

Exploring Quantitative Metagenomics Studies using Oxford Nanopore Sequencing: A Computational and Experimental Protocol

Rohia ALILI

Unité de recherche mixte 1269 - NutriOmics <https://orcid.org/0000-0002-0158-4250>

Eugeni BELDA (✉ e.belda@integrative-phenomics.com)

Integrative Phenomics <https://orcid.org/0000-0003-4307-5072>

Karine CLEMENT

UMRS 1269 - NutriOmics

Phuong Le

Sorbonne Université,

Edi PRIFTI

IRD : UMMISCO

Jean-Daniel ZUCKER

IRD : UMMISCO

Thierry WIRTH

Institut de Systématique Evolution Biodiversité: Institut de Systematique Evolution Biodiversite

Methodology

Keywords: quantitative metagenomics, microbiome, obesity, gut microbiota, microbial DNA extraction, sequencing, Simulation, Oxford Nanopore Technologies, MinION

Posted Date: December 22nd, 2020

DOI: <https://doi.org/10.21203/rs.3.rs-131495/v1>

License: © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Exploring quantitative metagenomics studies using Oxford Nanopore sequencing : a computational and experimental protocol

Rohia Alili^{1,2,5*}, Eugeni Belda^{3*#}, Phuong Le¹, Thierry Wirth^{4,5}, Jean-Daniel Zucker^{1,6}, Edi Prifti^{1,6}, Karine Clément^{1,2}

*These authors contributed equally

Institutional addresses

¹ Sorbonne Université, INSERM, Nutrition and obesities; systemic approaches (NutriOmics), Paris, France.

² Assistance Publique Hôpitaux de Paris, Pitié-Salpêtrière Hospital, Nutrition department, CRNH Ile de France, Paris France.

³ Integrative Phenomics, Paris, France.

⁴ Institut de Systématique, Evolution, Biodiversité (ISYEB), Muséum national d'Histoire naturelle, CNRS, Sorbonne Université, Université des Antilles, EPHE, Paris 75005, France.

⁵ PSL University, EPHE, Paris 75014, France.

⁶ IRD, Sorbonne Université, UMMISCO, Unité de Modélisation Mathématique et Informatique des Systèmes Complexes, F-93143, Bondy, France.

Email addresses :

Rohia ALILI : rohia.alili@aphp.fr

Eugeni BELDA : e.belda@integrative-phenomics.com

Karine CLEMENT : karine.clement@inserm.fr

1 Phuong LE : phuongleee@gmail.com

2 Edi Prifti : edi.prifti@gmail.com

3 Thierry WIRTH : thierry.wirth@mnhn.fr

4 Jean-Daniel ZUCKER : jdzucker@gmail.com

5 **# Corresponding authors**

6 #Eugeni Belda, Integrative Phenomics, 8 Rue des Pirogues de Bercy, 75012 Paris

7 Paris, France.

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

Abstract

1
2 Background : The gut microbiome plays a major role in chronic diseases, several of
3 which are characterized by an altered diversity and composition of bacterial
4 communities. Large-scale sequencing projects allowed the characterization of these
5 microbial community perturbations. However, a gap remains in how these discoveries
6 can be translated into clinical applications. To facilitate routine implementation of
7 microbiome profiling in clinical settings, portable, real-time, and low-cost sequencing
8 technologies are needed.

9 Results : Here, we propose a computational and experimental protocol for whole
10 genome quantitative metagenomics studies of the human gut microbiome with Oxford
11 Nanopore sequencing technology (ONT). We developed a bioinformatic pipeline to
12 process ONT sequences based on the evaluation of different alignment parameters in
13 the estimation of microbial diversity and composition. We also optimized stool
14 collection and DNA extraction methods to maximize read length, a critical parameter
15 for the sequence alignment and classification. Our analytical pipeline was evaluated
16 using simulations of metagenomic communities to reflect naturally occurring
17 compositional variations. We then validated our experimental and analytical pipeline
18 with stool samples from a bariatric surgery cohort sequenced with ONT and Illumina,
19 revealing comparable diversity and microbial composition profiles. These results were
20 compared to those previously obtained with SOLiD sequencing, where differences
21 were observed, possibly explained by variations in library preparation steps. Finally,
22 we found that sequences obtained with ONT allowed assembly of complete genomes
23 for disease-related species.

24 Conclusion : This protocol can be implemented in the clinical or individual setting,
25 bringing rapid personalized whole genome profiling of target microbiome species.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25

Keywords: quantitative metagenomics, microbiome, obesity, gut microbiota, microbial DNA extraction, sequencing, Simulation, Oxford Nanopore Technologies, MinION.

1 **Background**

2 In recent years, there has been a burst in knowledge related to gut microbiota
3 screening in chronic diseases. The increasing access to high throughput sequencing
4 has led to the discovery of alterations in the composition of intestinal microbiota in
5 many human disorders, including metabolic diseases. Currently, there is a challenge
6 to discover reproducible gut microbial signatures for diseases in order to develop
7 diagnostic and prognostic tools.

8 In the field of metabolic diseases, microbial diversity and richness is generally
9 representative of microbiome and host health, as exemplified in previous studies such
10 as MetaHit[1], HMP[2], Metacardis[2] and others covering severe obesity, bariatric
11 surgery[3], diabetes, NAFLD/NASH[4], and cirrhosis[5][6]. In mild[7] and severe
12 obesity[8] for instance, we previously showed that reduced microbial richness linked
13 to altered composition was found in 40% to 75% of the subjects and was associated
14 with a more deleterious host phenotype. Even with these established signatures in
15 metabolic diseases, the gut microbiome varies greatly in composition and abundance
16 from one individual to another.

17 In addition to biological variation, gut microbiome quantification is subject to technical
18 variation along the experimental process starting from sample collection and extending
19 to DNA extraction, library preparation and sequencing, and bioinformatics analytical
20 protocols[9]. This is observed in the literature where results are frequently
21 irreproducible[10][11][12][13]. Reports highlight the need for technical standardization
22 [14]. Even though progress has been made with the work of different international
23 consortia[15][16] to standardize protocols, there is still a need to develop methods for
24 fast-track and affordable microbiome screening in clinical settings. Among critical steps

1 is DNA extraction prior to sequencing. Costea et al.[16] for example reported biases in
2 microbial composition and diversity varying with different DNA extraction protocols.
3 Extraction protocols with or without bead-beating increase the representation of gram-
4 positive bacteria as is also the case for different DNA extraction kits. For example, the
5 richness is higher and reads are longer with the Qiagen compared to Magnapure
6 kits[17]. The library preparation has also an impact on relative abundances of bacterial
7 features and functional annotation[10]. Finally, the bioinformatic pipelines can yield
8 consequent variability in microbial ecosystem description[18].

