generated significantly more DNA than I and Q, while Z produced nearly the same amount as MN. Overall, Z and MN performed best in both sample types, clearly surpassing Q and I (Table 1).

Our analysis of fragment lengths revealed further differences in degradation patterns across methods. In MCS samples, Z, I, and Q produced the longest average fragment sizes, all above 60,000 bp, whereas MN produced a shorter average peak length of 54,740 bp (Table 1). In GMS samples, Z and Q again achieved the longest fragments, both above 60,000 bp. I followed with an average peak length of 58,897 bp, and MN gave 55,538 bp (Table 1). These results indicate that both the extraction protocol and the sample type strongly affect DNA quality (Supplementary Data 7, Supplementary Fig. 9b).

The Tukey HSD test confirmed significant differences in fragment length between several kit pairs, specifically MN/I, MN/Q, and MN/Z. These findings underline the influence of extraction technique on DNA integrity in MCS samples, as presented in Supplementary Fig. 9b.

4.4.2. Microbial composition, diversity, and dispersion

4.4.2.1. ZymoBIOMICS Microbial Community Standard (MCS)

The dispersion analysis showed that in the SRS library, consistent with our fecal DNA observations, the I kit produced the lowest variability. Interestingly, in the LRS library, the Q kit also demonstrated similarly low dispersion. When examining microbial community profiles, we found that unlike the strong similarity seen *in vivo*, the MN and I kits diverged considerably in the MCS. The Q method remained the most distinct from I, while Z and MN showed closer resemblance to one another and a clear separation from I. This pattern was consistent across both library types and across both software tools and databases applied (minitax with Emu db, Emu with Emu db, and minitax with the NCBI genome collection, Fig. 7a, Fig. 7b,

Supplementary Fig. 10a, b and Supplementary Fig 11a, b), highlighting that *in vitro* standards behave differently than *in vivo* samples. The I kit produced results with the highest agreement to the expected values: no bacterial species showed significant differences using either bioinformatics tool or database, except *Pseudomonas aeruginosa* when analyzed with Emu and the Emu db. The other extraction methods showed between 1 and 5 significant differences depending on the wet lab approach and the bioinformatics workflow.

4.4.2.2. ZymoBIOMICS Gut Microbiome Standard (GMS)

For the more complex microbial standard, a different pattern was observed. The Z kit produced the fewest deviations from the expected composition (7–8 species), while the I kit produced the highest number of deviations (11 species). The MN and Q methods fell between these extremes, each identifying 8–10 species with significantly different abundances compared to the defined community. This trend was consistent regardless of the bioinformatics pipeline or database used. Z and MN results were more similar to each other than to I, while Q differed the most from the other three. These trends were visible using minitax with Emu db, Emu with Emu db, and minitax with the NCBI genome collection (Fig. 7c, d and Supplementary Figs. 10c, d and 11c, d).

The I kit, although successful with MCS, performed differently with the more complex GMS, indicating that robustness and sample characteristics strongly affect accuracy. Interestingly, I and MN behaved similarly when analyzing the *in vivo* samples (dog feces), showing that context influences performance. DNA concentrations were especially low for the GMS I kit (Supplementary Data 7b) and for MCS using Q and I, whereas *in vivo* only the Q kit gave insufficient amounts. These results emphasize that DNA yield must be taken into account, since low concentrations can impair downstream steps. Kits that consistently produce higher yields may be more suitable for analyzing complex microbial communities. Overall, none of the kits matched theoretical expectations equally well across different sample types. Portik et al. also concluded that MCS (D6300) is more reliable than GMS (D6331)⁸⁹, suggesting that synthetic standards are not always the best indicators for selecting DNA isolation kits or sequencing strategies for *in vivo* applications. Careful evaluation is recommended before starting large-scale studies, especially when working with less familiar sample sources.

4.4.3. Difference in ratio of Gram(+) and Gram(-) bacteria

In this part of the analysis, we found that the Q kit tended to overrepresent Gram– bacteria in MCS samples and also produced higher ratios of Gram– species in *in vivo* datasets. The MN kit showed a similar overrepresentation in fecal samples, but this bias did not appear in MCS.

For both kits, the main underrepresented Gram⁺ taxa were *Lactobacillus*/*Limosilactobacillus*, while the main overrepresented Gram[−] taxa were *Escherichia* and *Salmonella*. Our results suggest that the MN method favors Gram[−] bacteria in a controlled environment such as MCS, but this trend may be masked in fecal samples due to their complexity, which affects the lysis of both Gram⁺ and Gram[−] species. In GMS, the I kit underestimated Gram[−] bacteria, whereas MN, Q, and Z provided ratios close to the expected values.

4.5. Laboratory work: key findings

Our assessment of the DNA extraction kits and library preparation protocols revealed several important points (Table 2 and Table 3). The Z kit provided very high DNA yields and the longest fragments, making it the most effective overall for high-quality diversity analysis. The MN kit produced high yields and results close to theoretical compositions but gave the shortest fragments. The I method generated consistent data of good quality but showed more variability in yield and fragment length compared with Z and MN. The Q kit produced lower yields and higher host contamination, which reduced its reliability for microbial diversity analysis (Supplementary Fig. 3b). Regarding library preparation, Illumina MiSeq with Illumina DNA Prep (WGS) produced the most accurate and high-quality data, although it required significant labor and was expensive. Illumina MiSeq with PerkinElmer V1-V3 was less consistent and often generated incorrect fragments. ONT MinION with ONT V1-V9 performed well with minimal hands-on work but produced reads of lower quality compared to Illumina and PacBio. PacBio Sequel IIe with PacBio V1-V9 generated high-quality reads but was both expensive and labor-intensive.

Kit	Advantages	Disadvantages
Qiagen	<ul style="list-style-type: none"> • Suitable for certain applications 	<ul style="list-style-type: none"> • Lowest DNA yield among the kits • Higher levels of host DNA contamination • Less accurate in representing microbial diversity
Invitrogen	<ul style="list-style-type: none"> • Good overall quality • High amount of DNA from stool samples • Consistent results in various libraries and sequencing methods • Closest match to theoretical composition in MCS samples 	<ul style="list-style-type: none"> • Higher variability in yield and fragment length compared to Z and MN kits • Less optimal performance with GMS samples • Less effective in certain sample types (e.g., MCS and GMS) compared to Zymo and MN
Macherey-Nagel	<ul style="list-style-type: none"> • High overall DNA yield • Better performance in GMS samples • Good alignment with theoretical microbial composition in MCS samples • Comparable yield to Z kit when adjusted for initial volume 	<ul style="list-style-type: none"> • Shorter average DNA fragment length
Zymo Research	<ul style="list-style-type: none"> • High DNA yield in canine fecal and MCS samples, highest in GMS sample • Provides high-quality DNA • Longest average DNA fragment length • Superior performance in MCS samples • Closest match to theoretical composition in GMS samples • Consistent results and performance for LRS and WGS libraries • Stable performance with low relative dispersion 	<ul style="list-style-type: none"> • Less accurate representation in MCS

Table 2: Summary of DNA isolation methods: yield, fragment quality, and suitability for various sample types

This table details the advantages and disadvantages of various DNA isolation kits, including Zymo Research, Qiagen, Macherey-Nagel, and Invitrogen. Key factors include DNA yield, fragment length, contamination levels, and overall performance across different sample types.

Sequencing platform / library prep kit	Advantages	Disadvantages
Illumina MiSeq/Illumina DNA Prep (WGS)	<ul style="list-style-type: none"> • Known for high-quality data and accuracy in determining true microbial compositions • Shows high consistency across libraries • High-quality libraries 	<ul style="list-style-type: none"> • Labor-intensive and complex library preparation • Highest hands-on time • Generally more costly
Illumina MiSeq/PerkinElmer V1-V3	<ul style="list-style-type: none"> • Widely used V-region 	<ul style="list-style-type: none"> • The applied kit often produces nonspecific or erroneous fragments • Shows the least similarity to other libraries • Performance can be inconsistent across different samples
ONT MinION/ONT V1-V9	<ul style="list-style-type: none"> • Requires less experimental work compared to other methods • High-quality, high-consistency libraries • Shows high similarity with PacBio V1-V9 results and moderate similarity to Illumina WGS data 	<ul style="list-style-type: none"> • Generally lower read quality compared to Illumina and PacBio • May not consistently match well with other libraries, particularly short-read amplicon libraries
PacBio Sequel Hi/PacBio V1-V9	<ul style="list-style-type: none"> • Provides high-quality reads • High-quality, high-consistency libraries • Shows strong similarity with ONT results • Shows good moderate similarity with V1-V2 and V1-V3 libraries 	<ul style="list-style-type: none"> • Labor-intensive preparation • Typically more costly
Illumina MiSeq/Zymo Research V1-V2	<ul style="list-style-type: none"> • Similarity with PE V1-V3 data • High-quality libraries • High consistency, reliable data quality 	<ul style="list-style-type: none"> • Moderate costs and labor intensity
Illumina MiSeq/Zymo Research V3-V4	<ul style="list-style-type: none"> • Strong similarity with Illumina WGS (with I and MN DNA samples) • High-quality libraries • High consistency, reliable data quality 	<ul style="list-style-type: none"> • Moderate costs and labor intensity

Table 3: Comparison of library preparation methods: performance metrics and practical considerations

This table provides an overview of library preparation methods, including Illumina DNA Prep, PerkinElmer V1-V3, Zymo Research V1-V2 and V3-V4, ONT V1-V9, and PacBio V1-V9. It highlights the benefits and drawbacks related to sequencing quality, consistency, labor intensity, and cost for each preparation technique.

4.6. Comparison of bioinformatics techniques

To reduce the variability in canine microbiome composition caused by different DNA extraction methods, we used only the I kit for comparing bioinformatics approaches. To minimize differences linked to database choice, we selected the NCBI genome collection as the reference (Fig. 5e).

In this part of the study, we examined the canine microbiome at the phylum, order, genus, and species levels, comparing the influence of different databases (Emu, NCBI genomes, and SILVA) and bioinformatic tools (Emu, minitax, DADA2, and sourmash). A full-model PERMANOVA showed that databases and programs strongly affected the results, explaining 59.4% of the observed variation in microbial community composition. The “library” factor explained 20.1%, and the remaining 20.4% was unexplained. Pairwise comparisons revealed that DADA2 produced significantly different outcomes with the Emu database compared to the other tools, while minitax closely matched Emu. Using DADA2 with SILVA produced even larger differences (Fig. 5e). Although several combinations of databases and programs yielded significant differences ($p < 0.05$), the minitax and NCBI genome pairing was highly consistent. Significant differences were only observed when minitax with NCBI was compared to minitax with Emu for the ONT V1–V9 library, and to sourmash with NCBI for WGS. Aside from these, results remained stable across libraries. Therefore, the minitax and NCBI genomes combination was the most reliable choice for evaluating DNA extraction kits and library protocols (Fig. 5e).

We also analyzed MinION datasets prepared with the ONT 16S Barcoding Kit using the 16S module of the EPI2ME Labs software (EPI2ME Desktop, version 2021.09.09). This program provides a simplified workflow and rapidly generates Sankey diagrams on its online platform. Our initial analysis indicated that *Blautia* was the dominant genus (Fig. 8a, b), with an average abundance of around 80%. This was consistent across all DNA extraction kits but conflicted with both published data and sequencing from other methods (Illumina WGS, V1–V2, V3–V4, V1–V3, PacBio V1–V9). When analyzed with Emu and minitax, the *Blautia* proportion aligned with the ~15% reported in earlier studies¹⁰¹, showing that EPI2ME had strongly inflated this genus. However, closer inspection showed that EPI2ME did not directly overestimate *Blautia* but instead excluded some genera and underrepresented others. To clarify this, we compared the ONT V1–V9 datasets processed with EPI2ME and Emu. In our canine samples, EPI2ME failed to detect *Peptacetobacter*, *Faecalimonas*, and *Mediterraneibacter*, discarding about 75% of the reads and altering the community ratios. To confirm whether these discrepancies were entirely due to bioinformatics filtering, we expanded our analysis beyond the original dog

sample (Fig. 8a, b). We included fecal samples from six additional dogs (Fig. 8c), neonatal⁹⁵ (Fig. 8d) and adult¹⁰⁵ (Fig. 8e) human stool, adult human skin⁹⁴ (Fig. 8f), and the MCS standard (Fig. 8g). All ONT V1-V9 datasets were analyzed with both EPI2ME and Emu.

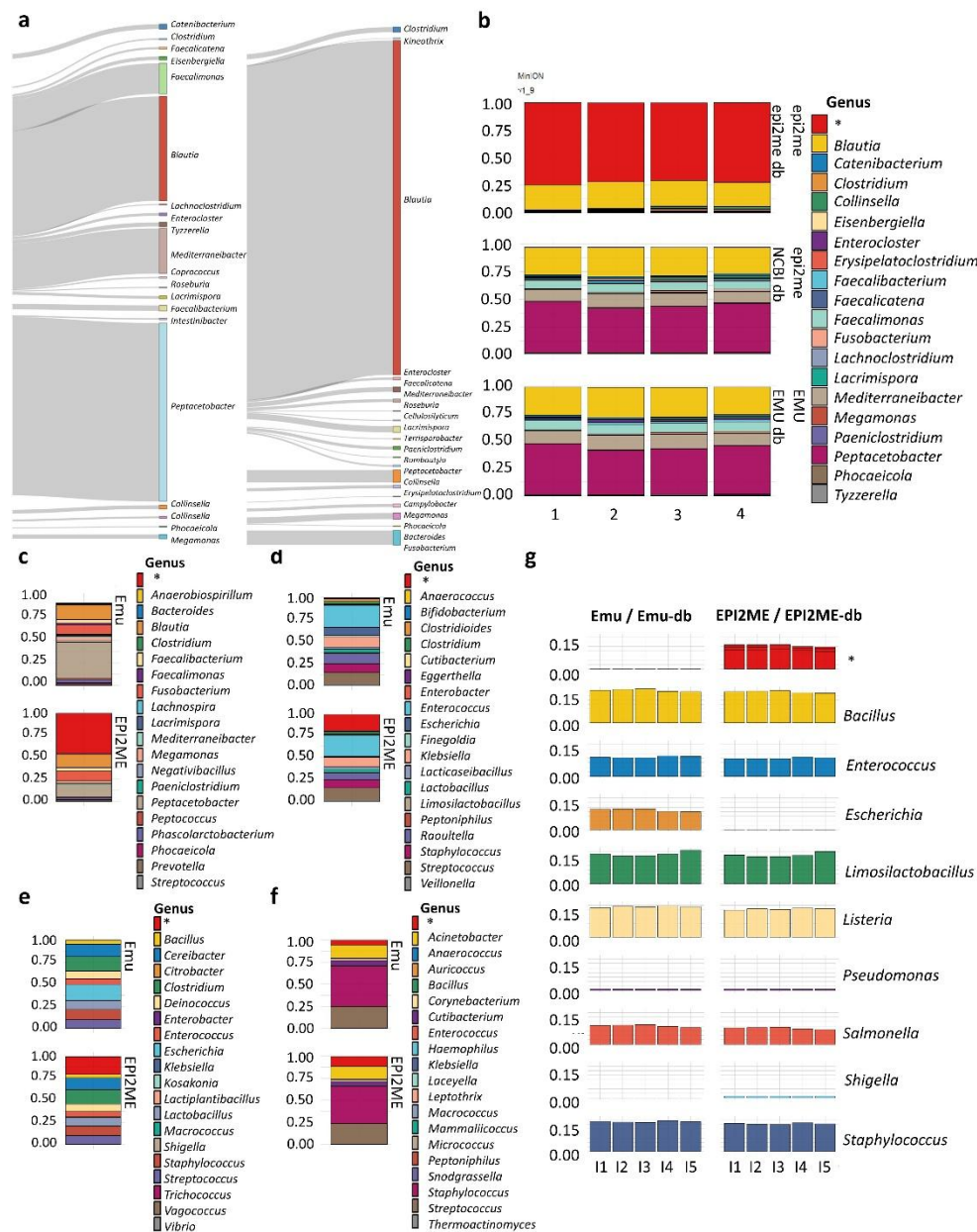


Figure 8: The EPI2ME analysis platform yields inaccurate genus-level data.

Results from the 16S module of the EPI2ME Labs program by ONT employed on canine stool sample (a,b), fecal samples from six additional dogs (c), neonatal⁹⁵ (d) and adult¹⁰⁵ (e) human stool specimens, adult human skin specimen⁹⁴ (f), and MCS (g). In our canine samples, the software failed to yield abundances to the *Peptacetobacter*, *Faecalimonas*, and *Mediterraneibacter* genera, discarding 75% of the reads and skewing the actual compositional ratios. * In b-g parts of this figure the blocks highlighted in red represent the reads that the EPI2ME program filters out and excludes from the analysis due to the LCA tag.

For the six control dog fecal samples, EPI2ME discarded around 50% of the reads (Fig. 8c). Similar to the original dog sample, *Peptacetobacter* was not detected and *Faecalimonas* and *Mediterraneibacter* were underestimated. Using ONT V1-V9 data from the European Nucleotide Archive (ENA), we compared EPI2ME and Emu again. For fecal samples, EPI2ME failed to process 15–25% of the reads and completely excluded *Escherichia* (Fig. 8d, e). For skin samples, the differences were minimal, with about 5% of the reads excluded by EPI2ME, most of which belonged to *Cutibacterium* and *Staphylococcus* (Fig. 8f).

In the microbial community mix from Zymo Research, EPI2ME excluded about 15% of the reads and failed to detect *Escherichia*. At the same time, it identified *Shigella*, a genus absent from the mixture (Fig. 8g).

We found that EPI2ME discarded large portions of data due to the use of Lowest Common Ancestor (LCA) tags. In the BLAST module, the LCA tag removes reads if the top three predicted genera do not match. When we reassigned these reads using their NCBI taxon IDs, the results produced by EPI2ME became very similar to those from Emu. This showed that the LCA filtering step is unnecessary and lowers accuracy, making EPI2ME less reliable when applied in its default configuration (Fig. 8b).

4.7. Evaluation of minitax: benchmarking across various sequencing methods and data types

To allow consistent comparisons across different sequencing datasets, we created minitax, a tool tailored for metagenome sequencing. We tested it thoroughly on multiple platforms and datasets, and compared its performance with other bioinformatic tools based on how accurately they could reproduce the reference microbial composition, measured by the correlation (r^2 values) between observed and theoretical communities. In addition, we applied Chi-square tests to determine whether the reconstructed and expected distributions showed statistically significant differences.

4.7.1. Comparing minitax with Emu using ONT V1-V9 sequencing of MCS and GMS

We compared minitax to Emu, which also uses minimap2 for alignment but applies an expectation-maximization step afterwards. For the MCS dataset, minitax and Emu produced highly similar results with the Emu database, confirming the robustness of minitax (Fig. 9a and Supplementary Fig. 12a). Even though full genome databases are usually not applied in 16S rRNA-Seq, minitax with the NCBI genome collection provided accurate reconstructions up to the genus level. This flexibility is valuable for researchers wishing to use the same database

across both WGS and 16S rRNA analyses. However, species-level accuracy dropped notably when the NCBI database was applied, emphasizing the importance of choosing a database that matches the resolution required. This pattern was also visible with GMS (Fig. 9a and Supplementary Fig. S12b). Since both programs performed comparably across the two synthetic standards, the differences seen between MCS and GMS reconstructions are more likely due to extraction methods or experimental variables rather than the bioinformatics tools or databases. The Chi-square results showed that only for the I kit with MCS did the reconstructed composition, whether with Emu or minitax using the Emu database, not differ significantly from the expected values (Supplementary Fig. S12a, b).

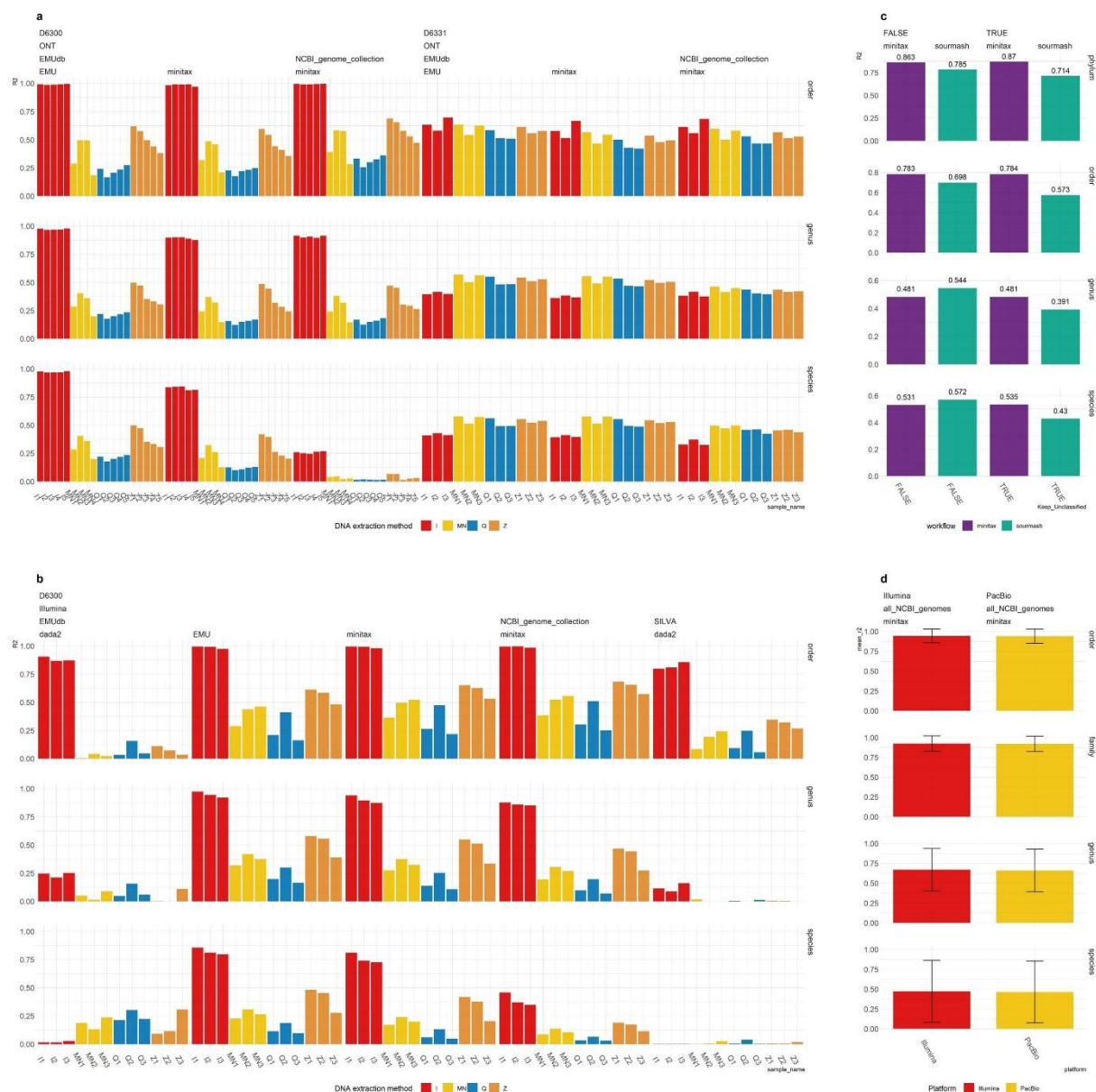


Figure 9: Comparative analysis of minitax performance across different sequencing platforms and datasets. The comparison between the theoretical and the observed microbial compositions, was generated using different software and databases. r^2 values were calculated across four taxonomic levels (phylum, order, genus, species) and plotted. **a** ONT V1-V9 Sequencing: performance comparison of minitax against Emu for the Zymo D6300 microbial community reference. **b** Illumina V1-V2 sequencing: minitax benchmarking against Emu and DADA2 for Zymo D6300 MCS. **c** PacBio HiFi WGS: minitax vs. sourmash performance comparison on Zymo D6331 MCS dataset. **d** CAMISIM mouse gut dataset: comparative analysis of minitax performance on Illumina and PacBio platforms. Ten samples were utilized from the CAMISIM database. Y-axis represent the mean r^2 -value for each group, error bars show standard deviation ($n = 10$).

4.7.2. Comparing minitax with Emu and DADA2 using Illumina V1-V2 sequencing of MCS

We next applied Illumina-sequenced V1-V2 data from MCS DNA. Here, both minitax and Emu produced substantially higher r^2 values than DADA2 at both genus and species levels, even when DADA2 was paired with the standard SILVA database (Fig. 9b). It is worth noting that although both Emu and DADA2 are tailored for amplicon sequencing, minitax remained competitive, particularly because of its ability to handle both 16S rRNA and WGS datasets. Emu achieved slightly higher r^2 values than minitax at the species level when used with the Emu database. In the MCS samples sequenced from the V1-V2 region, Chi-square testing confirmed that both Emu and minitax produced microbial compositions not significantly different from the expected community (Supplementary Fig. 12b), similar to the results with ONT V1-V9.

4.7.3. Comparing minitax with sourmash using MCS data of PacBio HiFi WGS

We also compared minitax with sourmash, one of the most widely used tools for analyzing both LRS and SRS WGS data¹⁰⁶. This comparison was performed on a PacBio HiFi dataset (NCBI accession: SRX9569057) from GMS. We calculated r^2 values both including and excluding unclassified reads. Minitax surpassed sourmash when unclassified reads were included in abundance estimates but performed slightly worse at the species level when these reads were excluded (Fig. 9c), since this adjustment shifted the relative abundances of identified taxa. Chi-square tests with three different detection thresholds showed that with minitax, the reconstructed and expected compositions were not significantly different when 0.1% and 0.01% thresholds were applied (Supplementary Fig. 12c).

4.7.4. CAMISIM: simulated mouse gut datasets

To broaden the evaluation, we also tested CAMISIM-simulated mouse gut datasets, which included 10 samples each from PacBio and Illumina. In both cases, minitax showed stable performance, achieving $r^2 = 0.96$ at the phylum level. This value decreased to a mean of $r^2 = 0.46$ with Illumina and $r^2 = 0.55$ with PacBio at the species level (Fig. 9d).

5. Discussion

In this comprehensive study, we compared metagenomic strategies using canine fecal samples alongside two synthetic microbial community standards. Our main goal was to assess the efficiency of different methodological steps, including DNA extraction, library preparation, and bioinformatic processing. To support this, we developed “minitax,” a flexible bioinformatic tool designed to work with a wide range of metagenomic laboratory protocols.

Selecting the right DNA extraction method is critical for maximizing yield, minimizing fragmentation, and maintaining the integrity needed for downstream analyses.²⁰ An ideal extraction kit produces high-quality, high-yield DNA to lower the chance of false negatives⁹⁶, removes PCR inhibitors that are common in fecal material^{107–111}, ensures consistent results, and efficiently lyses Gram+ bacterial cell walls to give a faithful picture of the microbial community.^{21,98,99,112,113}

Our evaluation of four DNA isolation kits—Zymo (Z), Qiagen (Q), Macherey-Nagel (MN), and Invitrogen (I)—showed significant differences in yield, fragment size, and overall performance. In canine fecal samples, the Z kit provided the highest DNA yield and the longest fragments with very low variation, highlighting its stability (Table 1, Fig. 4a–f). The Q kit yielded the least DNA and had a higher proportion of host contamination. The MN kit produced the largest overall yield but shorter fragments, while the I kit generated acceptable quality but with more variability in both yield and fragment length. For MCS, the Z kit produced considerably more DNA with longer fragments, outperforming MN, which had shorter peaks (Table 1). For GMS, MN slightly outperformed Z in yield, although both performed strongly (Table 1). Altogether, Z and MN showed the best yields, with Z excelling in MCS and MN performing better in GMS.

When examining microbial diversity in canine fecal samples, I and MN consistently provided similar outcomes across different libraries, sequencing methods, and bioinformatic pipelines. Z differed somewhat from these, while Q consistently performed the weakest in capturing microbial diversity (Fig. 4g–i).

For MCS, I provided the closest match to the expected community, while Z and MN were less accurate. Q deviated most strongly from the theoretical composition. In GMS, MN and Z produced results closest to expectations, although neither was ideal. I, which performed well with MCS, did not perform as effectively with GMS (Fig. 7).

These findings suggest that sample type and purification strategy play a key role in kit performance. The protective medium used for MCS and GMS may not be optimal for any kit, possibly reducing efficiency. Extremely low yields from some kits may also explain discrepancies in microbial representation. This shows that synthetic standards are not always appropriate for choosing DNA extraction kits.

Overall, Z is recommended for most purposes, since it delivers high yield, stable output, and long DNA fragments, which are particularly important for long-read sequencing and diversity analysis. MN and I also perform well, but each has strengths and weaknesses that vary with sample type and application (Table 2). Q proved to be the least reliable, with low-quality DNA and poor representation of the microbial community.

Therefore, the choice of extraction kit should depend on the type of sample and the requirements of the sequencing application.

We also observed that ONT 16S rRNA amplicon libraries required the least experimental work but had lower read quality than Illumina or PacBio. Illumina WGS and PacBio amplicon libraries required more effort but consistently produced high-quality results. The PerkinElmer (PE) V1–V3 protocol often generated nonspecific DNA fragments, reducing performance, while the other methods showed strong quality and consistency.

Our analysis, consistent with other reports^{19,114}, showed that WGS captures more taxonomic diversity than 16S sequencing. However, full-length 16S sequencing provides a clearer view of bacterial communities compared with SRS targeting only one or a few variable regions. For this reason, we included both SRS and LRS approaches, as well as 16S rRNA amplicon sequencing and metagenomic WGS.

We found that WGS data obtained from I and MN extractions closely matched the microbial profiles seen in the Z V3-V4 libraries. Similarly, ONT V1-V9 libraries based on I and MN DNA showed comparable community structures. Libraries made from Z DNA were generally consistent with those from the other kits, demonstrating that Z produces stable microbial community profiles across different library preparations.

The V1-V3 libraries showed considerable variability across all extraction kits, regardless of the method used. This highlights the importance of carefully selecting library preparation methods for reliable microbial profiling. A detailed evaluation of 16S V1-V9 sequencing with ONT and PacBio showed that only the Z kit gave consistent results across both platforms. We also

compared V1-V2 and V1-V9 libraries on MCS and found that the I kit produced the closest match to theoretical expectations in both. MN resembled Z more closely and differed significantly from I in both SRS and LRS. These observations emphasize the different outcomes seen between *in vitro* and *in vivo* experiments and suggest that MCS may not be suitable for validation in other systems.

The α -diversity analysis showed that SRS libraries produced generally consistent values, except where differences came from library preparation. PacBio and Illumina methods displayed distinct diversity profiles that influenced microbial community composition. β -diversity analyses further indicated that library preparation shaped clustering more than DNA extraction. The clear separation between V1-V2 and V3-V4 underlined how strongly library methods affect microbial profiles. The contrasting patterns in microbial diversity and composition between V1-V3 and V3-V4 showed that sequencing method plays a decisive role. The consistent clustering by library rather than extraction method highlighted the importance of selecting the right preparation protocol for accurate analysis.

The importance of the bioinformatic pipeline in metagenomics has been noted before⁸⁹. Our findings also show that roughly 60% of the variation in microbial profiles comes from computational choices. Emu and minitax consistently produced closely matching results, and for MCS they aligned well with theoretical compositions of all amplicon libraries. Sourmash and minitax also showed strong agreement for WGS, demonstrating the broad applicability of minitax. In contrast, results from DADA2 differed greatly from both minitax and Emu *in vivo* and *in vitro*.

For specific sequencing data types, Emu was the most reliable for amplicon sequencing, performing well across V-regions. Although this study did not focus in depth on WGS bioinformatic pipelines, sourmash—previously highlighted by Portik and colleagues⁸⁹ as highly accurate—performed effectively in our work. This suggests that Emu is best suited for amplicon data, while sourmash is highly reliable for WGS. Minitax proved to be a versatile option, consistently working across library types and sequencing methods. Although it did not always give the very highest correlations, it often outperformed other pipelines by producing stable results. It was also effective in handling both amplicon and WGS data, reaching genus-level resolution even with genome-wide databases. This flexibility makes it a strong candidate for comparative studies spanning different sequencing strategies.

In conclusion, our study provides a detailed evaluation of gut microbiome analysis and shows that optimizing and standardizing methods is crucial for accuracy and reproducibility. At the same time, applying multiple complementary methods can help overcome the limitations of each individual approach, leading to a deeper understanding of the microbiome¹¹⁵.

We recommend the following wet-lab pipelines for gut microbiome profiling:

1. For cost-effective studies, use Z DNA extraction with ONT V1-V9 libraries. This provides an affordable option, suitable for longitudinal designs where maximum precision is not essential.
2. For a balanced approach, apply I DNA extraction with V3-V4 libraries on Illumina. This combines moderate cost with solid accuracy, offering dependable microbial profiling.
3. For the most precise projects, use MN DNA extraction with Illumina DNA Prep (WGS). Though more expensive, this provides the most detailed and accurate results.

For bioinformatics analysis, we recommend:

1. For cross-platform comparisons, use minitax with the NCBI genome collection, as it performed well across both amplicon and WGS data.
2. For amplicon sequencing, Emu provides the most reliable results.
3. For WGS, sourmash is highly effective, as noted by previous studies and supported by our findings.

We did not fully exhaust the dataset in all analyses, making it a potentially valuable resource for further investigation by the research community.

4. Summary

In this study, various laboratory and bioinformatic methods for metagenomic profiling of the gut microbiome were compared in order to determine how technical decisions influence the obtained results. The investigations were carried out on different sample types, including dog fecal samples and artificially constructed microbial communities, using several DNA extraction kits, library preparation strategies, and sequencing platforms. The data obtained in this way were processed with different bioinformatic approaches, allowing for a detailed evaluation of the effects of laboratory and computational steps. Based on the results, it can be concluded that a universally best method cannot be identified, as the chosen strategy depends on the sample type and the research objective.

During DNA extraction, the Zymo Research Quick-DNA HMW MagBead Kit proved to be the most reliable, as with its high yield and good quality it provided a stable basis for both short- and long-read sequencing procedures. Other kits, such as Qiagen QIAamp Fast DNA Stool Mini Kit, produced less accurate results. In the case of comparing sequencing platforms, the accuracy of 16S sequencing, combined with long-read sequencing, enables more in-depth analyses.

In terms of bioinformatic processing, it was demonstrated that the applied software and reference database explain a significant part of the variability, meaning that the computational methodology plays a decisive role in the final results. Our developed tool, “minitax,” showed consistent and stable performance, while with the Oxford Nanopore EPI2ME software, frequent biases were observed.

Overall, it can be concluded that in microbiome research, the combined application of multiple complementary methods is justified, since only in this way can sufficient accuracy and reproducibility be ensured. The results highlight that the success of metagenomic studies is determined by the entire methodological chain, therefore in future research standardized procedures and the critical comparison of different techniques are of key importance.

5. Funding

The project was funded by the HAS (MTA) Lendület Programme LP2020-8/2020.

6. Acknowledgements

First and foremost, I wish to express my sincere gratitude to my supervisor, Dr. Dóra Tombácz, for her invaluable guidance and support both throughout and prior to my doctoral studies. I am equally grateful to Prof. Dr. Zsolt Boldogkői, Professor and Head of the Institute of Medical Biology at University of Szeged Albert Szent-Györgyi Medical School, for granting me the opportunity to conduct my research within his institute, as well as for the generous support and numerous opportunities he has provided.

I would like to express my gratitude to Ákos Dörmő, Tamás Járay, Dr. Zsolt Csabai, Dr. Balázs Kakuk, and Dr. István Prazsák for their contributions and assistance, which greatly supported the realization of the publication forming the foundation of this thesis.

I wish to extend my thanks to all the staff members of the Institute of Medical Biology at University of Szeged Albert Szent-Györgyi Medical School for their work and assistance.

And finally, above all, I wish to express my deepest gratitude and heartfelt thanks to my love, Dr. Diána Szűcs, and to my entire loving family, without whose endless love and support I could not have reached this point. Thank you.

7. References

1. Ursell, L. K. *et al.* The intestinal metabolome: An intersection between microbiota and host. *Gastroenterology* **146**, 1470–1476 (2014).
2. Grice, E. A. & Segre, J. A. The human microbiome: Our second genome. *Annual Review of Genomics and Human Genetics* vol. 13 151–170 (2012).
3. Qin, J. *et al.* A human gut microbial gene catalogue established by metagenomic sequencing. *Nature* **464**, 59–65 (2010).
4. Huttenhower, C. *et al.* Structure, function and diversity of the healthy human microbiome. *Nature* **486**, 207–214 (2012).
5. Nayfach, S. *et al.* A genomic catalog of Earth’s microbiomes. *Nat Biotechnol* **39**, 499–509 (2021).
6. Almeida, A. *et al.* A new genomic blueprint of the human gut microbiota. *Nature* **568**, 499–504 (2019).
7. Franzosa, E. A. *et al.* Sequencing and beyond: integrating molecular ‘omics’ for microbial community profiling. *Nat Rev Microbiol* **13**, 360–372 (2015).
8. Hou, K. Microbiota in health and diseases. *Sig. Transduct. Target Ther.* **7**, (2022).
9. Mendes, R. *et al.* Deciphering the Rhizosphere Microbiome for Disease-Suppressive Bacteria. *Science* **332**, 1097–1100 (2011).
10. Adrio, J. L. & Demain, A. L. Microbial Enzymes: Tools for Biotechnological Processes. *Biomolecules* **4**, 117–139 (2014).
11. The Microbial Engines That Drive Earth’s Biogeochemical Cycles | Science. <https://www.science.org/doi/10.1126/science.1153213>.
12. Hillman, E. T., Lu, H., Yao, T. & Nakatsu, C. H. Microbial ecology along the gastrointestinal tract. *Microbes and Environments* vol. 32 300–313 (2017).

13. Laterza, L., Rizzatti, G., Gaetani, E., Chiusolo, P. & Gasbarrini, A. The gut microbiota and immune system relationship in human graft-versus-host disease. *Mediterranean Journal of Hematology and Infectious Diseases* vol. 8 (2016).
14. Ley, R. E., Turnbaugh, P. J., Klein, S. & Gordon, J. I. Human gut microbes associated with obesity. *Nature* **444**, 1022–1023 (2006).
15. Turnbaugh, P. J. *et al.* An obesity-associated gut microbiome with increased capacity for energy harvest. *Nature* **444**, 1027–1031 (2006).
16. Roberfroid, M. B., Bornet, E., Bouley, C. & Cummings, J. H. *Colonic Microflora: Nutrition and Health Summary and Conclusions of an International Life Sciences Institute (ILSI) [Europe] Workshop Held in Barcelona, Spain.* (1995).
17. Fan, Y. & Pedersen, O. Gut microbiota in human metabolic health and disease. *Nature Reviews Microbiology* vol. 19 55–71 (2021).
18. Coelho, L. P. Similarity of the dog and human gut microbiomes in gene content and response to diet. *Microbiome* **6**, (2018).
19. Lewis, S. Comparison of 16S and whole genome dog microbiomes using machine learning. *BioData Min.* **14**, (2021).
20. Costea, P. I. Towards standards for human fecal sample processing in metagenomic studies. *Nat. Biotechnol.* **35**, (2017).
21. Yuan, S., Cohen, D. B., Ravel, J., Abdo, Z. & Forney, L. J. Evaluation of methods for the extraction and purification of DNA from the human microbiome. *PLoS One* **7**, (2012).
22. Zhang, B. *et al.* Impact of Bead-Beating Intensity on the Genus- and Species-Level Characterization of the Gut Microbiome Using Amplicon and Complete 16S rRNA Gene Sequencing. *Frontiers in Cellular and Infection Microbiology* **11**, (2021).
23. 1,500 scientists lift the lid on reproducibility | Nature.
<https://www.nature.com/articles/533452a>.

24. Ducarmon, Q. R., Hornung, B. V. H., Geelen, A. R., Kuijper, E. J. & Zwartink, R. D. Toward standards in clinical microbiota studies: comparison of three DNA extraction methods and two bioinformatic pipelines. *mSystems* **5**, (2020).
25. Tourlousse, D. M. Validation and standardization of DNA extraction and library construction methods for metagenomics-based human fecal microbiome measurements. *Microbiome* **9**, (2021).
26. Quince, C., Walker, A. W., Simpson, J. T., Loman, N. J. & Segata, N. Shotgun metagenomics, from sampling to analysis. *Nat. Biotechnol.* **35**, (2017).
27. Wood, D. E. & Salzberg, S. L. Kraken: ultrafast metagenomic sequence classification using exact alignments. *Genome Biol.* **15**, (2014).
28. Wood, D. E., Lu, J. & Langmead, B. Improved metagenomic analysis with Kraken 2. *Genome Biol.* **20**, (2019).
29. Lu, J. & Salzberg, S. L. Ultrafast and accurate 16S rRNA microbial community analysis using Kraken 2. *Microbiome* **8**, (2020).
30. Brown, C. T. & Irber, L. sourmash: a library for MinHash sketching of DNA. *J. Open Source Softw.* **1**, (2016).
31. Amarasinghe, S. L. *et al.* Opportunities and challenges in long-read sequencing data analysis. *Genome Biology* **21**, 30 (2020).
32. Watson, M. & Warr, A. Errors in long-read assemblies can critically affect protein prediction. *Nat Biotechnol* **37**, 124–126 (2019).
33. Nicholls, S. M., Quick, J. C., Tang, S. & Loman, N. J. Ultra-deep, long-read nanopore sequencing of mock microbial community standards. *Gigascience* **8**, giz043 (2019).
34. Bentley, D. R. *et al.* Accurate whole human genome sequencing using reversible terminator chemistry. *Nature* **456**, 53–59 (2008).

35. Ju, J. *et al.* Four-color DNA sequencing by synthesis using cleavable fluorescent nucleotide reversible terminators. *Proceedings of the National Academy of Sciences* **103**, 19635–19640 (2006).
36. Mardis, E. R. Next-Generation Sequencing Platforms. *Annual Review of Analytical Chemistry* **6**, 287–303 (2013).
37. Shendure, J. & Ji, H. Next-generation DNA sequencing. *Nat Biotechnol* **26**, 1135–1145 (2008).
38. Metzker, M. L. Sequencing technologies — the next generation. *Nat Rev Genet* **11**, 31–46 (2010).
39. Kircher, M. & Kelso, J. High-throughput DNA sequencing – concepts and limitations. *BioEssays* **32**, 524–536 (2010).
40. Goodwin, S., McPherson, J. D. & McCombie, W. R. Coming of age: ten years of next-generation sequencing technologies. *Nat Rev Genet* **17**, 333–351 (2016).
41. Quail, M. A. *et al.* A large genome center’s improvements to the Illumina sequencing system. *Nat Methods* **5**, 1005–1010 (2008).
42. Eid, J. *et al.* Real-Time DNA Sequencing from Single Polymerase Molecules. *Science* **323**, 133–138 (2009).
43. Travers, K. J., Chin, C.-S., Rank, D. R., Eid, J. S. & Turner, S. W. A flexible and efficient template format for circular consensus sequencing and SNP detection. *Nucleic Acids Res* **38**, e159 (2010).
44. Korlach, J. *et al.* Long, Processive Enzymatic Dna Synthesis Using 100% Dye-Labeled Terminal Phosphate-Linked Nucleotides. *Nucleosides Nucleotides Nucleic Acids* **27**, 1072–1083 (2008).
45. Levene, M. J. *et al.* Zero-Mode Waveguides for Single-Molecule Analysis at High Concentrations. *Science* **299**, 682–686 (2003).

46. Chin, C.-S. *et al.* Nonhybrid, finished microbial genome assemblies from long-read SMRT sequencing data. *Nat Methods* **10**, 563–569 (2013).
47. Fuller, C. W. *et al.* The challenges of sequencing by synthesis. *Nat Biotechnol* **27**, 1013–1023 (2009).
48. Wenger, A. M. *et al.* Accurate circular consensus long-read sequencing improves variant detection and assembly of a human genome. *Nat Biotechnol* **37**, 1155–1162 (2019).
49. Kasianowicz, J. J., Brandin, E., Branton, D. & Deamer, D. W. Characterization of individual polynucleotide molecules using a membrane channel. *Proceedings of the National Academy of Sciences* **93**, 13770–13773 (1996).
50. Bayley, H. Nanopore Sequencing: From Imagination to Reality. *Clin Chem* **61**, 25–31 (2015).
51. Clarke, J. *et al.* Continuous base identification for single-molecule nanopore DNA sequencing. *Nature Nanotech* **4**, 265–270 (2009).
52. Branton, D. *et al.* The potential and challenges of nanopore sequencing. *Nat Biotechnol* **26**, 1146–1153 (2008).
53. Manrao, E. A. *et al.* Reading DNA at single-nucleotide resolution with a mutant MspA nanopore and phi29 DNA polymerase. *Nat Biotechnol* **30**, 349–353 (2012).
54. Jain, M., Olsen, H. E., Paten, B. & Akeson, M. The Oxford Nanopore MinION: delivery of nanopore sequencing to the genomics community. *Genome Biology* **17**, 239 (2016).
55. Jain, M. *et al.* Nanopore sequencing and assembly of a human genome with ultra-long reads. *Nat Biotechnol* **36**, 338–345 (2018).
56. Wick, R. R., Judd, L. M. & Holt, K. E. Deepbinner: Demultiplexing barcoded Oxford Nanopore reads with deep convolutional neural networks. *PLOS Computational Biology* **14**, e1006583 (2018).

57. Rand, A. C. *et al.* Mapping DNA methylation with high-throughput nanopore sequencing. *Nat Methods* **14**, 411–413 (2017).
58. Garalde, D. R. *et al.* Highly parallel direct RNA sequencing on an array of nanopores. *Nat Methods* **15**, 201–206 (2018).
59. Biada, I., Santacreu, M. A., González-Recio, O. & Ibáñez-Escriche, N. Comparative analysis of Illumina, PacBio, and nanopore for 16S rRNA gene sequencing of rabbit's gut microbiota. *Front. Microbiomes* **4**, (2025).
60. Wagner, J. *et al.* Evaluation of PacBio sequencing for full-length bacterial 16S rRNA gene classification. *BMC Microbiology* **16**, 274 (2016).
61. Huson, D. H., Auch, A. F., Qi, J. & Schuster, S. C. MEGAN analysis of metagenomic data. *Genome Res.* **17**, (2007).
62. Huson, D. H. MEGAN community edition: interactive exploration and analysis of large-scale microbiome sequencing data. *PLOS Comput. Biol.* **12**, (2016).
63. O'Leary, N. A. *et al.* Reference sequence (RefSeq) database at NCBI: current status, taxonomic expansion, and functional annotation. *Nucleic Acids Res* **44**, D733–D745 (2016).
64. Quast, C. The SILVA ribosomal RNA gene database project: improved data processing and web-based tools. *Nucleic Acids Res.* **41**, (2013).
65. McDonald, D. *et al.* An improved Greengenes taxonomy with explicit ranks for ecological and evolutionary analyses of bacteria and archaea. *ISME J* **6**, 610–618 (2012).
66. Cole, J. R. *et al.* Ribosomal Database Project: data and tools for high throughput rRNA analysis. *Nucleic Acids Res* **42**, D633–D642 (2014).
67. Parks, D. H. *et al.* A complete domain-to-species taxonomy for Bacteria and Archaea. *Nat Biotechnol* **38**, 1079–1086 (2020).