9 Presently, most microbiome research studies are carried out using 16S ribosomal RNA
10 genes or whole genome shotgun sequencing (WGS), the latter requiring quite a long
11 time for data processing and result producing. Moreover, not all medical and research
12 centers are able to set up high-end shotgun sequencing platforms due to multiple
13 constraints. As opposed to previously existing technology, Oxford Nanopore
14 Technology (ONT) proposes real-time sequence data generation with fewer resources
15 and a small benchtop footprint. Due to the recent development of this technology, there
16 is a need to define standardized wet-lab and bioinformatics protocols for ONT in
17 metagenomics. Importantly, annotation and quantification of ONT-generated long
18 reads is properly needed as current bioinformatics and metagenomics pipeline are
19 adapted to short reads. Here, we have explored methods to optimize ONT for
20 microbiome analyses and propose a fine-tuned protocol, including wet-lab preparation
21 (i.e. sample collection, DNA extraction, and library preparation) and data processing
22 and analysis. In particular, we have set up a customized analytical pipeline to estimate
23 microbial composition and diversity as well as to classify ONT reads using current
24 bacterial gene catalogs along with functional profiling. This protocol is open-access,
25 allowing it to be replicated, and implemented by medical or research centers.

1 **Materials and methods**

2 **Study design**

3 To determine the optimal parameters to maximize sequence mapping and prepare for
4 wet-lab experiments, we first optimized our bioinformatics pipeline using a simulation
5 framework, a classical approach in computer science[19]. We simulated sequence
6 data based on a set of known bacterial genomes as well as their abundance
7 distribution based on existing data, which we varied in terms of composition, richness
8 and sequencing depth. The simulator we used took into account the particularity of
9 ONT sequences and biases. Next, we built and tuned our bioinformatics pipeline
10 considering the best hyperparameters to minimize the difference between the
11 estimated quantified features (abundance, richness) with those used to parameterize
12 the simulation (**Figure 1A**).

13 In parallel to the bioinformatic pipeline optimization, we conducted wet lab experiments
14 to establish an optimized protocol for stool collection and DNA extraction, DNA
15 fragmentation, and DNA end-repair on microbiome composition. Using the results of
16 these experiments in coordination with our bioinformatics pipeline, we proposed an
17 optimized protocol for DNA quality and increased ONT sequencing read length (**Figure**
18 **1B**). Finally, we validated our protocol and pipeline using human stool samples
19 sequenced in parallel with ONT and Illumina technology and previously sequenced
20 with SOLiD technology (**Figure 1C**).

21 **Simulated ONT microbiome data**

22 We designed the data simulation process to maximize representation of real human
23 gut microbial ecosystems. We used 506 reference genomes that were used with

1 metagenomic assemblies to build the IGC[1][20][21] human gut gene catalog. We
2 simulated ten samples (M1:M10) whose abundances were calculated using a Pareto
3 distribution, which was estimated using real metagenomic profiles[8] based on the
4 abundance of metagenomic species[20] computed on the same IGC catalog[21].
5 These empirical distributions were used to define generation probabilities of the 506
6 reference genomes. We included two important elements for the quantification of
7 microbial ecosystems into the simulation: richness (number of present species) and
8 the sequencing depth (i.e. the number of reads generated by the sequencing). We
9 simulated the variation in richness in microbial communities that could be observed
10 ranging from 50 to 450 species (R50:R450), as well as the sequencing depth ranging
11 from 1x to 5x the complete coverage of the genomes present. In total, 250 samples
12 were simulated using the CAMISIM software (option: Nanosim tool with default error
13 parameters for the E.coli example)[22] (**Figure 1A**).

14 **The bioinformatics workflow**

15 The proposed bioinformatic workflow for quantitative metagenomic (QM) analyses
16 from ONT shotgun sequencing starts with fast5 files generated by the MiniKNOWtm
17 software. The Fast5 file format is derived from the HDF5 standard and has a
18 hierarchical structure in which the metadata associated with individual reads (i.e.
19 position in the flow cell, start time of the sequencing of the DNA fragment, time elapsed
20 to pass through the flow cell pore and when it ends) as well as the events (i.e.
21 aggregated bulk current measurements for each DNA fragment) are stored and then
22 decoded by base calling algorithms into fastq format. The first step of the workflow
23 consists of base calling and demultiplexing the fast5 files into fastq files. Here we used
24 Albacore and Guppy ONT base callers from the community site available to ONT

1 customers[23]. Since the start of our work, ONT had discontinued development on
2 Albacore in favor of the more performant Guppy[24].

3 Secondly, the *sequence_summary.txt* files, generated during the base calling step, are
4 used to generate different QC visualizations along with information embedded in the
5 fast5 files. This allows evaluating the quality of the sequencing run in terms of number
6 of active channels in the flow cell, the distribution of active channels through time, the
7 yield in terms of reads of the run and the read length distribution. The third step of the
8 workflow consists in the taxonomic binning of ONT reads using two different reference
9 resources. The first one uses Centrifuge [25] for the taxonomic binning of individual
10 reads using a comprehensive reference database of more than 8000 reference
11 genomes from prokaryotes and viruses (including the reference human genome).
12 Thus, this step allows excluding human sequence reads. To remove spurious
13 taxonomic assignments, we additionally mapped against the corresponding reference
14 genome from centrifuge database using Minimap2 with *map-ont* option optimized for
15 ONT reads [26]. Based on simulation experiments results, only sequences with a
16 minimum mapQ score of 5 were retained for subsequent analyses. A species
17 abundance table was generated by summing the counts of each NCBI taxid from the
18 filtered Centrifuge results. This abundance table was combined with the experiment
19 metadata information and a reference taxonomic table reconstructed from Centrifuge
20 NCBI taxids using the R package *taxize*[27] into a S4 *phyloseq* object. This object was
21 used to carry out standard microbial ecology analyses (rarefaction, alpha-diversity,
22 beta-diversity, and differential abundance analysis) with *Phyloseq* or *VeganR*
23 packages [28].