68. Truong, D. T. *et al.* MetaPhlAn2 for enhanced metagenomic taxonomic profiling. *Nat Methods* **12**, 902–903 (2015).
69. Franzosa, E. A. *et al.* Species-level functional profiling of metagenomes and metatranscriptomes. *Nat Methods* **15**, 962–968 (2018).
70. Kanehisa, M., Sato, Y., Kawashima, M., Furumichi, M. & Tanabe, M. KEGG as a reference resource for gene and protein annotation. *Nucleic Acids Res* **44**, D457–D462 (2016).
71. Huerta-Cepas, J. *et al.* eggNOG 4.5: a hierarchical orthology framework with improved functional annotations for eukaryotic, prokaryotic and viral sequences. *Nucleic Acids Res* **44**, D286–D293 (2016).
72. The UniProt Consortium. UniProt: the universal protein knowledgebase in 2021. *Nucleic Acids Res* **49**, D480–D489 (2021).
73. Proctor, L. M. *et al.* The Integrative Human Microbiome Project. *Nature* **569**, 641–648 (2019).
74. Hou, K. *et al.* Microbiota in health and diseases. *Signal Transduct Target Ther* **7**, 135 (2022).
75. Liang, Y. *et al.* Systematic Analysis of Impact of Sampling Regions and Storage Methods on Fecal Gut Microbiome and Metabolome Profiles. *mSphere* **5**, (2020).
76. Tukey, J. W. *Exploratory Data Analysis*. (Addison-Wesley, Reading, MA, 1977).
77. Tukey, J. W. The problem of multiple comparisons. in *Proceedings of the 1953 Biennial Research Conference on the Teaching of Statistics* 15–20 (Princeton University, 1953).
78. Li, H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**, (2018).
79. Morgan, M., Pagès, H., Obenchain, V. & Hayden, N. Rsamtools: Binary alignment (BAM), FASTA, variant call (BCF), and tabix file import. *Bioconductor* Bioconductor <https://doi.org/10.18129/B9.bioc.Rsamtools> (2025).

80. Barrett, T. *et al.* data.table: Extension of 'data.frame'. CRAN (2024).
81. McMurdie, P. J. & Holmes, S. phyloseq: an R package for reproducible interactive analysis and graphics of microbiome census data. *PLoS One* **8**, (2013).
82. Callahan, B. J. DADA2: high-resolution sample inference from Illumina amplicon data. *Nat. Methods* **13**, (2016).
83. Curry, K. D. Emu: species-level microbial community profiling of full-length 16S rRNA Oxford Nanopore sequencing data. *Nat. Methods* **19**, (2022).
84. Oxford Nanopore Technologies. Guppy: Oxford Nanopore basecalling software (version 6.1.5, MinKNOW 20.05.8). (2025).
85. Oxford Nanopore Technologies. EPI2ME pipeline (version 3.6.1). (2025).
86. Biosciences, P. SMRT Link (version 10.2.0.133434).
87. Felix Krueger. Trim Galore.
88. Institute, I. L. R. BMTagger.
89. Portik, D. M., Brown, C. T. & Pierce-Ward, N. T. Evaluation of taxonomic classification and profiling methods for long-read shotgun metagenomic sequencing datasets. *BMC Bioinforma.* **23**, (2022).
90. Fritz, A. CAMISIM: simulating metagenomes and microbial communities. *Microbiome* **7**, (2019).
91. Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biology* **15**, 550 (2014).
92. Menzel, P., Ng, K. L. & Krogh, A. Fast and sensitive taxonomic classification for metagenomics with Kaiju. *Nat. Commun.* **7**, (2016).
93. Wickham, H. *Ggplot2: Elegant Graphics for Data Analysis*. (Springer-Verlag New York, 2016).

94. Rozas, M., Brillet, F., Callewaert, C. & Paetzold, B. MinIONTM nanopore sequencing of skin microbiome 16S and 16S-23S rRNA gene amplicons. *Front. Cell Infect. Microbiol.* **11**, (2022).
95. Cha, T. Gut microbiome profiling of neonates using nanopore MinION and Illumina MiSeq sequencing. *Front. Microbiol.* **14**, (2023).
96. Salonen, A. Comparative analysis of fecal DNA extraction methods with phylogenetic microarray: effective recovery of bacterial and archaeal DNA using mechanical cell lysis. *J. Microbiol. Methods* **81**, (2010).
97. Kerkhof, L. J. Is Oxford Nanopore sequencing ready for analyzing complex microbiomes? *FEMS Microbiol. Ecol.* **97**, (2021).
98. Boer, R. Improved detection of microbial DNA after bead-beating before DNA isolation. *J. Microbiol. Methods* **80**, (2010).
99. Knudsen, B. E. Impact of sample type and DNA isolation procedure on genomic inference of microbiome composition. *mSystems* **1**, (2016).
100. You, I. & Kim, M. J. Comparison of gut microbiota of 96 healthy dogs by individual traits: breed, age, and body condition score. *Animals* **11**, (2021).
101. Söder, J. Composition and short-term stability of gut microbiota in lean and spontaneously overweight healthy Labrador retriever dogs. *Acta Vet. Scand.* **64**, (2022).
102. Thomson, P., Santibáñez, R., Rodríguez-Salas, C., Flores-Yañez, C. & Garrido, D. Differences in the composition and predicted functions of the intestinal microbiome of obese and normal weight adult dogs. *PeerJ* **10**, (2022).
103. Li, Z. Analysis and comparison of gut microbiome in young detection dogs. *Front. Microbiol.* **13**, (2022).
104. Xu, J. The response of canine faecal microbiota to increased dietary protein is influenced by body condition. *BMC Vet. Res.* **13**, (2017).

105. Matsuo, Y. Full-length 16S rRNA gene amplicon analysis of human gut microbiota using MinION™ nanopore sequencing confers species-level resolution. *BMC Microbiol.* **21**, (2021).
106. Sinha, R. The microbiome quality control project: baseline study design and future directions. *Genome Biol.* **16**, (2015).
107. Monteiro, L., Cabrita, J. & Mégraud, F. Evaluation of performances of three DNA enzyme immunoassays for detection of *Helicobacter pylori* PCR products from biopsy specimens. *J. Clin. Microbiol.* **35**, (1997).
108. Flekna, G., Schneeweiss, W., Smulders, F. J. M., Wagner, M. & Hein, I. Real-time PCR method with statistical analysis to compare the potential of DNA isolation methods to remove PCR inhibitors from samples for diagnostic PCR. *Mol. Cell Probes* **21**, (2007).
109. Nechvatal, J. M. Fecal collection, ambient preservation, and DNA extraction for PCR amplification of bacterial and human markers from human feces. *J. Microbiol. Methods* **72**, (2008).
110. Li, X. Efficiency of chemical versus mechanical disruption methods of DNA extraction for the identification of oral Gram-positive and Gram-negative bacteria. *J. Int. Med. Res.* **48**, (2020).
111. Josefsen, M. H., Andersen, S. C., Christensen, J. & Hoorfar, J. Microbial food safety: potential of DNA extraction methods for use in diagnostic metagenomics. *J. Microbiol. Methods* **114**, (2015).
112. Santiago, A. Processing faecal samples: a step forward for standards in microbial community analysis. *BMC Microbiol.* **14**, (2014).
113. Teng, F. Impact of DNA extraction method and targeted 16S-rRNA hypervariable region on oral microbiota profiling. *Sci. Rep.* **8**, (2018).

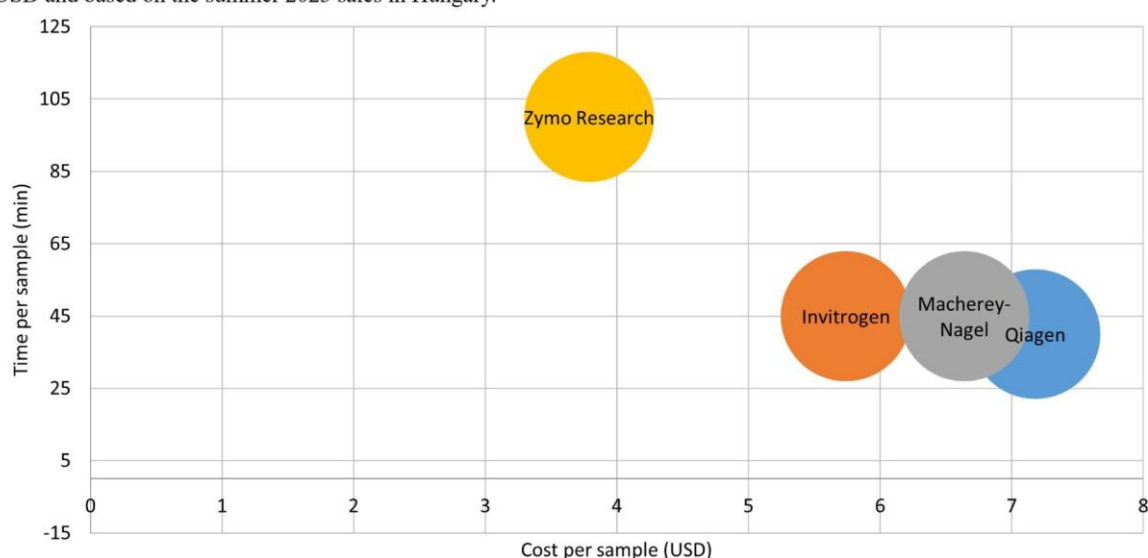
114. Johnson, J. S. Evaluation of 16S rRNA gene sequencing for species and strain-level microbiome analysis. *Nat. Commun.* **10**, (2019).
115. Walker, A. W. & Hoyles, L. Human microbiome myths and misconceptions. *Nat. Microbiol.* **8**, (2023).

8. Supplementary figures

8.1. Supplementary Figure 1

Supplementary Figure 1: Comparative analysis of cost and hands-on time for DNA isolation kits.

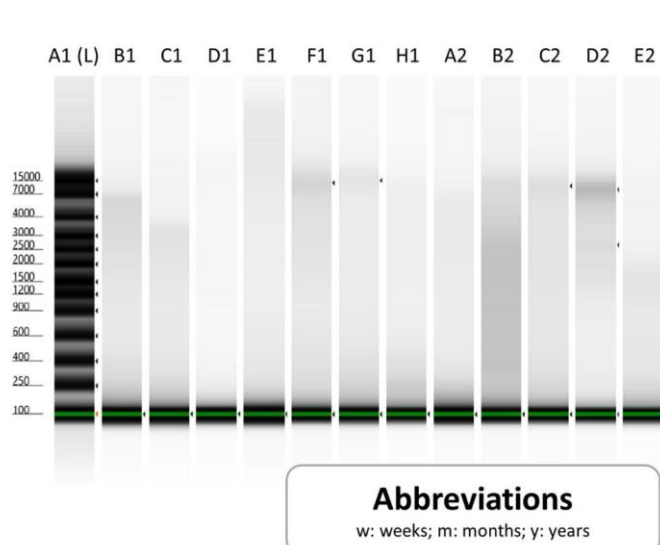
This figure contrasts the cost per sample with the hands-on time required for different DNA isolation kits, with prices presented in USD and based on the summer 2023 sales in Hungary.



8.2. Supplementary Figure 2

Supplementary Figure 2: Analysis of canine stool DNA extraction using Qiagen kit and Agilent ScreenTape Assay.

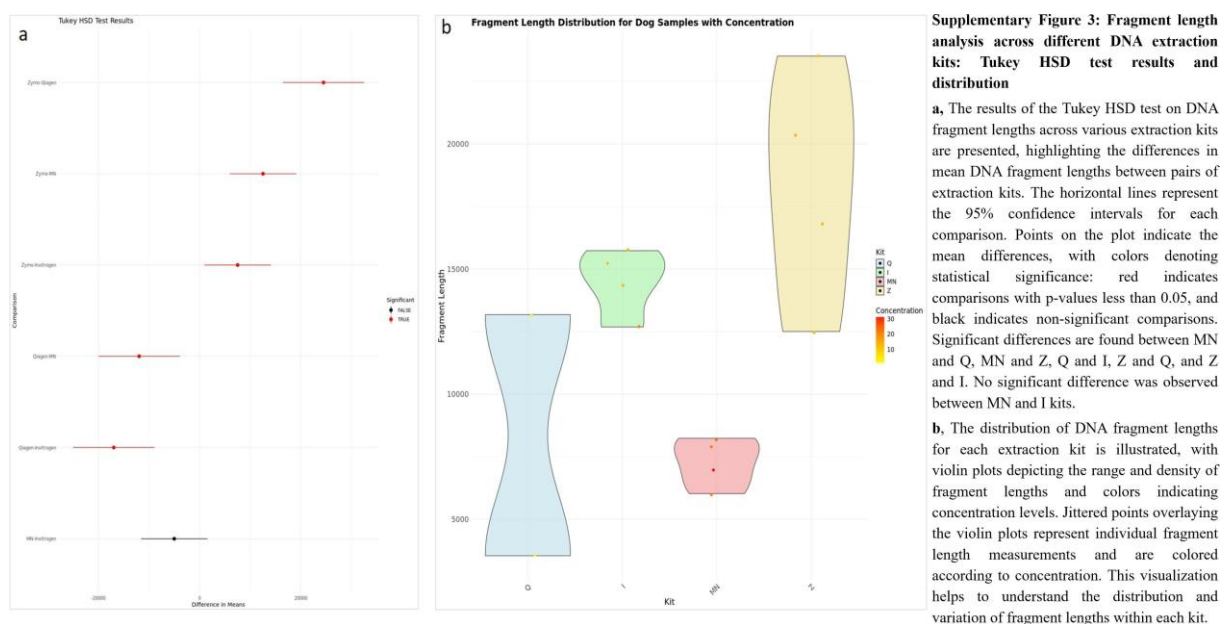
The figure illustrates DNA extracted from control canine fecal samples utilizing the Qiagen kit, which is then analyzed using the Agilent genomic ScreenTape assay. Lane A1 (L) represents the DNA ladder. Abbreviations are as follows: y: years, w: weeks.



	Sample	Age
B1	Female 1	6,5 y
C1	Female 1	6,5 y
D1	Male 1	6 y
E1	Male 1	6 y
F1	Female 2	4 w
G1	Female 2	5 w
H1	Male 2	5 w
A2	Male 2	5 w
B2	Male 3	7 w
C2	Male 3	5 w
D2	Female 3	4 w
E2	Female 3	5 w

Abbreviations
w: weeks; m: months; y: years

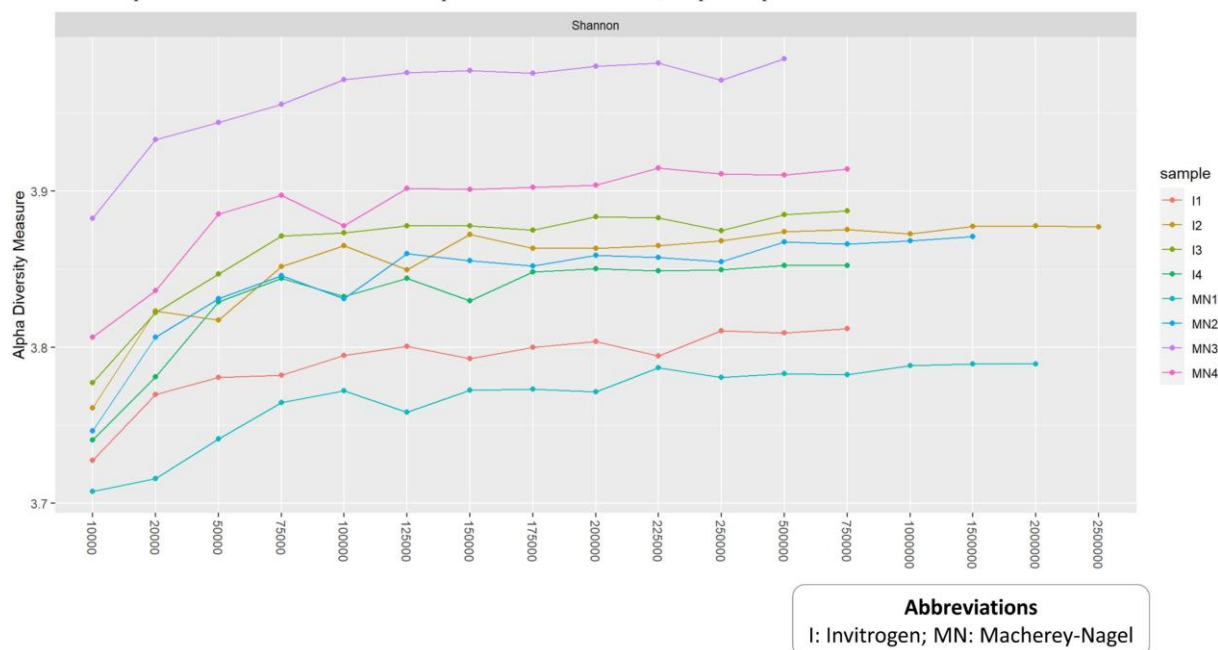
8.3. Supplementary Figure 3



8.4. Supplementary Figure 4

Supplementary Figure 4: Correlation between downsampling and diversity.

The relationship between read count and microbial diversity is explored in this figure, underscoring the opportunity for cost-efficiency in shotgun sequencing. Even with a reduction to approximately 200,000 reads, the diversity across samples remains stable. For assessments focusing on bacterial composition rather than detection of rare species, a read count of 200,000 per sample is sufficient.



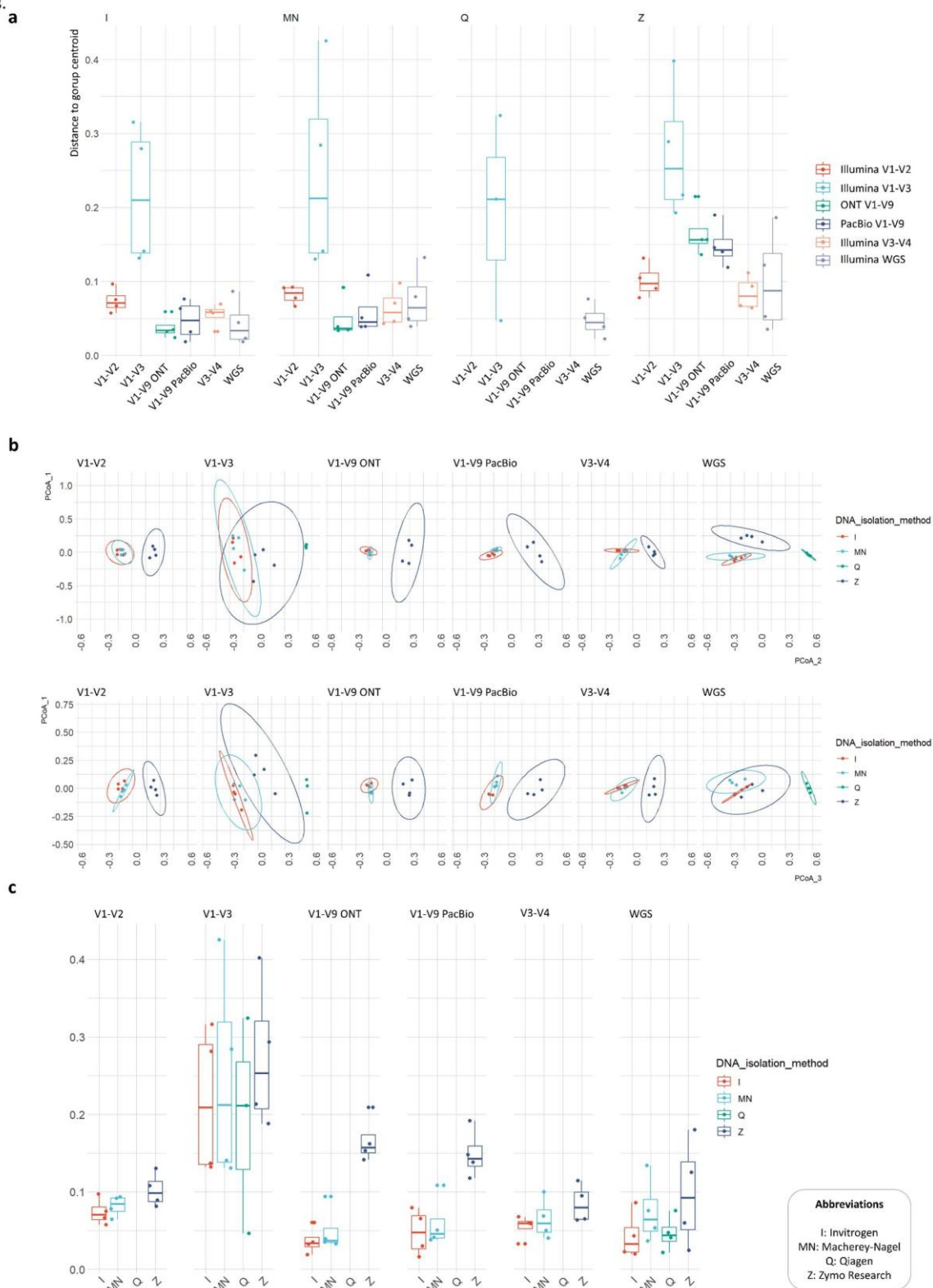
8.5. Supplementary Figure 5

Supplementary Figure 5: β -diversity assessment using Bray-Curtis distance and PERMIDISP analysis

a, The PERMIDISP results display a multivariate analysis of group dispersion homogeneity, indicating the distance to the centroid for each sample in relation to DNA isolation kits and library preparation protocols.

b, The PCoA results depict the variance and resemblance among the DNA isolation protocols for each library preparation.

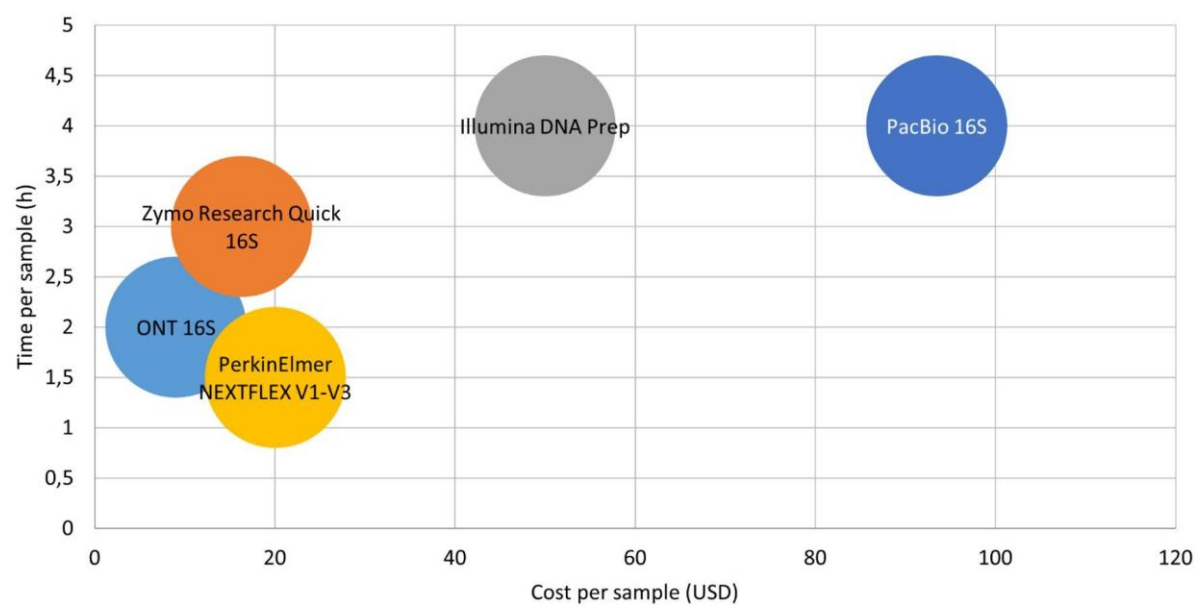
c, The results of the PERMIDISP analysis are presented, categorized by library preparation protocols and DNA isolation kits.



8.6. Supplementary Figure 6

Supplementary Figure 6: Assessment of cost and hands-on time for different library preparation kits

This figure contrasts the cost per sample and the hands-on time associated with various library preparation kits, denoted in USD, based on the summer 2023 sales prices in Hungary.



8.7. Supplementary Figure 7

Supplementary Figure 7: The TapeStation images of the various libraries.

a-j, represent samples from the primary test dog

k-l, indicate samples from the control dogs.

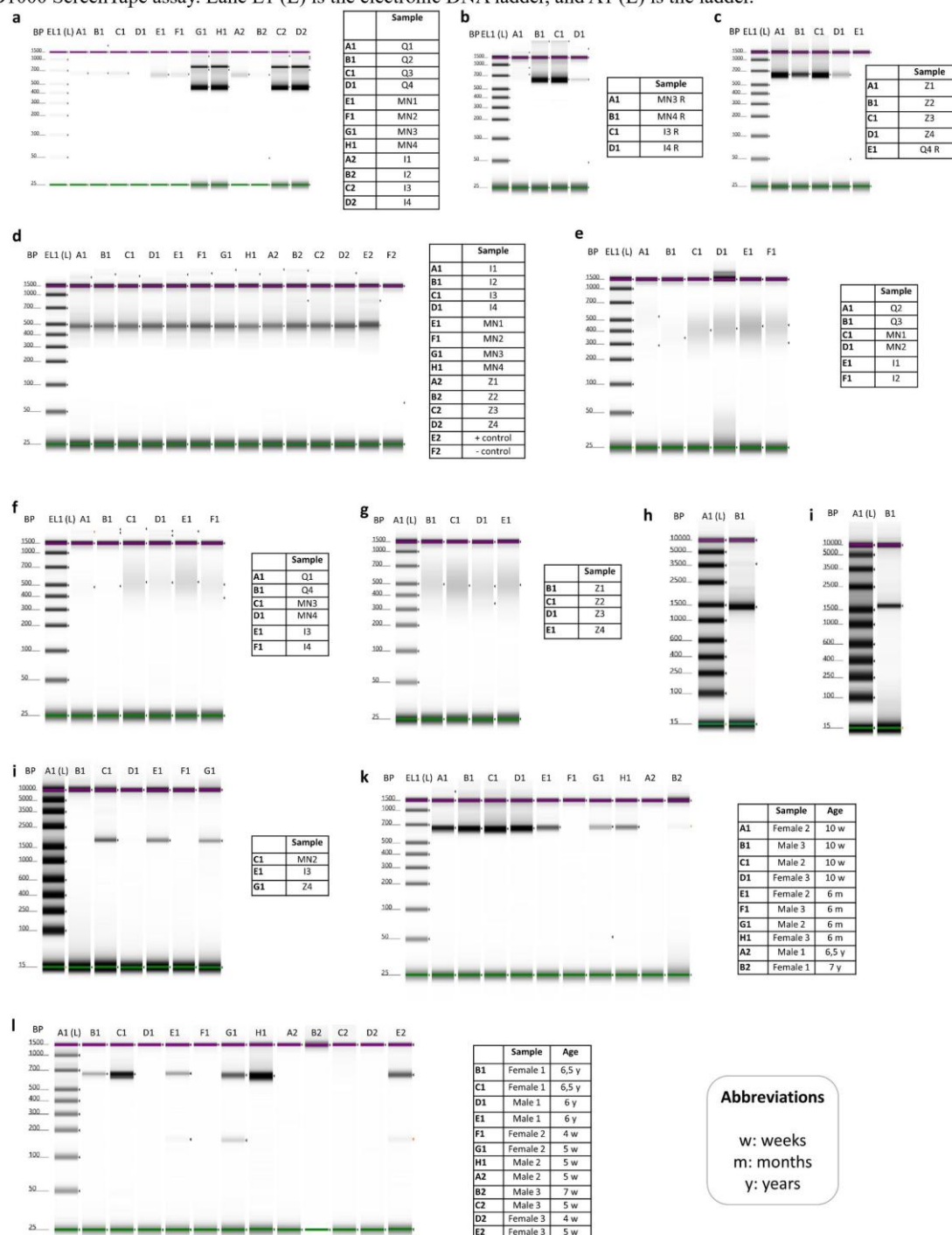
a-c, 16S V1-V3 libraries created using the PerkinElmer NEXTFLEX® Kit and analyzed with the Agilent D1000 ScreenTape assay. Lanes EL1 (L) represent electronic DNA ladders. R indicates a technical replicate (Dog stool).

d, Zymo Research V1-V2 libraries examined with the Agilent D1000 ScreenTape assay. Lane EL1 (L) represents the electronic DNA ladder.

e-g, Shotgun (mWGS, Illumina DNA Prep) libraries derived from canine stool samples.

h-j, V1-V9 library pools from dog stool prepared for ONT (h) and PacBio (i) sequencing. Panel h displays a pool of sixteen libraries, with four each prepared from DNA isolated using Q, I, MN, and Z DNA purification kits for ONT sequencing. Panel i shows twelve PacBio libraries (four each from I, MN, and Z-prepped DNAs). The PCR products were approximately 1.5 kb in size. Abbreviations: MN: Macherey-Nagel, I: Invitrogen, Z: Zymo Research.

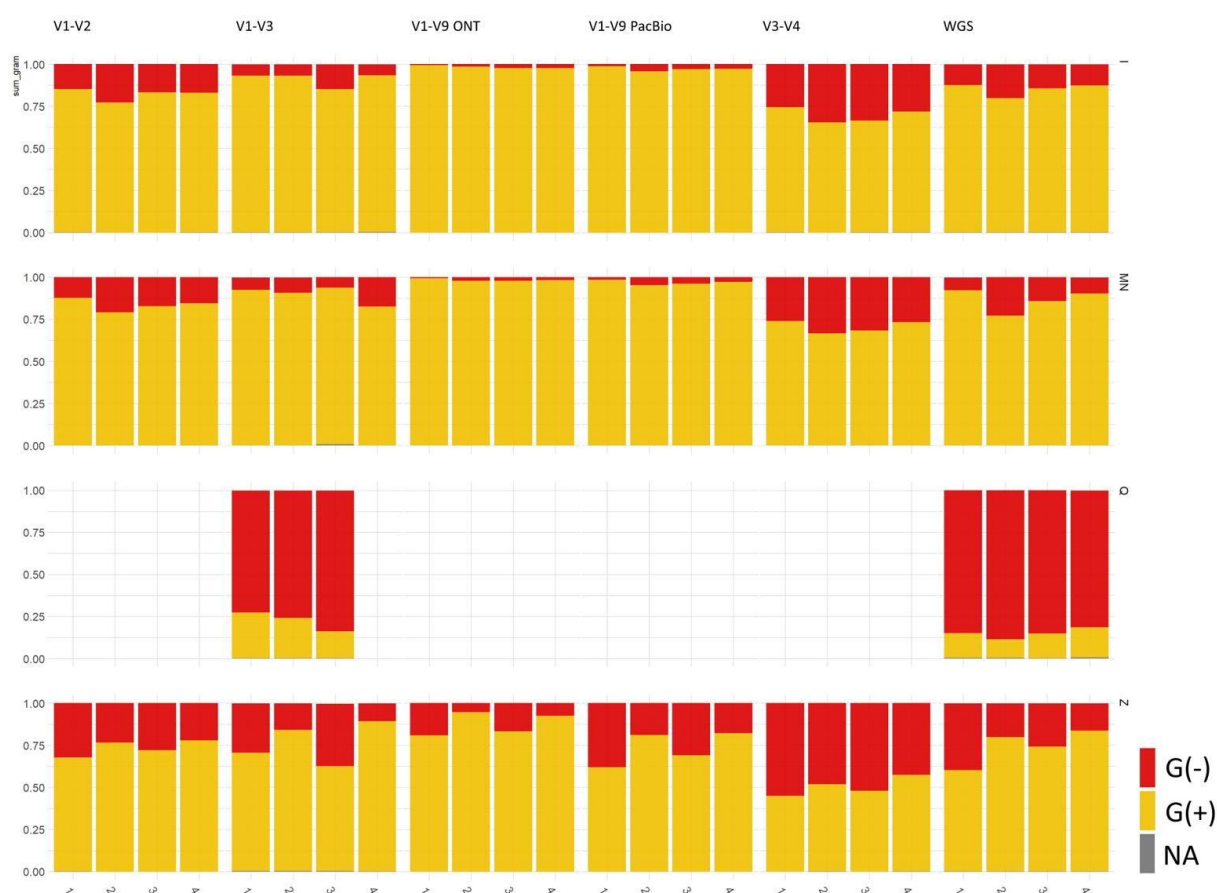
k-l, V1-V3 libraries made using the PerkinElmer NEXTFLEX® 16S Amplicon-Seq Kit and analyzed with the Agilent D1000 ScreenTape assay. Lane E1 (L) is the electronic DNA ladder, and A1 (L) is the ladder.



8.8. Supplementary Figure 8

Supplementary Figure 8: Ratio of Gram-positive and Gram-negative bacteria in canine stool as a result of various laboratory methods

The barplots show the ratio of Gram(+) and Gram(-) bacteria in each sample, according to library preparation protocols and DNA isolation methods. Zymo's results indicate a slight overrepresentation of Gram negatives compared to the Macherey-Nagel and Invitrogen methods, with the Qiagen method showing a significant overrepresentation.



8.9. Supplementary Figure 9

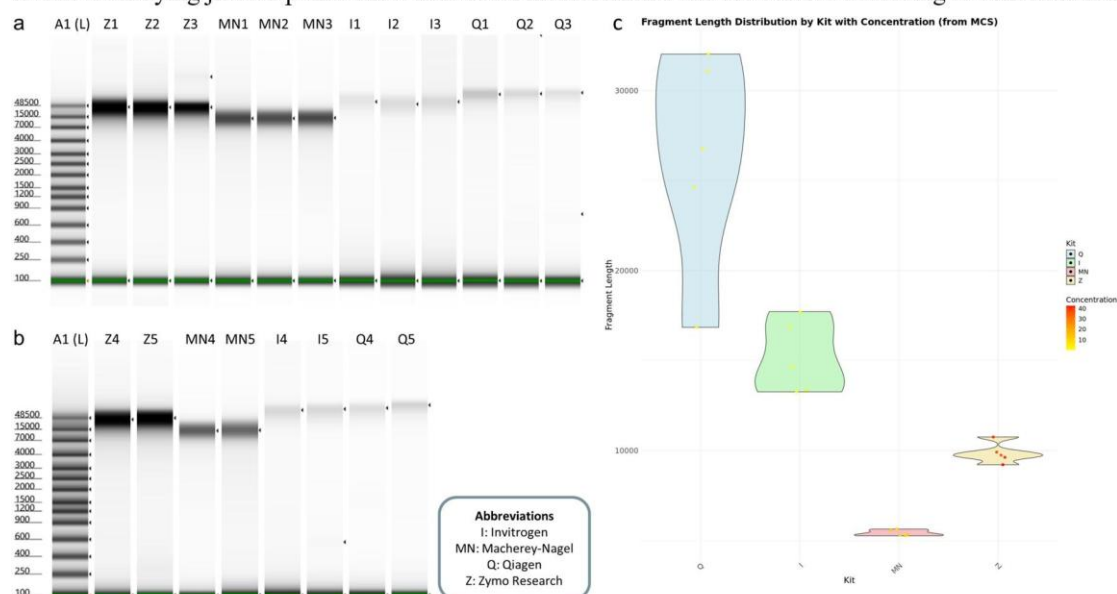
Supplementary Figure 9: Comparing DNA yields and fragment lengths from MCS samples

a, b The images depict the Agilent genomic DNA ScreenTape assay results for DNA extracted from MCS. Lanes A1 (L) represent DNA ladders. The Invitrogen and Qiagen samples showed lower DNA amounts compared to others, with Zymo yielding the highest amount. Despite Zymo Kit being designed for HMW DNA isolation, Invitrogen and Qiagen yielded longer DNA fragments from Zymo MCS samples.

a, First three replicates.

b, Two subsequent replicates.

c, The distribution of DNA fragment lengths for each extraction kit is depicted through a violin plot. The shape and width of each violin indicate the range and density of fragment lengths, with colors representing the concentration levels. Overlaying jittered points show individual measurements and are colored according to concentration.



8.10. Supplementary Figure 10

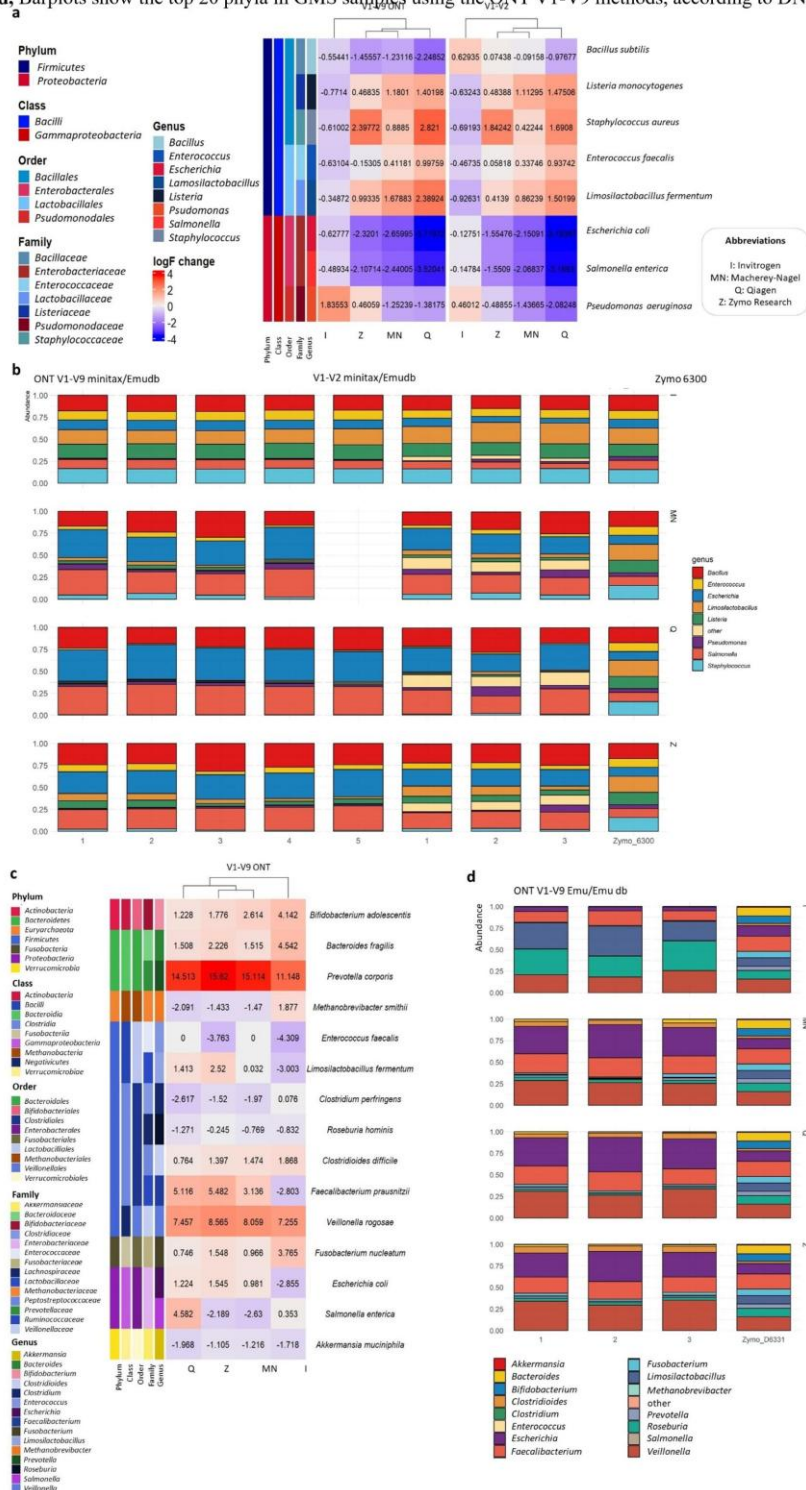
Supplementary Figure 10. Comparative analysis of DNA isolation and library preparation protocols in MCS and GMS using *Emu* with the *Emu* database

a, Heatmap of differences between the theoretical and experimental abundances on according to each species in the MCS. The abundance values were compared to the theoretical values provided by Zymo and \log_2 fold differences were estimated and are shown within the boxes. Deeper blue colors indicate lower experimental values compared to the theoretical, while more red colors indicate higher experimental values.

b, Barplots showing the top 20 phyla in MCS samples using the Illumina V1-V2 and ONT V1-V9 methods, according to DNA isolation kit.

c, Heatmap of differences between the theoretical and experimental abundances on according to each species in the GMS. The abundance values identified by the *Emu* tool were compared to the theoretical values provided by Zymo Research and \log_2 fold differences were estimated and are shown within the boxes. Darker blue colors represent lower experimental values compared to the theoretical, while dark red colors indicate higher experimental values.

d, Barplots show the top 20 phyla in GMS samples using the ONT V1-V9 methods, according to DNA isolation kit.



8.11. Supplementary Figure 11

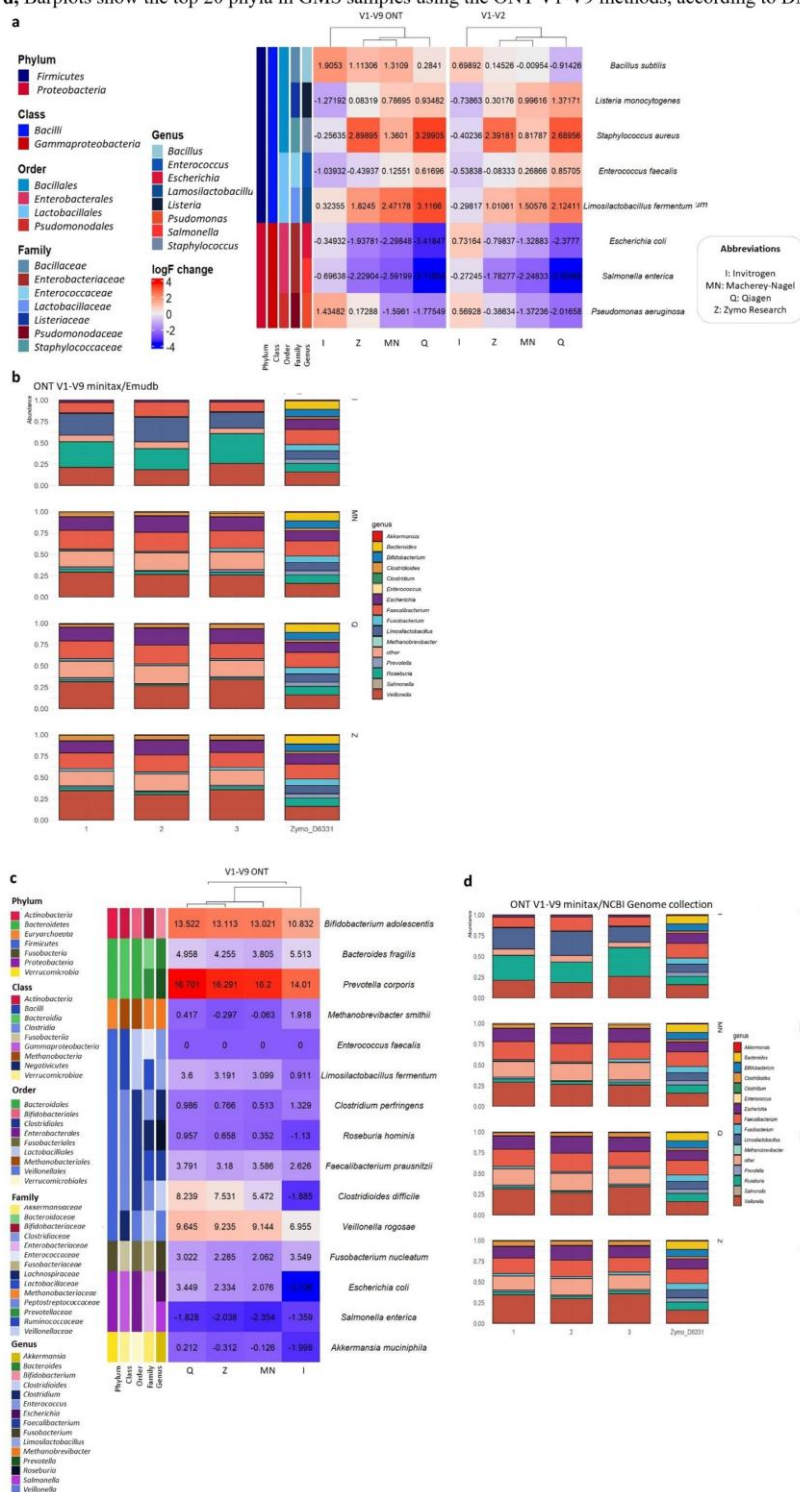
Supplementary Figure 11. Comparative analysis of DNA isolation and library preparation protocols in MCS and GMS using *minitax* with the NCBI genome collection database

a, Heatmap of differences between the theoretical and experimental abundances on according to each species in the MCS. The abundance values were compared to the theoretical values provided by Zymo and \log_2 fold differences were estimated and are shown within the boxes. Deeper blue colors indicate lower experimental values compared to the theoretical, while more red colors indicate higher experimental values.

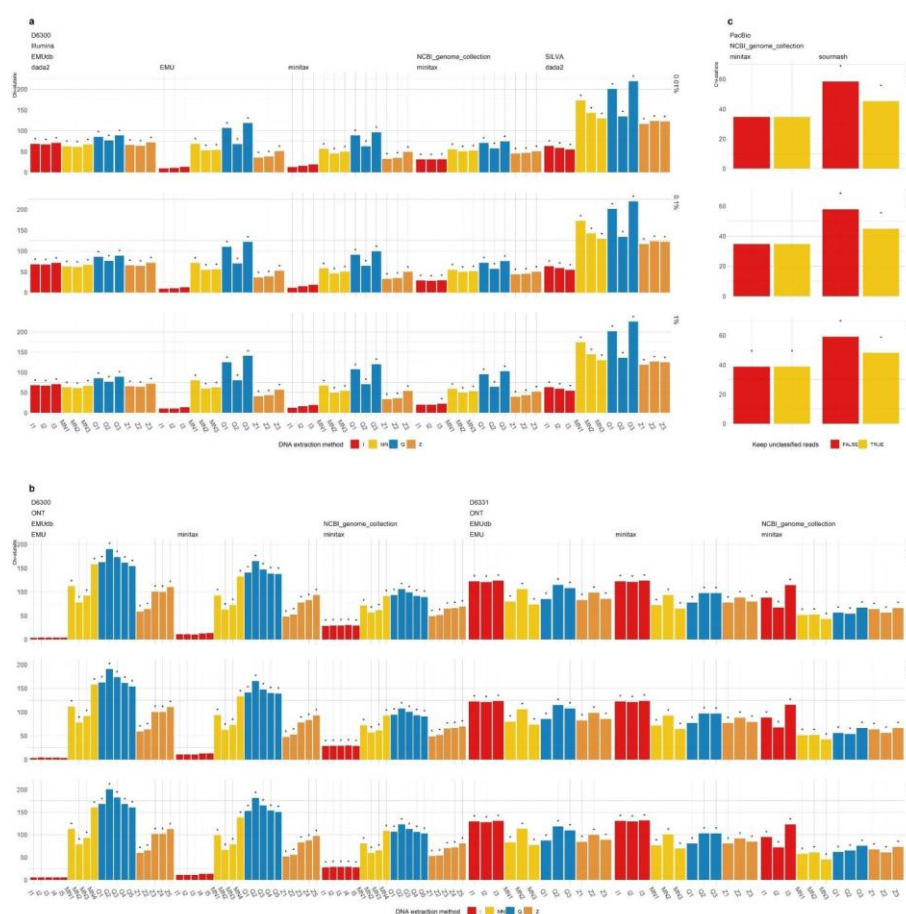
b, Barplots showing the top 20 phyla in GMS samples using the ONT V1-V9 methods, according to DNA isolation method.

c, Heatmap of differences between the theoretical and experimental abundances on according to each species in the GMS. The abundance values were compared to the theoretical values provided by Zymo Research and \log_2 fold differences were estimated and are shown within the boxes. Darker blue colors represent lower experimental values compared to the theoretical, while dark red colors indicate higher experimental values.

d, Barplots show the top 20 phyla in GMS samples using the ONT V1-V9 methods, according to DNA isolation method.