24 A complementary approach consisted of quantifying the abundance of microbial
25 genes. For this, ONT reads were aligned against the IGC [21] gene catalog using

1 Minimap2 with *map-ont* option [26]. The alignment of long ONT reads over short or
2 fragmented IGC genes provided a preliminary structural annotation of the reads in the
3 sense that a long read could cover a genomic region harboring more than one gene,
4 but also provided situations of multiple genes mapping in overlapping regions of a read.
5 Overlaps were filtered out using *GenomicRanges* and *plyrRangesR* packages[29][30]
6 allowing to quantify genes with the highest mapQ score and sequence identity across
7 each alignment region. The raw gene abundance table was reconstructed by counting
8 the number of times each gene was mapped by read.

9 The final step consisted in quantifying functional features (i.e. taxonomic binning).
10 Functional features were KEGG orthology groups (KO groups), quantified from the
11 gene abundance table, using available reference annotation from the IGC catalog [31].
12 For taxonomic results produced from Centrifuge quantification, we retrieved the KO
13 content of KEGG genomes from the KEGG API [32] for which species-level pan-
14 genomes were reconstructed for all species-level bins based on NCBI taxonomy and
15 matched with genomic sequences in the Centrifuge database. Based on this matching,
16 the abundance of KO groups from Centrifuge results were computed as the sum of the
17 abundances of the species containing these KO groups. The pan-genome strategy fits
18 with the compressed nature of Centrifuge genomes at the species level, followed to
19 reduce the size of the indexes and improve the overall performance of the classification
20 process [25].

21 **Study participants for wet-lab experiments**

22 Stool samples used for wet-lab protocol optimization were collected from healthy
23 volunteers (n=15; men=8, BMI 18-25 kg/m²) from the European "Metacardis" cohort[2].
24 For the comparison between sequencing technologies, we used 33 samples from the

1 Microbaria study, where the gut microbiome of subjects with severe obesity was
2 characterized before and after bariatric surgery[8].

3 **Sample collection and bacterial DNA extraction for pre-analytic protocol** 4 **experiments**

5 Fresh stools were collected with two different methods: 1) a dry spoon tube
6 (SARSTED), which requires storage at -80°C and 2) a tube containing DNA/RNA
7 stabilizing solution, which can be kept at room temperature, -20°C, or -80°C depending
8 on the duration of storage. For the latter collection method, we tested four available
9 commercial kits including "DNA/RNA Shield-Fecal Collection Tube" (Zymo marketed
10 by Ozyme), "Stool Nucleic Acid Collection and Preservation Tubes" (Norgen Biotek)
11 and "Omnigen Gut for Microbiome" (DNA Genotek).

12 To extract bacterial DNA, we tested five different commercial kits using manual
13 extraction protocols: "PureLink™ Microbiome DNA Purification Kit" (Invitrogen),
14 "Qiam PowerFecal DNA Kit" (Qiagen), "ZimoBiomics DNA Mini Kit" (Ozyme) and
15 "Power Microbiome RNA/DNA isolation kit" (Mo Bio). We used the "MAXWELL
16 Instrument" a robotic station from Promega that extracts DNA from 16 samples
17 simultaneously. We tested automated extraction with two different kits: "Maxwell RSC
18 Buffy Coat DNA Kit" (Promega 1) and "Maxwell RSC PureFood GMO and
19 Authentication Kit" (Promega 2). Extracted stool DNA yield and quality was evaluated
20 with a fluorometer (Qubit, Life Technologies) and Nanodrop (Thermo Scientific),
21 respectively.

1 **Library preparation and sequencing**

2 We used 1.5 µg of DNA to perform the library construction. Extracted DNA was
3 fragmented in g-tubes from Covaris, and DNA end repair was performed using the
4 NEBNext FFPE Repair Mix from New England Biolabs (NEB). We used NEBext's
5 NEBNext Ultra II End Repair / dA-Tailing Module (NEB) for the "end prep" step, 1D
6 Native barcoding genomic DNA kit (ONT) and "NEB Blunt / TA Ligase Master Mix kit
7 (NEB) for DNA multiplexing and adapters ligation. We used Agentcourt AMPure XP
8 (Beckman Coulter) beads for DNA purification.

9 Whole genome metagenomic sequencing was performed with a ONT's MinION tool
10 using with flow cells on which 12 samples were loaded per run. Upon receipt of the
11 flow cells and, prior to sequencing, pore counts were measured using the Platform QC
12 script (MinKNOW from Version 1.4 and further). Flow cells were replaced into their
13 packaging, sealed with tape, and stored at 4°C until further use.

14 33 samples from the Microbaria study were sequenced in parallel with ONT and
15 Illumina Novaseq (2x150bp PE reads). Illumina sequences were processed following
16 the same procedure as described in the original Microbaria study [8] in order to
17 estimate microbial gene richness and the abundances of metagenomic species based
18 on the 9.9-million-gene integrated gene catalog (IGC catalog)[21].

19 **Statistical analyses**

20 All statistical analyses were performed on R v.3.6. Wilcoxon rank-sum tests (for 2-level
21 categorical variables) and Kruskal-Wallis tests (for categorical variables with more than
22 two levels) were used to compare differences in microbial diversity between
23 experimental conditions in different experiments. P-values<0.05 were considered as
24 significant. Spearman correlation tests were used to compare the abundance of

1 taxonomic and functional features between sequencing technologies (SOLiD, Illumina,
2 Nanopore) in Microbaria samples followed by correction for multiple comparison with
3 Benjamini-Hochberg method. Adjusted P-values <0.05 were considered as significant.
4 Permutational analyses of variance (PERMANOVA) with the *adonis* function of vegan
5 R package [33] were used to evaluate the impact of different covariates on microbiome
6 composition in different experiments using Bray-Curtis beta-diversity dissimilarity
7 matrix computed from genus-level abundance data.