8.12. Supplementary Figure 12



Supplementary Figure 12: Evaluating minitax: a comparison with other methods based on Pearson's correlations.

This figure compares minitax with other methods based on the Pearson's correlations between theoretical and observed compositions (r^2 values).

a, Comparison with Emu on ONT V1-V9 Sequencing of Zymo D3600 MCS data.

b, Comparison with DADA2 and Emu on Illumina V1-V2 Sequencing of Zymo D3600 MCS.

c, Comparison with sourmash on PacBio HiFi WGS of Zymo MCS D6331.

9. Supplementary methods

9.1. QIAGEN QIAamp Fast DNA Stool Mini Kit

Canine stool: A 200 mg stool sample was placed in a 2 ml microcentrifuge tube and kept on ice. InhibitEX Buffer was added to each sample, and they were mixed using a vortex until completely homogenized. The large stool pieces were removed by centrifugation. Six hundred μL of the supernatant was combined with 25 μL of proteinase K and 600 μL of Buffer AL. The mixture was thoroughly vortexed, then heated to 95°C for 5 min (this is a 95°C lysis incubation for 5 minutes, diverging from the 70°C recommended by the QIAamp kit). Six hundred μL of 100% ethanol was added to the lysate and then mixed. Six hundred μL of the lysate was loaded onto a QIAamp spin column and centrifuged at 20,000 \times g for 1 min. The QIAamp spin column was placed into a new 2 ml tube. The remainder of the lysate was then loaded onto the column. After centrifugation, 500 μL of Buffer AW1 was added to the column. This was followed by a 20,000 \times g centrifugation for 1 min, then we discarded the collection tube. Next, 500 μL of Buffer AW2 was added to the column, which was then placed into a new collection tube. A full-speed centrifugation (20,000 \times g) was performed for 3 min. To avoid any Buffer AW2 carryover, the spin column was set in a fresh 2 ml collection tube and the samples were spun down at full speed for 3 min. The spin columns were then transferred to new Eppendorf tubes, and 100 μL of Buffer ATE was directly loaded onto the QIAamp membrane. After letting it incubate at room temperature for 1 min, a centrifugation (20,000 \times g) step was carried out for 1 min to elute the DNA in 50 μL in elution buffer, followed by storing the DNA solution at -20°C.

MCS and GMS samples: 75 μL from the mixtures was used for DNA purification, following the same protocol as for the dog sample, with the DNA being eluted in a final volume of 50 μL .

9.2. Invitrogen PureLink™ Microbiome DNA Purification Kit

Canine stool: A 200 mg sample was combined with 600 μL of S1 Lysis Buffer in the Bead Tube (provided in the kit) and the mixture was homogenized by vortexing. Subsequently the 100 μL of S2 Lysis Enhancer (from the kit) was added, and the samples were vortexed again. The mixtures were then incubated at 65°C for 10 minutes. For homogenization, the samples were subjected to bead beating using a vortex mixer with horizontal agitation at maximum speed for 10 min. The samples were then centrifuged at 14,000 \times g for 5 min. Afterwards, 400 μL of the supernatant was transferred to a new Eppendorf tube and mixed with 250 μL of S3 Cleanup Buffer. The samples were centrifuged again at 14,000 \times g for 2 min, and 500 μL of the resulting supernatant was transferred to a clean tube and mixed with 900 μL of S4 Binding Buffer. After

brief vortexing, 700 μ L of the sample mixture was loaded onto a spin column and centrifuged at $14,000 \times g$ for 1 min. The spin column was then placed in a new tube, and the remaining sample mixture was loaded onto it for an additional 1 min centrifugation. The spin column was subsequently placed in a clean collection tube, and 500 μ L of S5 Wash Buffer was added, followed by centrifugation at $14,000 \times g$ for 1 min. To remove any residual S5 Wash Buffer, a second centrifugation was carried out at $14,000 \times g$ for 30 sec. Finally, DNA was eluted from the spin columns using 100 μ L of S6 Elution Buffer and stored at -20°C .

MCS and GMS samples: 75 μ L of the microbial mixes was utilized for DNA extraction, following the same steps as described above. The DNA was eluted in 50 μ L S6 Elution Buffer.

9.3. Macherey-Nagel NucleoSpin DNA Stool Mini kit

Canine stool: DNA isolation was performed using 200 mg of fecal sample, which was transferred to a Macherey-Nagel Bead Tube Type A, and then 850 μ L Buffer ST1 was added. The mixtures were shaken horizontally for 3 seconds before being placed in a heat incubator. Subsequently, the samples were incubated for 5 min at 70°C , agitated on a Vortex-Genie® 2 at full speed and room temperature for 10 min, and then centrifuged for 3 min at $13,000 \times g$. Six hundred μ L of the supernatant was transferred to a new 2 ml tube, and 100 μ L of Buffer ST2 was added and briefly vortexed. The mixtures were incubated for 5 min at 4°C and then centrifuged for 3 min at $13,000 \times g$. Five hundred fifty μ L of the lysate was loaded onto a NucleoSpin® Inhibitor Removal Column and centrifuged for 1 min at $13,000 \times g$. The Inhibitor Removal Column was discarded. Two hundred μ L of Buffer ST3 was added to the samples, which were then mixed. Seven hundred μ L of the sample mixture was loaded onto a NucleoSpin® DNA Stool Column and centrifuged for 1 min at $13,000 \times g$. The column was placed in a new tube. The sample was washed four times: first, 600 μ L of Buffer ST3 was added to the NucleoSpin® DNA Stool Column and centrifuged for 1 min at $13,000 \times g$. The column was then placed in a new tube, and 550 μ L of Buffer ST4 was added. After a 1 min centrifugation at $13,000 \times g$, the column was placed into a new tube, and 700 μ L of Buffer ST5 was added. Following a brief vortexing, the samples were centrifuged for 1 min at $13,000 \times g$. The column was placed in a new tube, and 700 μ L of Buffer ST5 was added, followed by a 1 min centrifugation at $13,000 \times g$. The flow-through was discarded, and the column was placed back onto the tube. The silica membrane of the column was dried by a 2 min centrifugation at $13,000 \times g$. One hundred μ L of Buffer SE was loaded onto the center of the column, and the DNA was eluted by centrifugation for 1 min at $13,000 \times g$. DNA samples were stored at -20°C .

MCS and GMS samples: 75 µl of the MCS or the GMS mixture was utilized for DNA isolation, with the sample being eluted in a final volume of 50 µl.

9.4. Zymo Research Quick-DNA™ HMW MagBead Kit

Canine stool: One hundred mg of the fecal sample was used as initial weight and resuspended in 200 µl of DNA/RNA Shield™, followed by incubation at room temperature (20-30°C) on a tube rotator for 5 minutes. Subsequently, 33 µl of MagBinding Beads were added to each sample, mixed and placed on a shaker for a 10 min. The sample was then placed on a magnetic stand until a clear separation of the beads and the solution was observed, after which the supernatant was removed. For the washing step, 500 µl of Quick-DNA™ MagBinding Buffer was added and the beads were resuspended and shaken for 5 min. The sample was returned to the magnetic stand, and the supernatant was discarded. Next, 500 µl of DNA Pre-Wash Buffer was added, and the beads were resuspended. The sample was placed on the magnetic stand again, and the supernatant was discarded. In the subsequent step, 900 µl of g-DNA Wash Buffer was added and mixed, and the entire liquid was transferred to a new tube. The magnetic stand was used to separate the beads from the solution, and the supernatant was discarded. These washing steps were repeated once more. The sample were left to air dry for 20 min. For the elution step, 50 µl of DNA Elution Buffer was added to each sample. After mixing, the solution was incubated at room temperature for 5 min. Finally, the sample was placed back on the magnetic stand until the beads separated from the solution. The eluted DNA was carefully transferred to a new tube and stored at -20°C for future use.

MCS and GMS samples: 75 µl of the mix was used for DNA purification, following the same protocol as described for the dog sample. The DNA was eluted in a final volume of 50 µl.

9.5. LIBRARY PREPARATION

9.5.1. From partial regions of the 16S rRNA gene

9.5.1.1. Zymo Research V1-V2

DNA from canine stool: Ten µl of Quick-16S™ qPCR Premix were mixed with 4 µl of Quick-16S™ Primer Set V1-V2 and 4 µl of ZymoBIOMICS® DNase/RNase Free Water. Additionally, 2 µl of DNA samples (2.5 ng/µl) were added. PCR was conducted in a Verity Thermal Cycler (Applied Biosystems) as per the Zymo Research Manual (Supplementary Data 9). After amplification, 1 µl of Reaction Clean-up Solution was added to the samples, which were then incubated at 37°C for 15 min. The reactions were terminated by heating to 95°C for

10 min, and the samples were subsequently cooled to 4°C. Next, 10 µl of Quick-16S™ qPCR Premix and 4 µl of ZymoBIOMICS® DNase/RNase Free Water were combined. Index primers (2 µl each from ZA5 and ZA7, Supplementary Data 10 for detailed pairs and sequences) and 2 µl of the amplified DNA were also measured into the mixture. Barcoded PCR reactions were performed as recommended by the manual (Supplementary Data 11). For purification of the PCR products, Select-a-Size MagBeads were used. First, the MagBeads were resuspended by shaking, and then 16 µl of the Select-a-Size MagBeads were mixed with each sample. The mixture was incubated at room temperature for 5 min and then placed on a magnetic rack for 3 to 10 min. The supernatant was discarded and the beads were washed twice with 200 µl of DNA Wash Buffer. The samples were removed from the magnet and were incubated for 3 min at room temperature to eliminate all traces of buffer. Libraries were eluted in 25 µl of DNA Elution Buffer and stored at -20°C until further use.

DNA from MCS samples: 2 ng/µl DNA was used to prepare the V1-V2 libraries from microbial mixture

9.5.1.2. Zymo Research V3-V4

The protocol used was the same as described in the ‘Zymo Research V1-V2’ section with the following modifications: In the initial PCR step, V3-V4 primers were utilized. Supplementary Data 9 provides details of the barcoded primers used, including the pairs and their sequences.

9.5.1.3. PerkinElmer NEXTFLEX® 16S V1-V3 Amplicon-Seq Kit for Illumina

Genomic DNA, having concentrations between 1.6 ng and 36 ng as specified in Supplementary Data 12, was diluted using Nuclease-free Water to maintain a total volume no greater than 36 µL. Subsequently, 12 µL of NEXTFlex™ PCR Master Mix and 2 µL of the 16S V1-V3 PCR I Primer Mix were added to the solution. The final reaction volume was adjusted to 50 µL. First amplification step of the PCR cycling was carried out using the settings outlined in Supplementary Data 13. PCR cleanup: fifty µL AMPure XP Beads was added to each sample. After mixing, the samples were incubated at room temperature for 5 min. Then, using a magnetic stand, the samples were left until the supernatant clarified. The supernatant was then discarded, and the beads were washed twice with 200 µL of freshly prepared 80% ethanol. Next, the samples were air-dried for 3 min and resuspended in 38 µL of Resuspension Buffer. After a further incubation of 2 minutes at room temperature, 36 µL from the clear supernatant was transferred to fresh tubes. This sample was then subjected to PCR amplification with the addition of 12 µL of NEXTFlex™ PCR Master Mix and 2 µL of NEXTFlex™ PCR II Barcoded

Primer Mix. The procedure was executed following the guidelines specified in Supplementary Data 14. PCR cleanup was carried out in accordance with the purification after the first PCR.

9.5.2. For the analysis of full-length 16S rRNA gene sequencing

9.5.2.1. ONT Rapid Sequencing 16S Barcoding Kit (SQK-RAB204)

Ten ng of high molecular weight genomic DNA (in a 10 μ l volume) was used for library preparation from both the canine, MCS and GMS samples. DNA isolated using the QIAGEN kit did not meet this criterion. The input DNA was mixed with 14 μ l of Nuclease-free water (Invitrogen), 1 μ l of 16S Barcode (1 μ M; Supplementary Data 15)) and 25 μ l of LongAmp Taq 2X master mix (New England Biolabs). PCR Amplification of the samples was carried out according Supplementary Data 16. Amplified DNA samples were transferred to clean 1.5 ml Eppendorf DNA LoBind tubes and mixed with 30 μ l of resuspended AMPure XP beads (Beckman Coulter). Next, they were incubated on a Hula mixer for 5 minutes at room temperature. Tubes were placed on a magnetic rack then the supernatant was discarded. The beads were washed with 200 μ l of freshly prepared 70% ethanol. Ethanol was removed and the washing was repeated once. After air drying of the beads, samples were removed from the magnet and beads were resuspended in 10 μ l of 10 mM Tris-HCl pH 8.0 with 50 mM NaCl. After 2 min incubation at room temperature, samples were placed on the magnet. Ten μ l of the clean supernatant, containing the ONT libraries, was transferred to a new Eppendorf DNA LoBind tube. Barcoded libraries were pooled in equal molar ratio and then, 1 μ l of RAP was added. The reaction was incubated for 5 minutes at room temperature. One hundred fmoles were loaded on a MinION flow cell.

9.5.2.2. PacBio Full-Length 16S Library Preparation Using SMRTbell Express Template Prep Kit 2.0 Sequel IIe System ICS v10.0 / Sequel II Chemistry 2.0 / SMRT Link v10.0

For each sample, 1.5 μ L of PCR-grade water and 12.5 μ L of 2X KAPA HiFi HotStart ReadyMix were mixed. Subsequently, 3 μ L of barcoded forward primer solution (2.5 μ M, sequences in Supplementary Data 17) was added. This was followed by the addition of 3 μ L of the respective reverse primer solution (Supplementary Data 17) and 5 μ L of the DNA sample. DNA amplification was carried out according to the parameters listed in Supplementary Data 18.

9.5.3. Shotgun sequencing

9.5.3.1. Illumina DNA Prep

Sixty ng of DNA (in 30 μ l) was used as total input per sample. Ten μ l Tagmentation Buffer 1 (TB1) was mixed with 10 μ l Bead-Linked Transposomes (BLT), and 20 μ l of this mixture was added to a DNA sample. The mixture was incubated at 55°C for 15 min and then held at 10 °C. Following this, 10 μ l of Tagment Stop Buffer (TSB) was added to the sample and gently mixed. The samples were incubated at 37°C for 15 min, and then kept at 10°C. Next, the samples were placed on a magnetic stand for 3 min, and the supernatant was discarded. The sample was removed from the magnet and 100 μ l of Tagment Wash Buffer (TWB) was slowly added directly onto the beads. The sample was placed back on the magnetic stand, the supernatant was discarded, and the wash step was performed again. A mixture of 20 μ l of Enhanced PCR Mix (EPM) and 20 μ l of nuclease-free water was prepared, and 40 μ l of this mixture was added to the washed beads. Index adapters (i5 and i7, 5 μ l each; Supplementary Data 19) were added, and PCR was carried out according to Supplementary Data 20. After amplification, the libraries were cleaned up. First, the samples were placed on a magnet for approximately 5 min. Forty-five μ l of supernatant from each PCR product was transferred to a new tube. Forty μ l of nuclease-free water and 45 μ l of Sample Purification Beads (SPB) were added to the supernatant and the samples were mixed at 1600 rpm for 1 min. They were incubated at room temperature for 5 min. After this step, the samples were placed on a magnet, and 15 μ l of SPB were added to new tubes. Next, 125 μ l of supernatant from each sample was added to the tubes containing 15 μ l of undiluted SPB. The samples were mixed at 1600 rpm for 1 min, and then incubated at room temperature for 5 min. The supernatant was discarded, and the washing step was carried out twice with 200 μ l of freshly prepared 80% ethanol. The sample was stored on magnetic stand for 30 secs, and then the ethanol was removed. After the second washing step, the pellet was air dried. Next, 32 μ l RSB was added to the beads, and they were mixed and incubated for 2 min. Finally, the sample was placed on a magnetic stand for 2 min, and the supernatant containing the prepared library was transferred to a new tube.

9.6. Data availability

All our sequencing data have been submitted to the ENA under the accession PRJEB59610. The datasets from other sources were downloaded from the ENA with the following accession numbers: PRJNA783735, PRJNA678365, PRJNA871395. Supplementary Data files (1–21): figshare, <https://doi.org/10.6084/m9.figshare.2723262673>. In addition, data used to generate figures can be found under the project's github repository:

<https://github.com/Balays/Microbiome-Method-Comparison>, specifically:—Figure 4 g: Fig4.G-H_MCM_all_Adiversity_estimates.tsv - Figure 4 i: Fig4_statistics_data.rds—Figure 5 a: Fig_5A_BrayCurtis_distances.tsv—Figure 5 b: Fig_5B_PCoA_data.tsv—Figure 5 c: Fig_5C_Heatmap_data.tsv - Figure 5 d: ./Dog_Feces/Microbial_abundances for each taxon level—Figure 5 e: DogFeces_all_methods_PS.rds - Figure 7: sig.freq.tsv and l2F_diff.dt.tsv separately for each group—Figure 7: ./MCM/.../sig.freq.tsv and l2F_diff.dt.tsv separately for each group - Figure 9 a-b: ./MCM/.../Detection_statistics.tsv for each each group—Figure 9 c: ./Portik_etal_2022/ Detection_statistics.tsv - Figure 9 d: ./CAMISIM/Detection_statistics.tsv. In addition, phyloseq objects that can be imported into R and contain the microbial abundances for the respective samples are also available: - MCM_ZymoD6300_all_PS.rds - MCM_ZymoD6331_all_PS.rds - DogFeces_all_methods_PS.rds. The data used to generate Fig. 4e, f are available in Supplementary Data 3, the data for Fig. 6 can be found in Supplementary Data 6, while the data for generate Fig. 8 are in Supplementary Data 21.

9.7. Code availability

minitax: <https://github.com/Balays/minitaxStatistics>: <https://github.com/Balays/Microbiome-Method-Comparison> other in-house scripts - <https://github.com/Balays/Microbiome-Method-Comparison> - <https://github.com/gabor-gulyas/Technical-article-downsample> DADA2: <https://benjjneb.github.io/dada2/> Emu: <https://gitlab.com/treangenlab/emu> Trim Galore: <https://github.com/FelixKrueger/TrimGalore> BMTagger: <https://hpc.ilri.cgiar.org/bmtagger-software>

10. Supplementary datas

10.1. Supplementary Data 1

Supplementary Data 1. This table summarizes the previously published articles focusing on the comparison of library prep kits or bioinformatic methods used at various stages of metagenomic analyses.

a. human and mammalian stool

#	sample	DNA isolation kits	Method	Library preparation kit	Approach	Sequencer	read length	Bioinformatics	Reference
1	human stool and pure culture	FastDNAR kit (Bio 101) NucleospinR C + T kit 1 (Macherey-Nagel) Quantum PrepR Aquapure Genomic DNA isolation kit QIAampR DNA stool minikit Boom et al., 1990	PCR						71. McOrist, A. L., Jackson, M. & Bird, A. R. A comparison of five methods for extraction of bacterial DNA from human faecal samples. <i>J. Microbiol. Methods</i> 50 , 131-9 (2002).
2	human stool	RNA/DNA Mini kit (Qiagen) QIAamp DNA Stool Mini kit (Qiagen) Fecal DNA Isolation kit (Mo BIO)	Real-time PCR						72. Nechvatal, J. M. et al. Fecal collection, ambient preservation, and DNA extraction for PCR amplification of bacterial and human markers from human feces. <i>J. Microbiol. Methods</i> 72 , 124-32 (2008).
3	human stool	Differential Centrifugation and Lysis (Apajalahti et al., 1998) Promega Genomic Wizard DNA Purification Kit (Promega) Repeated Bead Beating (Yu and Morrison) QiaAmp DNA Stool Mini Kit (Qiagen)	Real-time PCR, microarray						58. Salonen, A. et al. Comparative analysis of fecal DNA extraction methods with phylogenetic microarray: Effective recovery of bacterial and archaeal DNA using mechanical cell lysis. <i>J. Microbiol. Methods</i> 81 , 127-134 (2010).
4	stool from 6 months to 4 years kids with diarrhoea	QIAamp DNA stool MiniKit with or without Mini BeadBeater 8 (BioSpec Products Inc) or TissueLyser system (Qiagen Retsch GmbH)	PCR		V3				73. Smith, B., Li, N., Andersen, A. S., Slotved, H.C. & Krogfelt, K. A. Optimising bacterial DNA extraction from faecal samples: comparison of three methods. <i>Open Microbiol. J.</i> 5 , 14-7 (2011).
5	human stool	NucliSENS® easyMag® (BioMérieux) - semi-automatic QIAamp DNA Stool Mini Kit (Qiagen) - manual	PCR		V2-V3				74. Mirsepasi, H. et al. Microbial diversity in fecal samples depends on DNA extraction method: easyMag DNA extraction compared to QIAamp DNA stool mini kit extraction. <i>BMC Res. Notes</i> 7 , 50 (2014).

6	human stool	sample from outer an inner layer and from homogenized stool, with or without bead-beating	pyrosequencing		V4			9. Santiago, A. et al. Processing faecal samples: a step forward for standards in microbial community analysis. <i>BMC Microbiol.</i> 14 , 112 (2014).
7	Human stool	American Human Microbiome Project European MetaHIT project	SRS		WGS	Illumina HiSeq 2000	2x100 Mocat Blast+	14. Wesolowska-Andersen, A. et al. Choice of bacterial DNA extraction method from fecal material influences community structure as evaluated by metagenomic analysis. <i>Microbiome</i> 2 , 19 (2014).
8	human stool and synthetic DNA mock community	NA	SRS	Illumina Nextera XT kit Illumina TruSeq DNA PCR-free kit KAPA Biosystems Hyper Prep PCR system KAPA Biosystems Hyper Prep PCR-free system	WGS	Illumina HiSeq (v4 chemistry)		15. Jones, M. B. et al. Library preparation methodology can influence genomic and functional predictions in human microbiome research. <i>Proc. Natl. Acad. Sci. U S A.</i> 112 , 14024-9 (2015).
9	Sewage water, soil & human stool, biopsy	seems like an inhouse protocol DNA isolation kit (Qiagen, Germany) MagNA pure, Roche Diagnostics, Switzerland	SRS		V1-V3 V1-V5	454 GS FLX+ pyrosequencer		16. Bag, S. et al. An Improved Method for High Quality Metagenomics DNA Extraction from Human and Environmental Samples. <i>Sci. Rep.</i> 6 , 26775 (2016).
10	human stool	seven most commonly used (2017!) kit with modifications (altogether 21 extraction protocols) PSPStool (Invitex) PowerSoil (Mebio) EZNAstool (Omega Bio Tek) Maxwell (Promega) QLAmpStool Minikit (Qiagen) G'Nome (Bio101) FastDNA spin Soil (MP-Biomedicals) MagNAPureIII (Roche)	SRS		WGS	Illumina HiSeq	2x100	17. Costea, P.I. et al. Towards standards for human fecal sample processing in metagenomic studies. <i>Nat. Biotechnol.</i> 35 , 1069-1076 (2017).

11	human stool	TianLong Stool DNA/RNA Extraction Kit (Xi'an TianLong Science and Technology Co., Ltd.) with OR without bead-beating QIAamp DNA Stool mini kit (Qiagen) with OR without bead-beating QIAamp PowerFecal DNA Isolation kit (Qiagen)	SRS	16S Metagenomic Sequencing Library Preparation Illumina	V3-V4	Illumina MiSeq	2x300	QIIME	75. Lim, M. Y., Song, E. J., Kim, S. H., Lee, J. & Nam, Y. D. Comparison of DNA extraction methods for human gut microbial community profiling. <i>Syst. Appl. Microbiol.</i> 41 , 151-157 (2018).
12	human stool	QIAamp DNA Mini Kit (Qiagen) including bead-beating PowerFecal® DNA Isolation Kit (MO BIO)	SRS	PCR amplified	V3-V4	Illumina MiSeq	2x300	QIIME, UCHIME	76. Szopinska, J. W. et al. Reliability of a participant-friendly fecal collection method for microbiome analyses: a step towards large sample size investigation. <i>BMC Microbiol.</i> 18 , 110 (2018).
13	human stool and germ-free mice feces spiked with bacterial or fungal strains	QIAamp DNA Stool Mini Kit (Qiagen), and PureLink™ Microbiome DNA Purification Kit (ThermoFisher) Fecal DNA MiniPrep™ Kit (ZymoResearch) NucleoSpin® DNA Stool Kit (Macherey-Nagel) IHMS protocol Q (QIAamp DNA Stool Kit with bead beating)	qPCR and SRS	16S Metagenomic Sequencing Library Preparation protocol	V3-V4, ITS1F+ITS2	Illumina MiSeq	2x300	QIIME	77. Fiedorová, K. et al. The Impact of DNA Extraction Methods on Stool Bacterial and Fungal Microbiota Community Recovery. <i>Front. Microbiol.</i> 10 , 821 (2019).

14	paleofeces human and dog!!!	the Human Microbiome Project standard protocol using the PowerSoil kit (Qiagen) an aDNA-optimized modified MinElute (Qiagen) protocol for bone extraction following Dabney et al. (2013) Phenol-chloroform + modified MinElute protocol Split modified MinElute protocol HMP protocol + modified MinElute protocol	SRS	NEBNext DNA Library Prep Master Set	WGS	Illumina HiSeq	2x100	QIIME	78. Hagan, R. W., et al. Comparison of extraction methods for recovering ancient microbial DNA from paleofeces. <i>Am. J. Phys. Anthropol.</i> 171 , 275-284 (2020).
15	human stool and 36 bacterial strains represent microbes prevalent in the human body sites	DNeasy Blood and Tissue kit (Qiagen; bacteria) PowerSoil DNA Isolation Kit (Mobio; human)	SRS	TruSeq Nano DNA HT kit (Illumina)	WGS	Illumina MiSeq and HiSeq (bacteria) Illumina NextSeq (human)	2x150	MUSCLE	12. Johnson, J. S. et al. Evaluation of 16S rRNA gene sequencing for species and strain-level microbiome analysis. <i>Nat. Commun.</i> 10 , 5029 (2019).
			LRS	SMRTbell 1.0 Template Prep Kit (Pacific Biosciences)	V1-V9	PacBio RSII (bacteria)			
16	capuchin monkey stools	QIAamp DNA Stool Mini Kit (Qiagen) PowerSoil DNA Isolation kit (Mo Bio)	two-step PCR protocol amplicon library preparation on the Fluidigm Access Array	MyTaq HS Red Mix (Bioline) Accustart II PCR ToughMix (Quantiabio)	V3-V4 V4-V5	Illumina MiSeq		QIIME	31. Mallott, E. K., Malhi, R. S. & Amato, K. R. Assessing the comparability of different DNA extraction and amplification methods in gut microbial community profiling. <i>Access Microbiol.</i> 1 , e000060 (2019).
17	eight different biological specimens, including human stool, saliva, tissues etc. zymo microbial community standard	Quick-DNA fecal/soil microbe kit (Zymo Research) QIAampStool Minikit (Qiagen) MagNA Pure 96 (Roche Diagnostics)	SRS		V4	Illumina NextSeq	2x150	QIIME 2, NG-Tax 0.4	18. Ducarmon, Q. R., Hornung, B. V. H., Geelen, A. R., Kuijper, E. J. & Zwiitink, R. D. Toward Standards in Clinical Microbiota Studies: Comparison of Three DNA Extraction Methods and Two Bioinformatic Pipelines. <i>mSystems</i> 5 , e00547-19 (2020).

18	human stool & mock community (Zymo)	Mag-Bind® Universal Metagenomics Kit (Omega Bio-tek)	SRS	KAPA HyperPrep Kit (Kapa HiFi) with different sample inputs	WGS	Illumina HiSeq 4000	2x250	alpha-, beta-diversity, PCA	79. Peng, Z. et al. Comparative Analysis of Sample Extraction and Library Construction for Shotgun Metagenomics. <i>Bioinform. Biol. Insights</i> 14, 1177932220915459 (2020).
		DNeasy PowerSoil Kit (Qiagen)	shotgun	TruePrep DNA Library Prep Kit V2 (Vazyme Biotech) with different sample			2x350		
19	human stool	QIAamp DNA Stool Mini Kit (Qiagen)	SRS	16S Metagenomic Sequencing Library Preparation Illumina	V3-V4	Illumina MiSeq	2x300	DADA2 plugin [15] within QIIME2	80. Lim, M., et al. Evaluation of fecal DNA extraction protocols for human gut microbiome studies. <i>BMC Microbiol.</i> 20, 212 (2020).
		QIAamp PowerFecal Pro DNA Kit (Qiagen) QIAamp DNA Stool Mini Kit (Qiagen) with additional bead-beating step [5]							
20	simulation and human stool	PureLink Microbiome DNA Purification Kit (Invitrogen)	SRS		WGS	ONT MinION			81. Akili, R. et al. Exploring Semi-Quantitative Metagenomic Studies Using Oxford Nanopore Sequencing: A Computational and Experimental Protocol. <i>Genes</i> 12, 1496 (2021).
		QIamp PowerFecal DNA Kit (Qiagen) ZymoBionics DNA Mini Kit (Zymo Research) Power Microbiome RNA/DNA Isolation Kit (Mo Bio)	LRS		WGS	Illumina NovaSeq SOLiD	2x150		
21	human stool - 200 healthy Japanese	DNeasy Power Soil Kit (Qiagen)	SRS	16S library preparation protocol provided by Illumina	V1-V2	Illumina MiSeq	2x250	QIIME1 & QIIME2	82. Kameoka, S., et al. Benchmark of 16S rRNA gene amplicon sequencing using Japanese gut microbiome data from the V1–V2 and V3–V4 primer sets. <i>BMC Genomics</i> 22, 527 (2021).
			qPCR		V3-V4		2x300	UCLUST, DADA2	
22	dog stool (datasets)	PowerFecal DNA Isolation Kit (MoBio)	SRS		WGS	Illumina HiSeq 2500	2x125	MetaPhlAn2	20. Lewis, S. et al. Comparison of 16S and whole genome dog microbiomes using machine learning. <i>BioData Min.</i> 14, 41 (2021).
					V3-V4	Illumina MiSeq	2x300	DADA2, QIIME2	
23	human stool	Mechanical-Enzymatic Lysis method: Phenol: Chloroform: Isoamyl Alcohol	LRS		WGS	ONT MinION		EPI2ME	83. Sahu, S. et al. Fecal genomic DNA extraction method impacts outcome of MinION based metagenome profile of tuberculosis patients. <i>medRxiv</i> 11, 15.21266154 (2021).

24	human stool	Extrap Soil DNA Kit Plus ver.2 (NIPPON STEEL Eco-Tech Corporation) FastDNA SPIN Kit for Feces (MP Biomedicals) ISOSPIN Fecal DNA Kit (Nippon Gene) MagAttract PowerMicrobio me RNA/DNA EP Kit MORA-EXTRACT kit (Kyokuto Pharmaceutical) QIAmp PowerFecal Pro DNA kit (Qiagen) Quick DNA Fecal/Soil microbe Miniprep Kit (Zymo Research)	SRS	Accel NGS 2S Plus DNA Library Kit (Swift Biosciences) TruSeq DNA PCR-Free Library Prep Kit KAPA HTP Library Preparation Kit (Roche) KAPA HyperPrep Kit PCR-free (Roche) TruSeq Nano DNA Library Prep Kit NEBNext Ultra II DNA Library Prep Kit (New England Biolabs) QIAseq FX DNA Library Kit NEBNext Ultra II FS DNA Library Prep Kit Nextera DNA Flex Library Prep Kit SMARTer ThruPLEX DNA-Seq Kit (Takara Bio)	Illumina NextSeq 500	2x150 bp	Demultiplex: BBmap v38.46 Qc: fastp v0.20.0 for human samples: BMTagger v3.101 Data analysis and visualisation: R v4.0.2 Data handling: dplyr v1.0.2 Visualisation : ggplot2 v3.3.2 Metagenome assembly: MEGAHT v1.2.9. Statistics: Quast v5.0.0 Read annotate: kraken v2.0.8 For OTU taxonomy profiles: mOTUs2 v2.5.1	19. Tourlousse, D. M. et al. Validation and standardization of DNA extraction and library construction methods for metagenomics-based human fecal microbiome measurements. <i>Microbiome</i> 9 , 95 (2021).	
25	human stool	QIAamp DNA Stool Minikit (Qiagen) PSP Spin Stool DNA Plus Kit (Invitex) MoBio PowerSoil DNA Isolation Kit (Mo Bio)	SRS		V1-V3 V1-V2	454 Titanium pyrosequencing 454 FLX pyrosequencing		QIIME	84. Wu, G.D., C. et al. Sampling and pyrosequencing methods for characterizing bacterial communities in the human gut using 16S sequence tags. <i>BMC Microbiol.</i> 10 , 206 (2010).
26	human stool & saliva, conjunctiva, bile, sputum, plaque and water	DNeasy PowerSoil Pro (Qiagen) QIAamp DNA Microbiome Kit (Qiagen) ZymoBIOMIC S DNA Miniprep Kit (Zymo Research)	SRS LRS	MGIEasy Universal DNA Library Prep Set	WGS WGS	DNBSEQ-G400 ONT MinION	2x100		85. Rehner, J., et al. Systematic Cross-biospecimen Evaluation of DNA Extraction Kits for Long- and Short-read Multi-metagenomic Sequencing Studies. <i>Genomics Proteomics Bioinformatics</i> . 20 , 405-417 (2022).

27	10 diverse sample types (human body sites and environment, including stool)	MagAttract PowerSoil DNA Isolation Kit (Qiagen) MagAttract PowerSoil Pro DNA Isolation Kit (Qiagen) Norgen Stool DNA Isolation Kit (Norgen Biotek) MagMAX Microbiome Ultra Nucleic Acid Isolation Kit (Applied Biosystems) NucleoMag Food kit (Macherey-Nagel) ZymoBIOMICS 96 MagBead DNA Kit (Zymo Research)	SRS	NexteraXT	V4 ITS shallow WGS	Miseq Illumina HiSeq	2*150	QIIME2	21. Shaffer, J. P. et al. A comparison of six DNA extraction protocols for 16S, ITS and shotgun metagenomic sequencing of microbial communities. <i>Biotechniques</i> . 73, 34-46 (2022).
b. bird stool									
1	chicken stool	tissue kit genomic DNA template stool mini kit	Real-time PCR						60. Flekna, G., Schneeweiss, W., Smulders, F. J. M., Wagner, M. & Hein, I. Real-time PCR method with statistical analysis to compare the potential of DNA isolation
2	chicken stool	Kit for Feces Kit for Soil MiniPrep DNA Stool Stool Mini Kit DNA Isolation PowerSoil DNA Isolation Genomic DNA miniMAG	Real-time PCR						63. Josefsen, M. H., Andersen, S. C., Christensen, J. & Hoorfar, J. Microbial food safety: Potential of DNA extraction methods for use in diagnostic metagenomics. <i>J. Microbiol. Methods</i> 114, 30-4 (2015).
c. other human tissues									
1	human colonic tissue	DNA/RNA method of protocol method as	microarray			chip			86. Ó Cuív, P. et al. The effects from DNA extraction methods on the evaluation of microbial diversity associated with human colonic tissue. <i>Microb. Ecol.</i> 61, 353-62 (2011).
2	Biopsies from six anatomic regions in	Stool Mini Kit Mini Kit	SRS		V2	pyrosequenc		Greengenes	87. Momozawa, Y. et al. Characterization of bacteria in biopsies of colon and stools by high throughput sequencing of the V2
3	pediatric bronchoalveolar lavage & microbial community	trimeethylammo (saline) Tissue Kit MoBio PowerSoil DNA Isolation Kit (MoBio)	SRS		V8-V9	pyrosequenc	~ 250	QIIME	88. Willner, D. et al. Comparison of DNA extraction methods for microbial community profiling with an application to pediatric bronchoalveolar lavage samples. <i>PLoS One</i> . 7, e34605 (2012).

4	human lung tissue	Mini kit and QIAamp orm: Isoamyl orm: Isoamyl step and Phenol: Chloroform: Isoamyl alcohol step and Bead-beating and QIAamp DNA Mini kit (QIAGEN)	SRS	Index	V3-V4 ITS	MiSeq	2x300	QIIME	89. Pérez-Brocá, V. et al. Optimized DNA extraction and purification method for characterization of bacterial and fungal communities in lung tissue samples. <i>Sci. Rep.</i> 10 , 17377 (2020).
5	human biopsies (Colon)	Microbiome Technique developed by the group	SRS qPCR SRS		V3-V4 WGS	NovaSeq Illumina NovaSeq	2x250 2x150	Uparse Mothur	90. Bruggeling, C. E. et al. Optimized bacterial DNA isolation method for microbiome analysis of human tissues. <i>Microbiologyopen</i> . 10 , e1191 (2021).
6	human urine samples	BiOstic and Tissue PowerSoil UltraClean Maxwell RSC Purefood GMO and Authentication (Promega)	SRS	amplified	V4	MiSeq	2x150	DADA2	91. Karstens, L. et al. Benchmarking DNA isolation kits used in analyses of the urinary microbiome. <i>Sci. Rep.</i> 11 , 6186 (2021).
7	human upper airways (nose, saliva, pharynx) & mock	Mini Kit PowerViral Microbiome Stool DNA ZymoBIOMICS DNA Kit (Zymo Research)	SRS	DNA sample	WGS	MiSeq	2x300	orX2	92. Mancabelli, L. et al. Guideline for the analysis of the microbial communities of the human upper airways. <i>J. Oral. Microbiol.</i> 14 , 2103282 (2022).
8	human urine	Bacteremia kit Microbiome kit (Beckman- and Tissue kit MagNA Pure Compact Kit (Roche) - all the optional lysis steps were performed	SRS	amplified	V1-V3 V3-V4 V4-V5 V6-V8	MiSeq	2x300	SILVA	93. Vendrell, J. A. et al. Determination of the Optimal Bacterial DNA Extraction Method to Explore the Urinary Microbiota. <i>Int. J. Mol. Sci.</i> 23 , 1336 (2022).
9	human urine	Pre-Cell Lysis: Mechanical With Enzymatic Lysis (lytic enzyme solution (Qiagen) and	LRS		mNGS	MinION			94. Zhang, L. et al. Comparison Analysis of Different DNA Extraction Methods on Suitability for Long-Read Metagenomic Nanopore Sequencing. <i>Front. Cell. Infect. Microbiol.</i> 12 , 919903 (2022).

d. other tissues from animals									
1	cow and sheep rumen	extraction method [32] Spin Kit with Spin Kit with matrix (Bio-based method chloroform, chloroform, chloroform, chloroform Stool DNA Kit, Stool DNA Kit, DNA Stool beating plus MiniPrep	pyrosequencing qPCR		rRNA	Titanium		QIIME	95. Henderson, G. et al. Effect of DNA extraction methods and sampling techniques on the apparent structure of cow and sheep rumen microbial communities. <i>PLoS One</i> . 8 , e74787 (2013).
2	Chicken caeca and crop	procedure	SRS	DNA Library Metagenomic	WGS V3-V4	MiSeq	2x150 bp 2x150 bp	MG-RAST Silva SSU	96. Durazzi, F. et al. Comparison between 16S rRNA and shotgun sequencing data for the taxonomic characterization of the gut microbiota. <i>Sci Rep</i> . 11 , 2020
d. plant									
1	Arabidopsis thaliana, corn and soybean	Plant PCR Kit DNA Isolation kit (Qiagen) Fungal/Bacteria	SRS	amplified	V4 V3-V4 V5-V7	MiSeq	2x250	QIIME2	97. Giangacomio, C., Mohseni, M., Kovar, L. & Wallace, J. G. Comparing DNA Extraction and 16S rRNA Gene Amplification Methods for Plant Associated
e. bacteria strains and cultures									
1	Genomic DNA	chloroform	point detection specified, most		V2, V6			MEGALIGN	30. Chakravorty, S., Helb, D., Burday, M., Connell, N. & Alland, D. de Boer, R. et al. Improved
2	single	LC DNA III	rtPCR		(not)				98. Bowers, R.M., et al. Impact of library preparation protocols and template quantity on the metagenomic reconstruction of a mock microbial community. <i>BMC</i>
3	mock microbial community	genomic	SRS	Nextera XT Mondrian Illumina's MALBAC	WGS?	HiSeq 2000	2*150		99. Gand, M., et al. Comparison of 6 DNA extraction methods for isolation of high yield of high molecular weight DNA suitable for shotgun metagenomics Nanopore sequencing to detect bacteria. <i>BMC Genomics</i> 24 , 438 (2023).
4	bacterial cocktail mixes	HMW (Macherey-S DNA PowerFecal Pro kit (Claremont on Moss et al.:	LRS		WGS	ONT			
f. water and other nonliving source									
1	Lake Taihu's water (China)	Water DNA Kit (OMEGA,	SRS	V3-V4 regions were amplified from omic Sequencing Library Preparation with some minor modifications.	V3 V4 V6	Ion Torrent PGM		Statistical significance was tested in the differential taxa among V3, V4, and V6 using SPSS18.0 software	100. Zhang, J. Evaluation of different 16S rRNA gene V regions for exploring bacterial diversity in a eutrophic freshwater lake. <i>Sci. Total Environ.</i> 618 , 1254-1267 (2018).
2	Lake Baikal's water (Russia)	treatment,	SRS	Ultra II DNA	V2-V3 V3-V4	MiSeq			22. Bukin, Y. et al. The effect of 16S rRNA region choice on bacterial community metabarcoding results. <i>Sci. Data</i> 6 , 190007 (2019).
3	Equatorial Pacific & North Atlantic - water	MOBIO PowerSoil DNA Isolation Kit (Mo Bio)	SRS	V regions were amplified	V4 V6	Illumina MiSeq	2x250 & 2x150		23. Kerrigan, Z., Kirkpatrick, J. B. & DHondt, S. Influence of 16S rRNA Hypervariable Region on Estimates of Bacterial Diversity and Community Composition in Seawater and Marine Sediment

4	Activated Sludge, Biofilm, and Anaerobic Digestate	FastDNA Spin kit for Soil (MP Biomedicals) MicrobiomeT M Purification Kit (ThermoFisher Scientific) FavorPrep Soil DNA Isolation Mini Kit (Favorgen Biotech)	capillary electrophoresis						101. Florczyk, M., Cydzik-Kwiatkowska, A., Ziembinska-Buczynska, A. & Ciesielski A. Comparison of Three DNA Extraction Kits for Assessment of Bacterial Diversity in Activated Sludge, Biofilm, and Anaerobic Digestate. <i>Appl. Sci.</i> 12 , 9797 (2022).
5	Soil	not specified	SRS+LRS	TruSeq Nano SMRTbell Template Prep Kit	400 bp WGS, 20 000 bp	HiSeq 2000 PacBio Sequel	not specified	rk, BLAST, PB assembly:MetaFlye, IL: SPAdes	102. Xu, G. et al. Combined assembly of long and short sequencing reads improve the efficiency of exploring the soil metagenome. <i>BMC Genomics</i> 23 , 27 (2022).
6	On-site Wastewater	DNA Isolation	SRS	Ultra™ II	WGS V4	HiSeq MiSeq	2x150 2x150		103. de Vries, J. et al. Comparative Analysis of Metagenomic

10.2. Supplementary Data 2

Supplementary Data 2. Applied wet lab kits and bioinformatic pipelines			
a. DNA extraction kits used in our study. HMW: High-molecular weight DNA.			
Company	Kit	Bead-beating	HMW
Qiagen	QIAamp Fast DNA Stool Mini Kit	No	No
Macherey-Nagel	NucleoSpin DNA Stool Mini kit	Yes	No
Invitrogen	PureLink™ Microbiome DNA Purification Kit	Yes	No
Zymo Research	Quick-DNA™ HMW MagBead Kit	No	Yes
b. Library preparation kits utilized in this project. mWGS: metagenomic whole-genome sequencing.			
Company	Kit	16S rRNA region	mWGS
Zymo Research	Quick-16S NGS Library Prep Kit	Partial (V1-V2)	
Zymo Research	Quick-16S NGS Library Prep Kit	Partial (V3-V4)	
PerkinElmer	NEXTFLEX® 16S V1-V3 Amplicon-Seq Kit	Partial (V1-V3)	
Illumina	DNA Prep Kit		Yes
Oxford Nanopore Technologies	16S Barcoding Kit	Full-length (V1-V9)	
Pacific Biosciences	Full-Length 16S Library Preparation Using SMRTbell Express Template Prep Kit 2.0	Full-length (V1-V9)	
c. The applied bioinformatic tools			
Program	Partial 16S rRNA gene	Full-length 16S rRNA gene	Shotgun
DADA2	Yes		
Emu	Yes	Yes	
EPI2ME		Yes	
Kaiju			Yes
sourmash			Yes
minitax	Yes	Yes	Yes

10.3. Supplementary Data 3

Supplementary Data 3.								
a. The table summarizes the details of DNA purification from canine sample (the main tested dog), the obtained yield and quality of DNA								
Sample	Initial amount (g)	Elution (ul)	gDNA cc. (ng/ul)		Length of the highest detected peak (bp)	Length of the longest detected DNA fragment (bp)	Average of the highest detected peak (bp)	Average of the longest detected DNA fragment (bp)
Qiagen	0,2	100	Q1	2,25	-	-	-	-
			Q2	1,89	13178	15062	8359	9361
			Q3	1,96	-	-	-	-
			Q4	1,16	3540	3660	-	-
Invitrogen	0,2	100	I1	14,4	14293	26235	-	-
			I2	12	15162	26503	14469	26347
			I3	20,8	12685	24945	-	-
			I4	11,2	15734	27705	-	-
Macherey-Nagel	0,2	100	MN1	20	6023	18277	-	-
			MN2	23	8235	23199	7255	21325
			MN3	22	7867	21278	-	-
			MN4	31	6896	22547	-	-
Zymo Research	0,1	50	Z1	11,2	23514	35526	-	-
			Z2	13,2	16785	34863	18286	33867
			Z3	13,9	20340	33513	-	-
			Z4	11,3	12505	31568	-	-
b. DNA purification from canine sample (control dogs for validation) and								
Sample	Initial amount (g)	Elution (ul)	gDNA cc. (ng/ul)					
Qiagen	0,2	100	Female 1	1,85				
			Male 1	4,96				
			Female 2	41,2				
			Female 3	26,8				
			Male 2	43,4				
			Male 3	19				

c. Statistical Comparisons of DNA Fragment Lengths Across Different DNA Isolation Kits. This table shows the differences in DNA fragment lengths between pairs of DNA isolation kits, along with their corresponding p-values and significance levels. The "difference" column represents the mean difference in fragment lengths between kit pairs. The "p-value" column indicates the statistical significance of each comparison, with values less than 0.05 considered significant. The "Significance" column denotes whether the difference is statistically significant ("yes") or not ("no"). In all cases where significance is indicated, the 95% confidence intervals exclude zero, confirming the statistical significance of the differences observed. Conversely, no significant difference was detected between Macherey-Nagel and Invitrogen kits, as evidenced by a confidence interval that includes zero.