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

1 Results

2 Metagenome simulations identified key parameters for ONT microbiome 3 quantification

4 The metagenome simulation approach allowed the evaluation of the impact of different
5 bioinformatic parameters. The evaluation accuracy consisted of comparing the
6 estimated abundance of the microbial features (i.e. species abundance) with the
7 original values used before sequence generation. Over 100 million long reads were
8 generated out of the 250 simulated metagenomic samples from 10 different
9 microbiome compositions at 5 different levels of community complexity for species
10 richness and 5 different levels of sequencing depth (see methods; additional files
11 **Tables S1** and **S2**). These reads were aligned against the 506 reference genome
12 catalog using minimap2 aligner with the map-ont configuration, designed for optimal
13 performance and accuracy with ONT sequencing data [26]. On average, 381000 reads
14 per sample (94% of sample reads) were aligned against the reference genomes. We
15 observed that the reads that could be aligned were 2.5 longer in size (average read
16 length across 250 samples=8064 bp, sd=8) compared to those that could not (average
17 read length=3168 bp, sd=24) (**Figure 2A, Additional file Table S3**), suggesting that
18 read length was a key parameter to optimize ONT sequencing on reference genomes.

19 We evaluated the accuracy of the estimated species abundance and richness to the
20 reference values used for the simulation. Taxonomic profiles were obtained from all
21 minimap2 aligned reads as well as from all primary alignments (PA), defined as the
22 best alignment of a single reads among all possible multiple alignments.

1 For species richness, we observed that both raw and PA-filtered reads detected all
2 reference species in the simulated samples, reaching recall values close to 1 under all
3 simulated community compositions (**Figure 2B**). However, both raw and PA-filtered
4 reads overestimated the number of species especially in the low-richness community
5 compositions. For all community compositions, taxonomic profiles from PA-filtered
6 data reached higher precision values in species richness estimates compared to raw
7 alignments (**Figure 2C**). Importantly, PCoA analyses using a Bray-Curtis beta-diversity
8 dissimilarity matrix with reference and simulated samples showed that PA-filtered
9 samples were more similar to the corresponding reference distributions (R2 effect size
10 dissimilarities=0.04, pvalue<0.05; Permanova test) than raw alignment samples (R2
11 effect size dissimilarities=0.51, p<0.05; Permanova test) (**Figure 2D**). This suggests
12 that the noise introduced by secondary alignments decreased the precision of the
13 workflow.

14 **Alignment identity and read quality affected workflow precision**

15 We evaluated the impact of filtering at different thresholds of identity alignment on the
16 accuracy of the estimated microbiome profiles. Alignment sequence identity (i.e. the
17 ratio of the number of matching bases over the number of bases in a read alignment,
18 including gaps) is a commonly used parameter to filter read alignments, with values
19 around 0.8-0.9 to filter small reads from 2nd generation sequencing technologies
20 against reference genome sequences. This is particularly important for potential error-
21 prone reads generated with ONT sequencing. For species richness, the recall values
22 were close to 1 for identity levels up to 40%, meaning that all reference species in each
23 simulated sample were detected by the workflow. When progressively increasing
24 identity levels from 50% to 90%, a significant fraction of reference species was not
25 detected (**Figure 3A**), although increasing the precision of the estimated richness as

1 the filtering lowered the number of false positives (**Figure 3B**). However, when
2 considering overall microbial composition, the higher the stringency of the alignment
3 identity the more dissimilar the metagenomic profiles were from the reference
4 composition of simulated samples (**Figure 3C**), despite the presence of false positives.
5 Finally, as for the recall, we observed that the Spearman Rho correlation between the
6 estimated species abundance and the reference values decreased as the alignment
7 identity threshold increased. This was in agreement with ordination results (**Figure**
8 **3D**). Overall, these results showed that common approaches to filter read alignments
9 used in the context of second-generation NGS technologies were not directly
10 applicable to ONT sequencing data, probably a consequence of its high error rates.
11 Additional parameters were needed to be explored in order to improve the accuracy of
12 the resulting metagenomic profiles.

13 **Relevance of mapping quality scores in workflow precision**

14 The mapping quality score (mapQ) as computed by minimap2 provided relevant
15 information. High values corresponded to long reads and for which the scores assigned
16 to secondary alignments were weak when compared with primary alignments[26].
17 MapQ score distributions, based on primary alignments of simulated datasets
18 (**Additional file S1**) were evaluated with 11 different thresholds from [0:50], 0 being
19 no filtering. For species richness, we observed an inverse relationship between the
20 stringency of mapQ filtering scores and recall and positive with precision of the
21 workflow (**Figure 4A-B**). Importantly, we observed that mapQ score was more
22 sensitive than the percentage of identity in estimating of species richness. For
23 simulated samples with 50 species each at the lowest filtering threshold (mapQ=5),
24 the achieved precision values were higher (mean=0.55, sd=0.1; Figure 4B) than those
25 reached when filtering with the most stringent sequence identity threshold (90%

1 sequence identity; mean=0.49, sd=0.06; **Figure 3B**). Regarding community
2 composition, our results displayed an inverse relationship between the filtering
3 stringency and the similarity of the estimated communities with the reference simulated
4 samples based on PCoA ordination (**Figure 4C**).

5 Finally, we explored the similarity between estimated species abundance and the
6 reference values using Spearman correlations. We observed different results for
7 different species richness of simulated samples. Whereas in simulated samples with
8 low richness (R50, R150) the similarities between the estimated and the reference
9 species abundance vectors increased with the mapQ threshold from 5 to 30, this was
10 not the case for more complex samples, where the similarities did not improve as the
11 mapQ threshold was increased further than mapQ=5. On the contrary, the similarity
12 significantly decreased in simulated samples with 450 species (**Additional file S2**).

13 Based on these results and considering the trade-off between the accuracy of species
14 richness estimates and the accuracy of microbiome composition for complex microbial
15 communities, we proposed the threshold of mapq5 as the main criteria to process ONT
16 sequence alignments for quantitative metagenomic studies. This filtering was used to
17 process the subsequent experiments aiming to optimize the wet-lab protocols.
18 Moreover, the sequencing of the ZymoBIOMICS mock community and its
19 quantification combining Centrifuge taxonomic binning and filtering by minimap2
20 alignment of read bins vs. the corresponding reference genomes with parameters
21 derived from simulation experiments (primary alignments only, min. mapQ=5)
22 reproduced the composition of the mock community with high accuracy and reducing
23 the number of miss-assignments in comparison with classification based on Centrifuge
24 only (**Additional file S3**). This also led to a higher overall similarity of microbiome
25 composition (estimated as 1- Bray-Curtis beta diversity) with the reference mock

1 community with the combination of Centrifuge and minimap2 filtering (0.91) than with
2 raw Centrifuge results (0.88).