DNA isolation kit pairs	difference	p-value	Significance
Macherey-Nagel - Qiagen	-11,964.25	=0.005	yes
Macherey-Nagel - ZymoResearch	12,542.25	=0.001	yes
Macherey-Nagel - Invitrogen	-5,021.75	0.154	no
Zymo-Research - Invitrogen	7,520.50	= 0.025	yes
Zymo-Research - Qiagen	24,506.50	< 0.001	yes
Qiagen - Invitrogen	-16,986.00	< 0.001	yes

10.4. Supplementary Data 4

Supplementary Data 4. Sample-wise Shannon and Simpson indices.

sample	host	DNA extraction method	Target (V-region)	platform	program	database	sample name	Shannon	Simpson
TotI Illumina V1-V2 minitax NCBI.genome.collection.11	dog (TotI)	Invitrogen (I)	V1-V2	Illumina	minitax	NCBI.genome.collection	11	4.008	0.947
TotI Illumina V1-V2 minitax NCBI.genome.collection.12	dog (TotI)	Invitrogen (I)	V1-V2	Illumina	minitax	NCBI.genome.collection	12	4.904	0.951
TotI Illumina V1-V2 minitax NCBI.genome.collection.13	dog (TotI)	Invitrogen (I)	V1-V2	Illumina	minitax	NCBI.genome.collection	13	4.070	0.948
TotI Illumina V1-V2 minitax NCBI.genome.collection.14	dog (TotI)	Invitrogen (I)	V1-V2	Illumina	minitax	NCBI.genome.collection	14	3.887	0.945
TotI Illumina V1-V2 minitax NCBI.genome.collection.MN1	dog (TotI)	Machery-nagel (MN)	V1-V2	Illumina	minitax	NCBI.genome.collection	MN1	4.005	0.952
TotI Illumina V1-V2 minitax NCBI.genome.collection.MN2	dog (TotI)	Machery-nagel (MN)	V1-V2	Illumina	minitax	NCBI.genome.collection	MN2	4.260	0.952
TotI Illumina V1-V2 minitax NCBI.genome.collection.MN3	dog (TotI)	Machery-nagel (MN)	V1-V2	Illumina	minitax	NCBI.genome.collection	MN3	4.045	0.950
TotI Illumina V1-V2 minitax NCBI.genome.collection.MN4	dog (TotI)	Machery-nagel (MN)	V1-V2	Illumina	minitax	NCBI.genome.collection	MN4	3.962	0.953
TotI Illumina V1-V2 minitax NCBI.genome.collection.21	dog (TotI)	Zymo (Z)	V1-V2	Illumina	minitax	NCBI.genome.collection	21	4.113	0.949
TotI Illumina V1-V2 minitax NCBI.genome.collection.22	dog (TotI)	Zymo (Z)	V1-V2	Illumina	minitax	NCBI.genome.collection	22	3.880	0.950
TotI Illumina V1-V2 minitax NCBI.genome.collection.23	dog (TotI)	Zymo (Z)	V1-V2	Illumina	minitax	NCBI.genome.collection	23	3.898	0.950
TotI Illumina V1-V2 minitax NCBI.genome.collection.24	dog (TotI)	Zymo (Z)	V1-V2	Illumina	minitax	NCBI.genome.collection	24	3.884	0.947
TotI Illumina V1-V3 minitax NCBI.genome.collection.11	dog (TotI)	Invitrogen (I)	V1-V3	Illumina	minitax	NCBI.genome.collection	11	4.210	0.911
TotI Illumina V1-V3 minitax NCBI.genome.collection.12	dog (TotI)	Invitrogen (I)	V1-V3	Illumina	minitax	NCBI.genome.collection	12	3.670	0.912
TotI Illumina V1-V3 minitax NCBI.genome.collection.13	dog (TotI)	Invitrogen (I)	V1-V3	Illumina	minitax	NCBI.genome.collection	13	4.473	0.966
TotI Illumina V1-V3 minitax NCBI.genome.collection.14	dog (TotI)	Invitrogen (I)	V1-V3	Illumina	minitax	NCBI.genome.collection	14	4.513	0.966
TotI Illumina V1-V3 minitax NCBI.genome.collection.MN1	dog (TotI)	Machery-nagel (MN)	V1-V3	Illumina	minitax	NCBI.genome.collection	MN1	4.113	0.931
TotI Illumina V1-V3 minitax NCBI.genome.collection.MN2	dog (TotI)	Machery-nagel (MN)	V1-V3	Illumina	minitax	NCBI.genome.collection	MN2	3.855	0.925
TotI Illumina V1-V3 minitax NCBI.genome.collection.MN3	dog (TotI)	Machery-nagel (MN)	V1-V3	Illumina	minitax	NCBI.genome.collection	MN3	3.456	0.923
TotI Illumina V1-V3 minitax NCBI.genome.collection.MN4	dog (TotI)	Machery-nagel (MN)	V1-V3	Illumina	minitax	NCBI.genome.collection	MN4	4.214	0.934
TotI Illumina V1-V3 minitax NCBI.genome.collection.O1	dog (TotI)	Qiagen (Q)	V1-V3	Illumina	minitax	NCBI.genome.collection	O1	3.924	0.912
TotI Illumina V1-V3 minitax NCBI.genome.collection.O2	dog (TotI)	Qiagen (Q)	V1-V3	Illumina	minitax	NCBI.genome.collection	O2	4.032	0.948
TotI Illumina V1-V3 minitax NCBI.genome.collection.O3	dog (TotI)	Qiagen (Q)	V1-V3	Illumina	minitax	NCBI.genome.collection	O3	3.712	0.950
TotI Illumina V1-V3 minitax NCBI.genome.collection.O4	dog (TotI)	Qiagen (Q)	V1-V3	Illumina	minitax	NCBI.genome.collection	O4	4.034	0.951
TotI Illumina V1-V3 minitax NCBI.genome.collection.Z1	dog (TotI)	Zymo (Z)	V1-V3	Illumina	minitax	NCBI.genome.collection	Z1	4.159	0.950
TotI Illumina V1-V3 minitax NCBI.genome.collection.Z2	dog (TotI)	Zymo (Z)	V1-V3	Illumina	minitax	NCBI.genome.collection	Z2	3.756	0.919
TotI Illumina V1-V3 minitax NCBI.genome.collection.Z3	dog (TotI)	Zymo (Z)	V1-V3	Illumina	minitax	NCBI.genome.collection	Z3	3.827	0.918
TotI Illumina V1-V3 minitax NCBI.genome.collection.Z4	dog (TotI)	Zymo (Z)	V1-V3	Illumina	minitax	NCBI.genome.collection	Z4	3.143	0.844
TotI Illumina V1-V4 minitax NCBI.genome.collection.11	dog (TotI)	Invitrogen (I)	V1-V4	Illumina	minitax	NCBI.genome.collection	11	3.505	0.909
TotI Illumina V1-V4 minitax NCBI.genome.collection.12	dog (TotI)	Invitrogen (I)	V1-V4	Illumina	minitax	NCBI.genome.collection	12	3.548	0.921
TotI Illumina V1-V4 minitax NCBI.genome.collection.13	dog (TotI)	Invitrogen (I)	V1-V4	Illumina	minitax	NCBI.genome.collection	13	3.480	0.913
TotI Illumina V1-V4 minitax NCBI.genome.collection.14	dog (TotI)	Invitrogen (I)	V1-V4	Illumina	minitax	NCBI.genome.collection	14	3.504	0.912
TotI Illumina V1-V4 minitax NCBI.genome.collection.MN1	dog (TotI)	Machery-nagel (MN)	V1-V4	Illumina	minitax	NCBI.genome.collection	MN1	3.504	0.913
TotI Illumina V1-V4 minitax NCBI.genome.collection.MN2	dog (TotI)	Machery-nagel (MN)	V1-V4	Illumina	minitax	NCBI.genome.collection	MN2	3.343	0.911
TotI Illumina V1-V4 minitax NCBI.genome.collection.MN3	dog (TotI)	Machery-nagel (MN)	V1-V4	Illumina	minitax	NCBI.genome.collection	MN3	3.570	0.920
TotI Illumina V1-V4 minitax NCBI.genome.collection.MN4	dog (TotI)	Machery-nagel (MN)	V1-V4	Illumina	minitax	NCBI.genome.collection	MN4	3.608	0.922
TotI Illumina V1-V4 minitax NCBI.genome.collection.Z1	dog (TotI)	Zymo (Z)	V1-V4	Illumina	minitax	NCBI.genome.collection	Z1	3.576	0.925
TotI Illumina V1-V4 minitax NCBI.genome.collection.Z2	dog (TotI)	Zymo (Z)	V1-V4	Illumina	minitax	NCBI.genome.collection	Z2	3.403	0.915
TotI Illumina V1-V4 minitax NCBI.genome.collection.Z3	dog (TotI)	Zymo (Z)	V1-V4	Illumina	minitax	NCBI.genome.collection	Z3	3.522	0.923
TotI Illumina V1-V4 minitax NCBI.genome.collection.Z4	dog (TotI)	Zymo (Z)	V1-V4	Illumina	minitax	NCBI.genome.collection	Z4	3.857	0.947
TotI Illumina mWGS minitax NCBI.genome.collection.11	dog (TotI)	mWGS (M)	V1-V4	Illumina	minitax	NCBI.genome.collection	11	4.195	0.948
TotI Illumina mWGS minitax NCBI.genome.collection.12	dog (TotI)	mWGS (M)	V1-V4	Illumina	minitax	NCBI.genome.collection	12	4.804	0.951
TotI Illumina mWGS minitax NCBI.genome.collection.13	dog (TotI)	mWGS (M)	V1-V4	Illumina	minitax	NCBI.genome.collection	13	3.260	0.940
TotI Illumina mWGS minitax NCBI.genome.collection.14	dog (TotI)	mWGS (M)	V1-V4	Illumina	minitax	NCBI.genome.collection	14	4.203	0.947
TotI Illumina mWGS minitax NCBI.genome.collection.MN1	dog (TotI)	Machery-nagel (MN)	V1-V4	Illumina	minitax	NCBI.genome.collection	MN1	4.005	0.952
TotI Illumina mWGS minitax NCBI.genome.collection.MN2	dog (TotI)	Machery-nagel (MN)	V1-V4	Illumina	minitax	NCBI.genome.collection	MN2	4.414	0.957
TotI Illumina mWGS minitax NCBI.genome.collection.MN3	dog (TotI)	Machery-nagel (MN)	V1-V4	Illumina	minitax	NCBI.genome.collection	MN3	4.439	0.953
TotI Illumina mWGS minitax NCBI.genome.collection.MN4	dog (TotI)	Machery-nagel (MN)	V1-V4	Illumina	minitax	NCBI.genome.collection	MN4	3.993	0.950
TotI Illumina mWGS minitax NCBI.genome.collection.O1	dog (TotI)	Qiagen (Q)	V1-V4	Illumina	minitax	NCBI.genome.collection	O1	4.414	0.958
TotI Illumina mWGS minitax NCBI.genome.collection.O2	dog (TotI)	Qiagen (Q)	V1-V4	Illumina	minitax	NCBI.genome.collection	O2	4.375	0.956
TotI Illumina mWGS minitax NCBI.genome.collection.O3	dog (TotI)	Qiagen (Q)	V1-V4	Illumina	minitax	NCBI.genome.collection	O3	3.394	0.943
TotI Illumina mWGS minitax NCBI.genome.collection.O4	dog (TotI)	Qiagen (Q)	V1-V4	Illumina	minitax	NCBI.genome.collection	O4	4.472	0.957
TotI Illumina mWGS minitax NCBI.genome.collection.Z1	dog (TotI)	Zymo (Z)	V1-V4	Illumina	minitax	NCBI.genome.collection	Z1	4.228	0.954
TotI Illumina mWGS minitax NCBI.genome.collection.Z2	dog (TotI)	Zymo (Z)	V1-V4	Illumina	minitax	NCBI.genome.collection	Z2	4.273	0.948
TotI Illumina mWGS minitax NCBI.genome.collection.Z3	dog (TotI)	Zymo (Z)	V1-V4	Illumina	minitax	NCBI.genome.collection	Z3	4.272	0.943
TotI Illumina mWGS minitax NCBI.genome.collection.Z4	dog (TotI)	Zymo (Z)	V1-V4	Illumina	minitax	NCBI.genome.collection	Z4	4.005	0.940
TotI Illumina mWGS minitax NCBI.genome.collection.11	dog (TotI)	Invitrogen (I)	V1-V3	Illumina	minitax	NCBI.genome.collection	11	4.018	0.938
TotI Illumina mWGS minitax NCBI.genome.collection.12	dog (TotI)	Invitrogen (I)	V1-V3	Illumina	minitax	NCBI.genome.collection	12	4.160	0.942
TotI Illumina mWGS minitax NCBI.genome.collection.13	dog (TotI)	Invitrogen (I)	V1-V3	Illumina	minitax	NCBI.genome.collection	13	4.050	0.945
TotI Illumina mWGS minitax NCBI.genome.collection.14	dog (TotI)	Invitrogen (I)	V1-V3	Illumina	minitax	NCBI.genome.collection	14	4.011	0.938
TotI Illumina mWGS minitax NCBI.genome.collection.MN1	dog (TotI)	Machery-nagel (MN)	V1-V3	Illumina	minitax	NCBI.genome.collection	MN1	4.218	0.948
TotI Illumina mWGS minitax NCBI.genome.collection.MN2	dog (TotI)	Machery-nagel (MN)	V1-V3	Illumina	minitax	NCBI.genome.collection	MN2	4.219	0.949
TotI Illumina mWGS minitax NCBI.genome.collection.MN3	dog (TotI)	Machery-nagel (MN)	V1-V3	Illumina	minitax	NCBI.genome.collection	MN3	4.230	0.950
TotI Illumina mWGS minitax NCBI.genome.collection.MN4	dog (TotI)	Machery-nagel (MN)	V1-V3	Illumina	minitax	NCBI.genome.collection	MN4	4.196	0.946
TotI Illumina mWGS minitax NCBI.genome.collection.O1	dog (TotI)	Qiagen (Q)	V1-V3	Illumina	minitax	NCBI.genome.collection	O1	3.169	0.837
TotI Illumina mWGS minitax NCBI.genome.collection.O2	dog (TotI)	Qiagen (Q)	V1-V3	Illumina	minitax	NCBI.genome.collection	O2	2.960	0.800
TotI Illumina mWGS minitax NCBI.genome.collection.O3	dog (TotI)	Qiagen (Q)	V1-V3	Illumina	minitax	NCBI.genome.collection	O3	4.007	0.949
TotI Illumina mWGS minitax NCBI.genome.collection.O4	dog (TotI)	Qiagen (Q)	V1-V3	Illumina	minitax	NCBI.genome.collection	O4	3.245	0.843
TotI Illumina mWGS minitax NCBI.genome.collection.Z1	dog (TotI)	Zymo (Z)	V1-V3	Illumina	minitax	NCBI.genome.collection	Z1	3.956	0.925
TotI Illumina mWGS minitax NCBI.genome.collection.Z2	dog (TotI)	Zymo (Z)	V1-V3	Illumina	minitax	NCBI.genome.collection	Z2	3.043	0.819
TotI Illumina mWGS minitax NCBI.genome.collection.Z3	dog (TotI)	Zymo (Z)	V1-V3	Illumina	minitax	NCBI.genome.collection	Z3	4.045	0.937
TotI Illumina mWGS minitax NCBI.genome.collection.Z4	dog (TotI)	Zymo (Z)	V1-V3	Illumina	minitax	NCBI.genome.collection	Z4	4.104	0.938
TotI ONT V1-V9 minitax NCBI.genome.collection.11	dog (TotI)	Invitrogen (I)	V1-V9	MinION (ONT)	minitax	NCBI.genome.collection	11	2.765	0.881
TotI ONT V1-V9 minitax NCBI.genome.collection.12	dog (TotI)	Invitrogen (I)	V1-V9	MinION (ONT)	minitax	NCBI.genome.collection	12	2.488	0.893
TotI ONT V1-V9 minitax NCBI.genome.collection.13	dog (TotI)	Invitrogen (I)	V1-V9	MinION (ONT)	minitax	NCBI.genome.collection	13	2.910	0.897
TotI ONT V1-V9 minitax NCBI.genome.collection.14	dog (TotI)	Invitrogen (I)	V1-V9	MinION (ONT)	minitax	NCBI.genome.collection	14	2.865	0.889
TotI ONT V1-V9 minitax NCBI.genome.collection.MN1	dog (TotI)	Machery-nagel (MN)	V1-V9	MinION (ONT)	minitax	NCBI.genome.collection	MN1	3.013	0.907
TotI ONT V1-V9 minitax NCBI.genome.collection.MN2	dog (TotI)	Machery-nagel (MN)	V1-V9	MinION (ONT)	minitax	NCBI.genome.collection	MN2	3.045	0.908
TotI ONT V1-V9 minitax NCBI.genome.collection.MN3	dog (TotI)	Machery-nagel (MN)	V1-V9	MinION (ONT)	minitax	NCBI.genome.collection	MN3	3.987	0.903
TotI ONT V1-V9 minitax NCBI.genome.collection.MN4	dog (TotI)	Machery-nagel (MN)	V1-V9	MinION (ONT)	minitax	NCBI.genome.collection	MN4	3.046	0.912
TotI ONT V1-V9 minitax NCBI.genome.collection.O1	dog (TotI)	Qiagen (Q)	V1-V9	MinION (ONT)	minitax	NCBI.genome.collection	O1	2.389	0.872
TotI ONT V1-V9 minitax NCBI.genome.collection.O2	dog (TotI)	Qiagen (Q)	V1-V9	MinION (ONT)	minitax	NCBI.genome.collection	O2	2.210	0.875
TotI ONT V1-V9 minitax NCBI.genome.collection.O3	dog (TotI)	Qiagen (Q)	V1-V9	MinION (ONT)	minitax	NCBI.genome.collection	O3	3.735	0.915
TotI ONT V1-V9 minitax NCBI.genome.collection.O4	dog (TotI)	Qiagen (Q)	V1-V9	MinION (ONT)	minitax	NCBI.genome.collection	O4	2.287	0.886
TotI ONT V1-V9 minitax NCBI.genome.collection.Z1	dog (TotI)	Zymo (Z)	V1-V9	MinION (ONT)	minitax	NCBI.genome.collection	Z1	2.504	0.897
TotI ONT V1-V9 minitax NCBI.genome.collection.Z2	dog (TotI)	Zymo (Z)	V1-V9	MinION (ONT)	minitax	NCBI.genome.collection	Z2	2.999	0.914
TotI ONT V1-V9 minitax NCBI.genome.collection.Z3	dog (TotI)	Zymo (Z)	V1-V9	MinION (ONT)	minitax	NCBI.genome.collection	Z3	3.319	0.926
TotI ONT V1-V9 minitax NCBI.genome.collection.Z4	dog (TotI)	Zymo (Z)	V1-V9	MinION (ONT)	minitax	NCBI.genome.collection	Z4	3.014	0.910
TotI PacBio V1-V9 minitax NCBI.genome.collection.11	dog (TotI)	Invitrogen (I)	V1-V9	PacBio	minitax	NCBI.genome.collection	11	2.409	0.799
TotI PacBio V1-V9 minitax NCBI.genome.collection.12	dog (TotI)	Invitrogen (I)	V1-V9	PacBio	minitax	NCBI.genome.collection	12	2.581	0.835
TotI PacBio V1-V9 minitax NCBI.genome.collection.13	dog (TotI)	Invitrogen (I)	V1-V9	PacBio	minitax	NCBI.genome.collection	13	2.949	0.880
TotI PacBio V1-V9 minitax NCBI.genome.collection.14	dog (TotI)	Invitrogen (I)	V1-V9	PacBio	minitax	NCBI.genome.collection	14	2.677	0.869
TotI PacBio V1-V9 minitax NCBI.genome.collection.MN1	dog (TotI)	Machery-nagel (MN)	V1-V9	PacBio	minitax	NCBI.genome.collection	MN1	4.697	0.885
TotI PacBio V1-V9 minitax NCBI.genome.collection.MN2	dog (TotI)	Machery-nagel (MN)	V1-V9	PacBio	minitax	NCBI.genome.collection	MN2	2.739	0.884
TotI PacBio V1-V9 minitax NCBI.genome.collection.MN3	dog (TotI)	Machery-nagel (MN)	V1-V9	PacBio	minitax	NCBI.genome.collection	MN3	2.739	0.885
TotI PacBio V1-V9 minitax NCBI.genome.collection.MN4	dog (TotI)	Machery-nagel (MN)	V1-V9	PacBio	minitax	NCBI.genome.collection	MN4	2.727	0.885
TotI PacBio V1-V9 minitax NCBI.genome.collection.Z1	dog (TotI)	Zymo (Z)	V1-V9	PacBio	minitax	NCBI.genome.collection	Z1	2.846	0.890
TotI PacBio V1-V9 minitax NCBI.genome.collection.Z2	dog (TotI)	Zymo (Z)	V1-V9	PacBio	minitax	NCBI.genome.collection	Z2	3.692	0.928
TotI PacBio V1-V9 minitax NCBI.genome.collection.Z3	dog (TotI)	Zymo (Z)	V1-V9	PacBio	minitax	NCBI.genome.collection	Z3	2.883	0.919
TotI PacBio V1-V9 minitax NCBI.genome.collection.Z4	dog (TotI)	Zymo (Z)	V1-V9	PacBio	minitax	NCBI.genome.collection	Z4	3.001	0.921