3 **Steps to optimize DNA extraction, DNA fragmentation and End Repair**

4 DNA extraction tests were performed from stool samples collected in dry tubes from
5 three healthy subjects from the MetaCardis cohort (BMI<20kg/m²) at three sampling
6 times for subject 01. After collection, stool samples were aliquoted and immediately
7 stored at -80°C. Each sample was extracted according to the protocols proposed by
8 the manufacturer. After extraction, the samples were evaluated using the "Qubit"
9 fluorometer to estimate the DNA yield obtained in ng/μl and using Nanodrop to
10 evaluate DNA quality.

11 **DNA extraction kits influenced read length distribution**

12 When testing all DNA extraction kits on dry spoon stool samples from healthy
13 volunteers, all kits except the "ZimoBionics DNA Mini Kit" (Ozyme) provided sufficient
14 DNA quantity and quality for sequencing. Thus, we examined library preparation and
15 sequencing all kits except the Ozyme kit. The Invitrogen and Mo Bio kits showed a
16 higher portion of long reads (> 1.1 kb, **Figure 5A**). However, the Mo Bio kit led to the
17 production of sequences with a mean of 8.5kb in size while the Invitrogen kit produced
18 sequences up to 24 kb. The Qiagen kit and the two Promega (Promega 1 and Promega
19 2) kits yielded sequences up to 17kb but with a higher portion of short reads. The
20 examination of alpha diversity showed a significantly higher number of present species
21 for the Invitrogen kit compared to the other kits tested (**Figure 5E**). Based on these
22 observations, the Invitrogen kit was selected as the preferred extraction kit for
23 sequencing.

24 **DNA fragmentation and end-repair steps did not influence alpha diversity**

1 The first step in ONT's library preparation protocol DNA fragmentation to generate 8kb
2 fragments [34]. Using 3 different samples from one subject, DNA fragmentation had
3 no effect on read length distribution (**Figure 5B**), alpha diversity (Wilcoxon rank-sum
4 test $P\text{-value} > 0,05$; **Figure 5F**), or microbiome composition based on PCoA ordination
5 (**Figure 5I**). Therefore, we decided to exclude this step from our experimentation
6 framework. The ONT DNA preparation protocol also recommends *DNA-end repair*. We
7 evaluated the effect of DNA-end repair on read length and microbial diversity by
8 extracting DNA from stools of the same three subjects using the Invitrogen kit and
9 excluding the DNA fragmentation. We found that sequence length (**Figure 5C**) and
10 microbial diversity (**Figure 5G**) were not significantly affected by DNA-end repair. This
11 step was then omitted from our proposed protocol.

12 **DNA extraction kits affected workflow outcomes**

13 According to a PCoA and Permanova analyses, we observed that the differences
14 between microbiome composition of the different replicates were mainly explained by
15 the sampling time ($R^2=0.45$, $p < 0.001$), followed by donor ($R^2=0.03$, $P\text{-value} < 0.05$) and
16 DNA extraction kit ($R^2=0.02$; $P\text{-value} < 0.01$) (**Figure 5H**). Importantly, we observed no
17 effect of sequencing runs nor fragmentation as well as end repair tests (**Figure 5I &**
18 **5J**), suggesting no obvious batch effect.

19 We observed a bimodal distribution of read lengths across the Invitrogen and Mo Bio
20 DNA library preparation kits, whereas with the Promega and Qiagen kits the
21 distribution was skewed towards smaller reads (**Figure 5A**). Based on simulation
22 results, read length had a major impact on the fraction of reads that could be aligned
23 against reference sequences. This ultimately impacted the yield of sequencing runs
24 and the amount of information obtained for microbiome composition.

1 We observed that this was also the case with real data when we classified reads as
2 short or long based on the median of log₂-transformed read lengths and we compared
3 the fraction of classified reads between both groups (**Figure 5A**). We found that the
4 fraction of classified sequences was significantly higher for long reads (log₂-length >
5 9.96; P-value<0.05; Paired Wilcoxon rank-sum test), with on average 39% of long
6 reads successfully classified after the two-step's procedure based on Centrifuge
7 compared with the 24 % for shorter reads (log₂-length<9.96 (**Figure 5D**)). Based on
8 the simulated and real data, we hypothesized that longer ONT reads would increase
9 the accuracy of taxonomic binning. Consequently, in the next step, we aimed at
10 increasing the read length by optimizing DNA extraction and library preparation using
11 the Invitrogen kit.

12 **Optimized DNA extraction protocol improved microbial diversity estimation**

13 To optimize DNA extraction, we followed recommendations from the IHMS consortium
14 [16](**Figure 6A**). To increase the proportion of long reads, we improved the sequencing
15 library preparation protocol by modifying two main steps. The first one was the "End-
16 prep" step which prepares the binding of the adapter to the DNA after two incubation
17 periods. We used the "NEBNext Ultra II End Repair /dA-Tailing Module" from New
18 England Biolabs (NEB) company. In the ONT protocol, "End-prep" reaction incubating
19 is recommended for 5 minutes at 20°C followed by 5 minutes at 65°C. However, the
20 NEB kit recommends a first incubation at 20°C for 30 minutes followed by a second
21 incubation at 65°C for 30 minutes. Given the lack of effects of the ONT end prep
22 protocol, we attempted end repair using NEB kit and recommended protocol (**Figure**
23 **6B**).

1 In the library preparation step, DNA was purified by using "Agencourt AMPure XP"
2 beads (Beckman Coulter), which use SOLiD-phase reversible immobilization (SPRI)
3 paramagnetic bead technology that selectively binds nucleic acids according to type
4 and size. Agencourt AMPure XP utilizes an optimized buffer, Polyethylene glycol
5 (PEG), to selectively bind DNA fragments. The size of the fragments eluted from the
6 beads is determined by PEG concentration. For example, if 50 μ l of beads are added
7 to a 50 μ l DNA sample, a SPRI/DNA ratio of 1 is obtained. When this ratio was changed,
8 the length of the fragments binding and/or remaining in the solution also changed. The
9 SPRI/DNA ratio was disproportionately associated with the DNA fragment size, which
10 is due to fragment size affecting the total charge carried by the molecule. Thus, long
11 DNA fragments would have a greater proportion of negative charges, which promotes
12 their electrostatic interaction with the beads and allows a priority link to the carboxyl
13 molecules. The ONT protocol was developed based on DNA fragmentation of 8Kb
14 sequence length, and the SPRI/DNA ratio must be equal to 1. In order to promote the
15 selection of larger DNA fragments by paramagnetic beads, we reduced the SPRI/DNA
16 ratio to 0.4 (**Figure 6B**). The chosen ratio was based on the SPRI technology
17 documentation [35] and ONT users' recommendations from the "Community" forum
18 [23].

19 Thus, we performed two modifications (End-prep and DNA purification) with the
20 Invitrogen extraction protocol, referred to as "Optimized Invitrogen". This optimization
21 step was performed for six samples from one healthy subject (from the MetaCardis
22 study), collected at six time points. Each sample was extracted using the standard
23 "Invitrogen" protocol and with the "optimized" protocol. DNA yields extracted from this
24 optimized protocol were five-time greater than the ones obtained with the standard kit
25 (55 ng vs. 300 ng, P-value < 0.0001). The ratio of the absorbance at 260/230 was

1 higher with the optimized protocol 1.38 vs 2.11, respectively (P-value = 0.0007) and
 2 the absorbance ratio at 260/280 increased significantly 1.73 vs 1.89, respectively (P-
 3 value = 0.0046) (**Table 1**).

4 **Table 1: Improvement of DNA yield and quality with the Optimized Invitrogen**
 5 **protocol.**

Parameters	Invitrogen protocol	Optimized Invitrogen protocol	P-value
Yield (ng/μl)	55,27 ± 0,56	300,15 ± 0,75	<0,0001
260/280 ratio	1,73 ± 0,03	1,89 ± 0,02	0,0046
260/230 ratio	1,38 ± 0,05	1,87 ± 0,04	0,0007

6
 7 Nanodrop data Estimation of the amount of DNA, protein contamination by given by
 8 the 260/280 ratio and by impurities and solvents contamination estimated by the
 9 260/230 ratio (n = 6, Wilcoxon test).

10
 11 The read length was also improved (**Figure 7A**). The standard Invitrogen protocol
 12 produced two populations of reads with average read lengths of 500 and 6,000 bp
 13 while the optimized protocol produced a single read population with an average read
 14 length of 6,000 bp. With this new dataset, we confirmed that the fraction of classified
 15 reads was higher for long ONT reads (**Figure 7B**; P-value= 4.88e⁻⁰⁴, paired Wilcoxon
 16 rank-sum test). We observed that 29.72% of the reads with the optimized Invitrogen
 17 protocol were successfully classified after the processing with the Centrifuge strategy
 18 in comparison with 23.92% of the reads with the standard protocol (P-value=0.031,

1 paired Wilcoxon rank-sum test). We observed an increased microbial diversity in the
2 samples with the optimized protocol although these differences were not statistically
3 significant (P -value=0.31; Paired Wilcoxon rank-sum test) (**Figure 7C**). The PCoA
4 (**Figure 7D**) confirmed that differences in microbiome composition across samples are
5 explained mainly by subjects. Altogether, the optimized DNA extraction protocol
6 exhibited a better DNA yield and purity, longer sequences than the usual protocol,
7 leading to a significant yield improvement (fraction of classified reads) of the taxonomic
8 binning and to a more accurate microbial diversity assessment.

9 **Impact of stool sampling and storage on sequence length and diversity**

10 Subjects' stool samples were initially collected in a dry spoon tube and rapidly frozen
11 at -80°C to ensure the stability of the bacterial DNA. However, an increasing number
12 sampling systems contain a solution that can stabilize bacterial DNA at room
13 temperature for periods ranging from 60 days (DNA Genotek) up to 2 years (NORGEN
14 Biotek and OZYME). We evaluated the effects of room temperature stabilized samples
15 on bacterial DNA extraction, library preparation and sequencing. We prepared six DNA
16 libraries from stools of 12 healthy subjects collected by different protocols in three
17 different stabilizing kits: "Omnigen Gut for Microbiome" (DNA Genotek), "Stool Nucleic
18 Acid Collection and Preservation Tubes"(Norgen Biotek) and "DNA/RNA Shield-Fecal
19 Collection Tube" (Ozyme). For each sampling kit, we tested the effect of sample
20 storage temperature on DNA extraction and did not observe any significant difference
21 in the distribution of classified reads or microbial diversity in any of the three different
22 collection kits at -80°C (P -value=0.73; Kruskal-Wallis test) or at room temperature (P -
23 value=0.44; Kruskal-Wallis test) (**Figure 8B**), although we could notice a tendency with
24 Norgen and Omnigen kits to decrease microbial diversity at room temperature in
25 comparison with -80°C storage (**Figure 8C**). In contrast, we observed significant

1 variations in microbial diversity by donor (**Figure 8D**; P-value=0.0018, Kruskal-Wallis
2 test). Similarly, donor was the variable with the highest impact on microbiome
3 composition by PERMANOVA analyses (R²=0.82; P-value=0.001) in comparison with
4 collection kit (R²=0.11, P-value=0.009) and storage conditions (R²=0.07; P-
5 value=0.029) (**Figure 8E & 8F**). Thus, sampling methods at room temperature with
6 DNA stabilization performed similarly to snap frozen samples and the choice of the
7 sampling kit might depend on practical feasibility sampling for the subject as well as kit
8 price. We chose the Ozyme kit for its practicality for users to collect stool samples and
9 its relative costs.

10 **Optimized ONT protocol compared to main stream sequencers**

11 We compared ONT obtained QM profiles with those generated with other sequencing
12 technologies from the Microbaria study[8]. We selected 33 baseline samples covering
13 the extremes of microbiome diversity defined as microbial gene richness (13 samples
14 from individuals with High Gene Count (HGC) and 20 samples from individuals with
15 Low Gene Count (LGC)). ONT abundance profiles were generated using the two
16 bioinformatics workflows described above, based on Centrifuge and mapping over the
17 IGC catalog. DNA from 21 of the 33 samples were extracted with the optimized
18 Invitrogen protocol and sequenced using Illumina technology. Quantitative
19 metagenomic profiles were generated by mapping reads against the IGC gene catalog
20 [8].