sample	source	DNA extraction method	Target (V region)	platform	program	database	sample name	channel	Genomes
Zymo0600 Illumina V1-V2 minitax NCBI genome.collection 1	Zymo D6800	Invitrogen (I)	V1-V2	Illumina	minitax	NCBI genome collection	I1	3,772	0.952
Zymo0600 Illumina V1-V2 minitax NCBI genome.collection 12	Zymo D6800	Invitrogen (I)	V1-V2	Illumina	minitax	NCBI genome collection	I2	3,766	0.947
Zymo0600 Illumina V1-V2 minitax NCBI genome.collection 13	Zymo D6800	Invitrogen (I)	V1-V2	Illumina	minitax	NCBI genome collection	I3	3,762	0.941
Zymo0600 Illumina V1-V2 minitax NCBI genome.collection MM1	Zymo D6800	Machery-Nagel (MN)	V1-V2	Illumina	minitax	NCBI genome collection	MM1	4,012	0.953
Zymo0600 Illumina V1-V2 minitax NCBI genome.collection MM2	Zymo D6800	Machery-Nagel (MN)	V1-V2	Illumina	minitax	NCBI genome collection	MM2	3,990	0.955
Zymo0600 Illumina V1-V2 minitax NCBI genome.collection MM3	Zymo D6800	Machery-Nagel (MN)	V1-V2	Illumina	minitax	NCBI genome collection	MM3	3,959	0.953
Zymo0600 Illumina V1-V2 minitax NCBI genome.collection Q1	Zymo D6800	Qiagen (Q)	V1-V2	Illumina	minitax	NCBI genome collection	Q1	3,781	0.951
Zymo0600 Illumina V1-V2 minitax NCBI genome.collection Q2	Zymo D6800	Qiagen (Q)	V1-V2	Illumina	minitax	NCBI genome collection	Q2	3,934	0.954
Zymo0600 Illumina V1-V2 minitax NCBI genome.collection Q3	Zymo D6800	Qiagen (Q)	V1-V2	Illumina	minitax	NCBI genome collection	Q3	3,742	0.932
Zymo0600 ONT V1-V9 minitax NCBI genome.collection 11	Zymo D6800	Invitrogen (I)	V1-V9	MinION (ONT)	minitax	NCBI genome collection	I1	2,502	0.930
Zymo0600 ONT V1-V9 minitax NCBI genome.collection 12	Zymo D6800	Invitrogen (I)	V1-V9	MinION (ONT)	minitax	NCBI genome collection	I2	2,884	0.951
Zymo0600 ONT V1-V9 minitax NCBI genome.collection 13	Zymo D6800	Invitrogen (I)	V1-V9	MinION (ONT)	minitax	NCBI genome collection	I3	2,889	0.918
Zymo0600 ONT V1-V9 minitax NCBI genome.collection 14	Zymo D6800	Invitrogen (I)	V1-V9	MinION (ONT)	minitax	NCBI genome collection	I4	2,965	0.920
Zymo0600 ONT V1-V9 minitax NCBI genome.collection 15	Zymo D6800	Invitrogen (I)	V1-V9	MinION (ONT)	minitax	NCBI genome collection	I5	2,876	0.919
Zymo0600 ONT V1-V9 minitax NCBI genome.collection MM1	Zymo D6800	Machery-Nagel (MN)	V1-V9	MinION (ONT)	minitax	NCBI genome collection	MM1	3,500	0.897
Zymo0600 ONT V1-V9 minitax NCBI genome.collection MM2	Zymo D6800	Machery-Nagel (MN)	V1-V9	MinION (ONT)	minitax	NCBI genome collection	MM2	2,881	0.899
Zymo0600 ONT V1-V9 minitax NCBI genome.collection MM3	Zymo D6800	Machery-Nagel (MN)	V1-V9	MinION (ONT)	minitax	NCBI genome collection	MM3	2,888	0.885
Zymo0600 ONT V1-V9 minitax NCBI genome.collection Q1	Zymo D6800	Machery-Nagel (MN)	V1-V9	MinION (ONT)	minitax	NCBI genome collection	Q1	2,714	0.896
Zymo0600 ONT V1-V9 minitax NCBI genome.collection Q2	Zymo D6800	Qiagen (Q)	V1-V9	MinION (ONT)	minitax	NCBI genome collection	Q2	2,578	0.882
Zymo0600 ONT V1-V9 minitax NCBI genome.collection Q3	Zymo D6800	Qiagen (Q)	V1-V9	MinION (ONT)	minitax	NCBI genome collection	Q3	2,518	0.883
Zymo0600 ONT V1-V9 minitax NCBI genome.collection Q4	Zymo D6800	Qiagen (Q)	V1-V9	MinION (ONT)	minitax	NCBI genome collection	Q4	2,535	0.857
Zymo0600 ONT V1-V9 minitax NCBI genome.collection Q5	Zymo D6800	Qiagen (Q)	V1-V9	MinION (ONT)	minitax	NCBI genome collection	Q5	2,441	0.887
Zymo0600 ONT V1-V9 minitax NCBI genome.collection C5	Zymo D6800	Qiagen (Q)	V1-V9	MinION (ONT)	minitax	NCBI genome collection	C5	2,600	0.884
Zymo0600 ONT V1-V9 minitax NCBI genome.collection Z1	Zymo D6800	Zymo (Z)	V1-V9	MinION (ONT)	minitax	NCBI genome collection	Z1	2,998	0.910
Zymo0600 ONT V1-V9 minitax NCBI genome.collection Z2	Zymo D6800	Zymo (Z)	V1-V9	MinION (ONT)	minitax	NCBI genome collection	Z2	2,838	0.921
Zymo0600 ONT V1-V9 minitax NCBI genome.collection Z3	Zymo D6800	Zymo (Z)	V1-V9	MinION (ONT)	minitax	NCBI genome collection	Z3	2,714	0.875
Zymo0600 ONT V1-V9 minitax NCBI genome.collection Z4	Zymo D6800	Zymo (Z)	V1-V9	MinION (ONT)	minitax	NCBI genome collection	Z4	2,726	0.880
Zymo0600 ONT V1-V9 minitax NCBI genome.collection Z5	Zymo D6800	Zymo (Z)	V1-V9	MinION (ONT)	minitax	NCBI genome collection	Z5	2,725	0.880
Zymo0600 Illumina V1-V2 minitax NCBI genome.collection Z1	Zymo D6800	Zymo (Z)	V1-V2	Illumina	minitax	NCBI genome collection	Z1	3,921	0.954
Zymo0600 Illumina V1-V2 minitax NCBI genome.collection Z2	Zymo D6800	Zymo (Z)	V1-V2	Illumina	minitax	NCBI genome collection	Z2	3,864	0.947
Zymo0600 Illumina V1-V2 minitax NCBI genome.collection Z3	Zymo D6800	Zymo (Z)	V1-V2	Illumina	minitax	NCBI genome collection	Z3	3,997	0.957
Zymo06331 ONT V2-V9 minitax NCBI genome.collection 11	Zymo D6331	Invitrogen (I)	V1-V9	MinION (ONT)	minitax	NCBI genome collection	I1	2,959	0.913
Zymo06331 ONT V2-V9 minitax NCBI genome.collection 12	Zymo D6331	Invitrogen (I)	V1-V9	MinION (ONT)	minitax	NCBI genome collection	I2	2,999	0.823
Zymo06331 ONT V2-V9 minitax NCBI genome.collection 13	Zymo D6331	Invitrogen (I)	V1-V9	MinION (ONT)	minitax	NCBI genome collection	I3	2,961	0.786
Zymo06331 ONT V2-V9 minitax NCBI genome.collection MM1	Zymo D6331	Machery-Nagel (MN)	V1-V9	MinION (ONT)	minitax	NCBI genome collection	MM1	2,275	0.877
Zymo06331 ONT V2-V9 minitax NCBI genome.collection MM2	Zymo D6331	Machery-Nagel (MN)	V1-V9	MinION (ONT)	minitax	NCBI genome collection	MM2	2,238	0.859
Zymo06331 ONT V2-V9 minitax NCBI genome.collection MM3	Zymo D6331	Machery-Nagel (MN)	V1-V9	MinION (ONT)	minitax	NCBI genome collection	MM3	2,407	0.829
Zymo06331 ONT V2-V9 minitax NCBI genome.collection Q1	Zymo D6331	Qiagen (Q)	V1-V9	MinION (ONT)	minitax	NCBI genome collection	Q1	2,235	0.881
Zymo06331 ONT V2-V9 minitax NCBI genome.collection Q2	Zymo D6331	Qiagen (Q)	V1-V9	MinION (ONT)	minitax	NCBI genome collection	Q2	2,235	0.839
Zymo06331 ONT V2-V9 minitax NCBI genome.collection Q3	Zymo D6331	Qiagen (Q)	V1-V9	MinION (ONT)	minitax	NCBI genome collection	Q3	2,199	0.810
Zymo06331 ONT V2-V9 minitax NCBI genome.collection Z1	Zymo D6331	Zymo (Z)	V1-V9	MinION (ONT)	minitax	NCBI genome collection	Z1	2,236	0.821
Zymo06331 ONT V2-V9 minitax NCBI genome.collection Z2	Zymo D6331	Zymo (Z)	V1-V9	MinION (ONT)	minitax	NCBI genome collection	Z2	2,241	0.875
Zymo06331 ONT V2-V9 minitax NCBI genome.collection Z3	Zymo D6331	Zymo (Z)	V1-V9	MinION (ONT)	minitax	NCBI genome collection	Z3	2,230	0.811

10.5. Supplementary Data 5

Supplementary Data S5a. The statistical data of all sequencing performed during the project. a. Data obtained from a single sample of a 13-year-old dog.

[illegible]

10.6. Supplementary Data 6

Supplementary Data 6. Comparing our data with others' results. This dataset summarizes all available results related to the dog gut microbiome of this phylum (Nauzet-Guén, Coelho et al., Söder et al., Thomson et al., Li et al., Xue et al., Lemly-Thompson et al., Li et al., Xue et al.) and genus (Yu and Kim, Thomson et al., Li et al., Xue et al.), compared with our own data. Coelho and colleagues conducted Humana 1055 sequencing, while all other publications used the analysis of the 16S rRNA V3-V4 region. All of these publications, inclusive of our own dataset, cannot on a single aspect, the canine gut microbiome primarily consists of five phyla: Firmicutes, Bacteroidetes, Proteobacteria, Actinobacteria, and Fusobacteria, with the first contributing close to half the detected genes.

[illegible]

10.7. Supplementary Data 7

Supplementary Data 7. a. DNA purification from MCS. The table shows the yeald (quantity measured by Qubit) and the interval of the largest peaks detected by TapeStation.

Samples	Initial volume (μl)	Elution (μl)	gDNA cc. (ng/μl)		Largest peak	
					Highest average (bp)	Largest average (bp)
Qiagen	75	50	1	0,44	26271	>60000
			2	0,292		
			3	0,252		
			4	0,552		
			5	0,286		
Invitrogen	75	50	1	0,286	15170	>60000
			2	0,334		
			3	0,342		
			4	0,222		
			5	0,252		
Macherey-Nagel	75	50	1	9,86	5419	54740
			2	11		
			3	9,92		
			4	10,9		
			5	9,14		
Zymo Research	75	50	1	40,4	9849	>60000
			2	42,6		
			3	34,6		
			4	37		
			5	40		

b. DNA purification from GMS. The table shows the yeald (quantity measured by Qubit) and the interval of the largest peaks detected by TapeStation.

Samples	Initial volume (µl)	Elution (µl)	gDNA cc. (ng/µl)	Largest peak		
				Highest average (bp)	Largest average (bp)	
Qiagen	75	50	1	1,17	6895	>60000
			2	1,29		
			3	0,978		
Invitrogen	75	50	1	0,06	4746	58897
			2	0,07		
			3	0,075		
Macherey-Nagel	75	50	1	11,8	5340	55538
			2	12,7		
			3	0,038 *		
Zymo Research	75	50	1	4,36	5539	>60000
			2	4,14		
			3	3,76		

10.8. Supplementary Data 8

Supplementary Data 8. Sequencing platform-specific CIGAR scoring schemes for calculating CIGAR scores in minitax's alignment processing.

Platform	Match Score	Mismatch Score	Insertion Score	Deletion Score	Gap Opening Penalty	Gap Extension Penalty	Description
Illumina	2	-4	-3	-3	-4	-2	Optimized for high-accuracy, short reads. Higher penalties for mismatches and indels to reflect the platform's low error rate.
ONT	1	-3	-2	-2	-2	-1	Adjusted for longer reads with higher error rates. More lenient penalties to accommodate frequent indels and mismatches.
PacBio	2	-3	-3	-3	-3	-2	Balanced settings for long, high-fidelity reads (e.g., HiFi mode). Moderate penalties for indels to support accurate alignment in repetitive regions.

10.9. Supplementary Data 9

Supplementary Data 9. PCR amplification conditions for Zymo Research Quick-16S library preparation (V1-V2 and V3-V4 amplicon sequencing) – PCR I

PCR steps	Temperature	Time	Cycle number
Initial denaturation	95°C	10 min	1
Denaturation	95°C	30 sec	20 (dog) or 12 (MCS)
Annealing	55°C	30 sec	
Extension	72°C	3 min	

10.10. Supplementary Data 10

Supplementary Data 10. Index barcode (ZA7-ZA5) combinations used for library preparation using the Zymo Research Quick-16S NGS Library Prep Kit.

Abbreviations: I: Invitrogen; MN: Macherey-Nagel; Z: Zymo Research

a, Barcodes used for canine stool samples				
Sample	Index Barcode ZA7	Sequence ZA7	Index Barcode ZA5	Sequence ZA5
I1	ZA701	ACCTGGAT	ZA501	TTCTAGAC
I2	ZA702	GTGCCATA	ZA502	CCGATCTT
I3	ZA703	TGAATCCG	ZA503	TAAGATCC
I4	ZA704	CATGATGC	ZA504	AGGTCATT
MN1	ZA705	AATGTCCT	ZA501	TTCTAGAC
MN2	ZA706	ATAGGCTC	ZA502	CCGATCTT
MN3	ZA701	ACCTGGAT	ZA503	TAAGATCC
MN4	ZA702	GTGCCATA	ZA504	AGGTCATT
Z1	ZA703	TGAATCCG	ZA501	TTCTAGAC
Z2	ZA704	CATGATGC	ZA502	CCGATCTT
Z3	ZA705	AATGTCCT	ZA503	TAAGATCC
Z4	ZA706	ATAGGCTC	ZA504	AGGTCATT
b, Barcode sequences applied for MCS samples				
Sample	Index Barcode ZA7	Sequence ZA7	Index Barcode ZA5	Sequence ZA5
I1	ZA707	TTGCGGAG	ZA505	GTGTGTCA
I2	ZA705	AATGTCCT	ZA504	AGGTCATT
I3	ZA701	ACCTGGAT	ZA501	TTCTAGAC
Q1	ZA705	AATGTCCT	ZA503	TAAGATCC
Q2	ZA702	GTGCCATA	ZA502	CCGATCTT
Q3	ZA706	ATAGGCTC	ZA504	AGGTCATT
MN1	ZA703	TGAATCCG	ZA501	TTCTAGAC
MN3	ZA708	GCCTTCCA	ZA508	CCACAGGT
MN4	ZA704	CATGATGC	ZA502	CCGATCTT
Z1	ZA706	ATAGGCTC	ZA505	GTGTGTCA
Z2	ZA709	GTCAGTCT	ZA508	CCACAGGT
Z3	ZA705	AATGTCCT	ZA505	GTGTGTCA

10.11. Supplementary Data 11

Supplementary Data 11. PCR amplification conditions for Zymo Research Quick-16S library preparation (V1-V2 and V3-V4 amplicon sequencing) – PCR II (barcoded PCR)			
PCR steps	Temperature	Time	Cycle number
Initial denaturation	95°C	10 min	1
Denaturation	95°C	30 sec	5
Annealing	55°C	30 sec	
Extension	72°C	3 min	

10.12. Supplementary Data 12

Supplementary Data 12. Summary table for showing the amount of DNA used for PerkinElmer library preparation, the amount of libraries used for sequencing. The table also shows the barcode ID-s as well as the index sequences.

Sample	Genomic DNA cc (ng/ul)	Library cc (ng/ul)	Barcode	primer index
Q1	22,2	3,24	11	AAGCGTACGTCC
Q2	26,4	4,02	13	TCGGGAAGGTCC
Q3	25,5	8,56	37	TCGGGAAGGTCC
Q4	36	0,16	38	GAGGCATCGGCC
MN1	2,5	29,2	14	TTATCAGTCCTT
MN2	2,2	34,8	15	GTCATCGCGTCC
MN3	2,27	42,8	39	AATAATTGGTCC
MN4	1,6	40,6	40	GTCGTCAACCGG
I1	3,47	34,8	16	CCGTCTCTCCGG
I2	4,1	5,22	18	ACGCTCTTCCGG
I3	2,4	41,8	41	AAATCTCAGGCC
I4	4,46	45,4	42	GTGCGCGGCCGG
Z1	4,46	25,8	13	AAGCGTACGTCC
Z2	3,78	13,8	14	TTATCAGTCCTT
Z3	3,59	34,6	16	CCGTCTCTCCGG
Z4	4,42	4,1	18	ACGCTCTTCCGG

10.13. Supplementary Data 13

Supplementary Data 13. PCR Conditions for PerkinElmer NEXTFLEX® 16S V1-V3 Amplicon-Seq Kit for Illumina – PCR I

PCR steps	Temperature	Time	Cycle number
Initial denaturation	98°C	4 min	1
Denaturation	98°C	30 sec	8
Annealing	60°C	30 sec	
Extension	72°C	30sec	
Final extension	72°C	4 min	1

10.14. Supplementary Data 14

Supplementary Data 14. PCR Conditions for PerkinElmer NEXTFLEX® 16S V1-V3 Amplicon-Seq Kit for Illumina – PCR II

PCR steps	Temperature	Time	Cycle number
Initial denaturation	98°C	4 min	1
Denaturation	98°C	30 sec	varying according to the amount of the initial amount of DNA *
Annealing	60°C	30 sec	
Extension	72°C	30sec	
Final extension	72°C	4 min	1

10.15. Supplementary Data 15

Supplementary Data 15. Summary table for showing the barcode IDs and sequences for ONT 16S library preparation. IDs used for library preparation

a, from the main subject of the study			
Replicate #1		Replicate #2	
Sample	Barcode	Sample	Barcode
I1	bc09	I1	bc09
I2	bc10	I2	bc02
I3	bc11	I3	bc03
I4	bc12	I4	bc04
MN1	bc05	MN1	bc05
MN2	bc06	MN2	bc10
MN3	bc07	MN3	bc11
MN4	bc08	MN4	bc12
Z1	bc01	Z1	bc01
Z2	bc02	Z2	bc06
Z3	bc03	Z3	bc07
Z4	bc04	Z4	bc08
Q1	bc01		
Q2	bc02		
Q3	bc03		
Q4	bc04		
b, from the six additional dogs			
Sample	Barcode		
Female 1	bc21		
Male 1	bc17		
Female 2	bc1		
Female 3	bc5		
Male 2	bc9		
Male 3	bc13		
c, form MCS			
Sample	Barcode		
I1	bc7		
I2	bc8		
I3	bc9		
I4	bc6		
I5	bc7		
MN1	bc4		
MN2	bc5		
MN3	bc6		
MN4	bc3		
MN5	bc4		
Z1	bc1		
Z2	bc2		
Z3	bc3		
Z4	bc1		
Z5	bc2		
Q1	bc10		
Q2	bc11		
Q3	bc12		
Q4	bc8		
Q5	bc9		
d, barcode sequences			
ONT's Barcode ID	ID used in the study	Barcode sequence	
16S01	bc01	AAGAAAGTTGTCGGTGCTTTGTG	
16S02	bc02	TCGATTCCGTTTGTAGTCGTCTGT	
16S03	bc03	GAGTCTTGTTGCCAGTTACCAGG	
16S04	bc04	TTCGGATTCTATCGTGTTCCTTA	
16S05	bc05	CTGTCCAGGGTTTGTGTAACCTT	
16S06	bc06	TTCTCGCAAAGGCAGAAAGTAGTC	
16S07	bc07	GTGTTACCGTGGAATGAATCCTT	
16S08	bc08	TTCAGGGAACAAACCAAGTTACGT	
16S09	bc09	AACTAGGCACAGCGAGTCTTGTT	
16S10	bc10	AAGCGTTGAAACCTTTGTCCTCTC	
16S11	bc11	GTTTCATCTATCGGAGGAATGGA	
16S12	bc12	CAGGTAGAAAGAAGCAGAATCGGA	
16S13	bc13	AGAACGACTTCCATACTCGTGTGA	
16S17	bc17	ACCCTCCAGGAAAGTACCTCTGAT	
16S21	bc21	GAGCCTCTCATTGTCGGTTCTCTA	

10.16. Supplementary Data 16

Supplementary Data 16. PCR Conditions for preparation of ONT V1-V9 libraries

PCR step	Temperature	Time	No. of cycles
Initial denaturation	95 °C	1 min	1
Denaturation	95 °C	20 secs	25
Annealing	55 °C	30 secs	25
Extension	65 °C	2 mins	25
Final extension	65 °C	5 mins	1
Hold	4 °C	∞	

10.17. Supplementary Data 17

Supplementary Data 17. Summary table of primers used for PacBio 16S library preparation.

	Forward primer	Forward primer sequence	Reverse primer	Reverse primer sequence
I1	16S_Fw_1007	TCTGTATCTCTATGTG	16S_Rev_1056	ATGTGCGTGTGTGTCT
I2	16S_Fw_1008	ACAGTCGAGCGCTGCG	16S_Rev_1056	ATGTGCGTGTGTGTCT
I3	16S_Fw_1012	ACACTAGATCGCGTGT	16S_Rev_1056	ATGTGCGTGTGTGTCT
I4	16S_Fw_1015	CGCATGACACGTGTGT	16S_Rev_1056	ATGTGCGTGTGTGTCT
MN1	16S_Fw_1020	CACGACACGACGATGT	16S_Rev_1056	ATGTGCGTGTGTGTCT
MN2	16S_Fw_1022	CACTCACGTGTGATAT	16S_Rev_1056	ATGTGCGTGTGTGTCT
MN3	16S_Fw_1024	CATGTAGAGCAGAGAG	16S_Rev_1056	ATGTGCGTGTGTGTCT
MN4	16S_Fw_1005	CACTCGACTCTCGCGT	16S_Rev_1057	CTCTCAGACGCTCGTC
Z1	16S_Fw_1007	TCTGTATCTCTATGTG	16S_Rev_1057	CTCTCAGACGCTCGTC
Z2	16S_Fw_1008	ACAGTCGAGCGCTGCG	16S_Rev_1057	CTCTCAGACGCTCGTC
Z3	16S_Fw_1012	ACACTAGATCGCGTGT	16S_Rev_1057	CTCTCAGACGCTCGTC
Z4	16S_Fw_1015	CGCATGACACGTGTGT	16S_Rev_1057	CTCTCAGACGCTCGTC

10.18. Supplementary Data 18

Supplementary Data 18. PCR conditions for preparation of PacBio V1-V9 libraries

PCR steps	Temperature	Time	Cycle number
Initial denaturation	95°C	3 min	1
Denaturation	95°C	30 sec	25
Annealing	57°C	30 sec	
Extension	72°C	60sec	

10.19. Supplementary Data 19

Supplementary Data 19. Summary table of primers from Illumina DNA Prep Kit used for WGS library preparation.

Sample	Index Barcode i7	Barcode Sequence i7	Index Barcode i5	Barcode Sequence i5
Q1	H705	GGACTCCT	H505	GTAAGGAG
Q2	H705	GGACTCCT	H503	TATCCTCT
Q3	H706	TAGGCATG	H505	GTAAGGAG
Q4	H714	GCTCATGA	H517	GCGTAAGA
I1	H711	AAGAGGCA	H503	TATCCTCT
I2	H714	GCTCATGA	H505	GTAAGGAG
I3	H714	GCTCATGA	H506	ACTGCATA
I4	H711	AAGAGGCA	H505	GTAAGGAG
MN1	H707	CTCTCTAC	H506	ACTGCATA
MN2	H710	CGAGGCTG	H517	GCGTAAGA
MN3	H706	TAGGCATG	H506	ACTGCATA
MN4	H705	GGACTCCT	H517	GCGTAAGA
Z1	H711	AAGAGGCA	H506	ACTGCATA
Z2	H714	GCTCATGA	H503	TATCCTCT
Z3	H706	TAGGCATG	H517	GCGTAAGA
Z4	H707	CTCTCTAC	H503	TATCCTCT

10.20. Supplementary Data 20

Supplementary Data 20. PCR Conditions for preparation of Illumina WGS libraries

PCR steps	Temperature	Time	Cycle number
Pre-heating	68°C	3 min	1
Initial denaturation	98°C	3 min	1
Denaturation	98°C	45 sec	5
Annealing	62°C	30 sec	
Extension	68°C	2 min	
Final extension	68°C	1 min	1

10.21. Supplementary Data 21

Supplementary Data 21.

Genus	EPI2ME value	NCBI value
<i>Clostridium</i>	0.0137855341364841	0.0135811686783089
<i>Kineothrix</i>	0.00103065605367019	N/A
<i>Collinsella</i>	0.0341797461460624	0.0359682063346239
<i>Terrisporobacter</i>	0.00232185350281726	N/A
<i>Paeniclostridium</i>	0.00827780884220112	N/A
<i>Romboutsia</i>	0.00209802830920327	N/A
<i>Blautia</i>	0.921588848109184	1.0123767434625
<i>Enterocloster</i>	0.00636597606432633	0.0260434289128637
<i>Campylobacter</i>	0.00104069475650225	N/A
<i>Erysipelatoclostridium</i>	0.00866333190917401	N/A
<i>Faecalicatena</i>	0.0128830781652849	0.0202894786939917
<i>Mediterraneibacter</i>	0.00822035492654341	0.438629625707591
<i>Phocaeicola</i>	0.0169983976390598	0.0183188779189124
<i>Megamonas</i>	0.00813308189787587	0.0088669435189919
<i>Roseburia</i>	0.00140255975128613	0.00950004122302111
<i>Cellulosilyticum</i>	0.00263681138534562	N/A
<i>Bacteroides</i>	0.00183680964869668	N/A
<i>Lacrimispora</i>	0.0167242255407207	0.0271792333362727
<i>Fusobacterium</i>	0.0399184899712927	0.0409118529388183
<i>Peptacetobacter</i>	0.0045492181048344	1.73183878618779
<i>Catenibacterium</i>	N/A	0.0462073780062417
<i>Eisenbergiella</i>	N/A	0.0336467399273302
<i>Faecalimonas</i>	N/A	0.296900378156428
<i>Lachnoclostridium</i>	N/A	0.0100645523848315
<i>Intestinibacter</i>	N/A	0.0140135036154961
<i>Tyzzerella</i>	N/A	0.0400242743675637
<i>Coprococcus</i>	N/A	0.0135261646295627