21 First, we compared the estimates of microbial diversity from ONT (GeneRichness from
22 IGC mapping and Observed Species from Centrifuge classification) and Illumina
23 sequencing (Gene richness from IGC mapping) with the gene richness inferred from
24 the original SOLiD sequencing of these samples. SOLiD sequencing generated

1 4.38e+07 single reads of 35 bases (sd=1.86e+07) per sample on average,
2 representing 1.53e+09 base pairs overall (sd=6.52e+08). With the ONT, we generated
3 an average of 1.53e+05 reads per sample between 200bp and 24 Kb (sd=6.07e+04),
4 representing 4.19e+08 bp overall (sd=1.607e+08).

5 We observed significant positive associations between diversity estimates based on
6 ONT sequencing with the reference gene richness estimates from SOLiD sequencing
7 based on IGC gene catalog (Spearman Rho=0.59, P-value=3e-04 for Observed
8 Species based on Centrifuge results; Spearman Rho=0.74, P-value=2e-06 for Gene
9 Richness based on ONT read mapping over IGC catalog; **Figure 9A & 9B**). These
10 similarities increased with the use of reference databases maximizing the genomic
11 information of the gut microbiome ecosystem (gut microbiome gene catalogs vs.
12 generic Centrifuge database). However, the similarity was higher with gene richness
13 estimates based on Illumina sequencing despite differences in library preparation
14 (Spearman Rho=0.86, $p < 2.2e-16$; **Figure 9C**). When we integrated the scaled diversity
15 profiles of different sources available for these 33 samples ordering them based on the
16 reference gene richness[8], we observed that DNA extraction had an impact on
17 diversity. Both ONT and Illumina sequencing using the same DNA extraction method
18 showed similar variations in microbial diversity estimates compared to the reference
19 SOLiD data. This included a switch in the sample showing the highest diversity (i.e.
20 MB12 with Illumina and ONT sequencing; MB21 with SOLiD sequencing; **Figure 9D**).

21 This was confirmed in an ordination framework where we integrated the genus-level
22 abundance profiles from IGC quantification with the three sequencing technologies
23 (ONT, Illumina, SOLiD), where we observed that samples product of the same DNA
24 extraction method (Illumina and ONT) are closer in a PCoA ordination (**Additional file**
25 **S4A**) and in hierarchical clustering analyses (**Additional file S4B**).

1 **Potential of the ONT workflow for the detection of target species and functional** 2 **profiles**

3 Regarding taxonomic feature quantification, we found a good agreement between
4 ONT sequencing and both Illumina and SOLiD sequencing data. Based on Centrifuge,
5 we observed a positive correlation in 91 of the 95 common taxonomic features with
6 SOLiD quantifications (96%), of which 72 (76%) were significantly associated
7 (FDR<0.05, Spearman correlations). Similarly, we observed a positive correlation in
8 94 of the 101 common taxonomic features with Illumina quantifications (93%), of which
9 78 (77 %) were significantly associated (FDR<0.05, Spearman correlations) (**Figure**
10 **9E**). Using IGC results at MGS level, 137 common taxonomic features with SOLiD
11 quantifications were positively associated, 128 of which (93%) were significantly
12 associated (FDR<0.05, Spearman correlations), whereas 133 of the 137 common
13 taxonomic features with Illumina quantifications were positively associated (98 %), 122
14 of which (90%) were significantly associated (FDR<0.05, Spearman correlations)
15 (**Additional file, Table 4**). Importantly, we observed that the similarities between
16 taxonomic features between ONT and Illumina quantifications were higher than
17 between ONT and SOLiD sequencing (**Additional file S5**; $p < 0.05$ for comparisons at
18 species and genus level with ONT Centrifuge results; $p < 0.005$ for comparisons at
19 species, genus, and family level with ONT IGC results).

20 We made similar observations with functional profiles based on KEGG modules. Using
21 Centrifuge, 76% and 72% of the functional modules were positively associated with
22 the equivalent modules quantified with Illumina and SOLiD sequencing respectively,
23 whereas this fraction substantially increased to 98% and 98% with ONT abundance
24 data based on IGC quantifications (**Additional file S6A**). This difference may be
25 related to the different content of both genomic reference spaces (Centrifuge genomes;

1 IGC gene catalog), which can have a major impact on the quantification of functional
2 modules if differences in composition also result in differences in gene content.
3 Importantly, we observed that DNA extraction also impacted the functional profiles,
4 with ONT functional profiles being more similar to Illumina functional profiles based
5 both on Centrifuge and IGC quantifications (**Additional file 6B**; P-value=0.0046 for
6 Centrifuge-based quantifications; P-value=6.1e-14 for IGC-based quantifications;
7 Wilcoxon rank-sum test of Spearman's Rho distributions between ONT-SOLID
8 comparisons and ONT-Illumina comparisons).

9 We reproduced these associations at the functional level. We found significant positive
10 associations between the sporulation module md:M00485 (KinABCDE-Spo0FA
11 (sporulation control) two-component regulatory system) and microbial diversity
12 (**Additional file S7**), which was in agreement with estimations of 50%-60% of bacteria
13 from gut microbiome of healthy individuals producing resilient spores. This sporulation
14 phenotype is an unappreciated and basic feature of the human microbiome with a key
15 impact in bacterial persistence and the spread of microbes between individuals [36].
16 This was also the case for the negative association between modules involved in the
17 biosynthesis of bacterial Lipopolysaccharide (LPS) and microbial diversity (**Additional**
18 **file S8**), in line with the association of obesity and other metabolic disorders with an
19 increase of blood LPS concentration[37].

20 **The ONT workflow enabled de novo genome assembly**

21 We finally explored the potential of ONT sequencing for de-novo assembly of bacterial
22 genomes from shotgun metagenomics data. We carried out de-novo assembly of the
23 genomes of *Bacteroides vulgatus* (a bacterial species with the highest mean
24 abundance based on Centrifuge quantifications, **Additional file S9A**) and
25 *Akkermansia muciniphila*, a bacterium whose abundance in this cohort was lower than

1 in other overweight or obese cohorts [38]. We observed that de-novo assemblies from
2 the corresponding read bins product of Centrifuge workflow gave genomic scaffolds
3 that covered 88% of the *Bacteroides vulgatus* ATCC 8482 genome at 99% Average
4 Nucleotidic Identity (ANI) and 78% of the *Akkermansia muciniphila* isolate Urmite at
5 95% ANI (**Additional file S9B & S9C**). These reference genomes were the closest
6 relatives to the assembled ONT genomes based on mash genomic distances vs.
7 RefSeq genomes[39]. The structural annotation of these genome assemblies also
8 revealed a highly fragmented proteome, which is reflected in the smaller gene lengths
9 observed in ONT assemblies in comparison with the corresponding reference
10 genomes (**Additional file S9D & S9E**) and in the large number of genes in ONT
11 assemblies in comparison with the corresponding reference genomes (15532 CDS in
12 *Bacteroides vulgatus* assembly vs. 4104 in *Bacteroides vulgatus* ATCC 8482 genome;
13 2900 CDS in *Akkermansia muciniphila* assembly vs. 2101 CDS in *Akkermansia*
14 *muciniphila* isolate Urmite genome). Thus, we were able to recover near complete
15 genomes for individual species bins directly from ONT metagenome sequencing of
16 stool samples but the error-prone character of this sequencing technology prevented
17 accurate structural annotation of these genomes uniquely with ONT sequencing data.

18
19
20
21
22
23

1 **Discussion**

2 We present here a novel protocol and analytical pipeline enabling the quantification of
3 the gut microbiome features using Oxford Nanopore Technologies. This technology
4 potentially supports easy access and use of high throughput sequencing at competitive
5 costs as well as fast data production and analyses of the results. We believe this
6 established protocol enables the study of the gut microbiome in the context of clinical
7 applications or group studies. To generate optimal results, we optimized protocols for
8 the wet-lab (from DNA extraction to sequencing) and data analysis. We also compared
9 the final results to state-of-the-art sequencing methods (Illumina and SOLiD) in an
10 already described patient cohort. This was driven by 1) an initial assessment of the
11 best parameters in terms of alignment of ONT sequencing reads from simulated
12 metagenomic datasets with different levels of complexity and, 2) the development of a
13 bioinformatic pipeline that combines rapid k-mer based classification of ONT reads
14 with read alignments vs. reference genomes to improve the quantification of
15 microbiome species diversity and composition. The simulation experiments revealed
16 that filtering strategies commonly used with second generation sequencing
17 technologies, such as high sequence identity thresholds, could not be extrapolated to
18 highly error-prone reads such as those produced by ONT. Also, the alignment quality
19 based on Nanopore-adapted sequence aligners like minimap2 allowed for improved
20 accuracy in the estimation of microbial diversity and composition in complex
21 ecosystems such as the human gut.

22 Regarding sample processing, a first step required the elaboration of a DNA extraction
23 protocol from human stools that provide high DNA quality. Several studies over the
24 past years have used bacterial DNA or RNA to explore microbial communities in
25 diverse ecosystems including stool samples from large cohorts [2][40][41]. Authors in

1 these areas have used different DNA extraction protocols and different sequencing
2 techniques (Illumina, SOLiD, Ion Proton). Multiple studies have also noted “batch”[42]
3 effects and differences in data analyses [43][44], which introduce biases in study
4 analytical comparisons. Thus, the need for procedure standardization has been
5 highlighted by several reports, as illustrated by the IHMS consortium[45]. In this study,
6 the authors compared 21 DNA extraction methods using whole genome metagenomic
7 shotgun sequencing with Illumina HiSeq2000 technology. They assessed the
8 taxonomic profile and functional variability while standardizing the stages of stool
9 collection, bacterial DNA stabilization, library preparation and sequencing. This
10 resulted in the generation of "Standard Operating Procedure "SOPs" with
11 recommendations that would improve DNA extraction in terms of yield and quality.

12 We applied these IHMS recommendations to the PureLink™ Microbiome DNA
13 Purification kit from Invitrogen, and we further optimized the microbial DNA extraction
14 protocol. The DNA yield from the optimized protocol was significantly improved,
15 compared to the conventional kit in our hands. We worked on two critical steps,
16 bacterial wall lysis and protein/RNA elimination. Two rounds of lysis process gave two
17 supernatants with higher DNA yield. DNA recovery by isopropanol precipitation
18 followed by elimination of protein and RNA contamination allowed further DNA purity.
19 We also explored other methods of DNA stabilization. Sampling conditions, storage
20 and harmonization have previously been shown to be critical in affecting microbiome
21 results. Although storing fecal samples at 4°C appeared to protect bacterial DNA from
22 degradation to some extent, a reduction in microbial diversity was observed [46]. A
23 previous study showed that prior storage of stool samples at 4°C (one hour) before
24 placing them at -20°C, had a large impact on the taxonomic composition at the genus
25 and species level[47]. However, these studies were conducted before the development

1 and widespread use of commercially available fecal collection kits with stabilizing
2 solution. In our current study, our results suggest that sample storage temperature is
3 not a significant factor as long as guidelines from manufacturers are followed.
4 Furthermore, the effect of sample storage kit type or temperature on sequencing and
5 microbiome results are largely outweighed by inter-donor variation.

6 Since we identified read length as a critical criterion for subsequent bioinformatics
7 analyses, we selected the DNA extraction protocol and improved the library
8 preparation protocol to increase the proportion of long reads. For this purpose, we
9 optimized the end-prep and DNA purification steps in the library. Applying this
10 approach, we obtained a unimodal read length distribution curve and observed an
11 increase of the proportion of long reads. PCoA of different wet-lab experiments showed
12 that individual microbiome composition drove most of the variation observed in
13 microbial diversity and quantitative metagenomic profiles obtained from ONT
14 sequencing data, with no apparent batch effects associated to different wet-lab steps
15 (DNA fragmentation, DNA end-repair, Collection kits, DNA library preparation or
16 sequencing run).

17 To examine the relevance of our pipeline in human cohorts, we performed comparison
18 of the results obtained with ONT sequencing with those obtained with SOLiD and
19 Illumina sequencing on human stool samples collected in the "Microbaria" study. We
20 observed that for gene richness, microbiome composition and functional modules, the
21 similarity was higher between ONT and Illumina sequencing compared to SOLiD. ONT
22 and Illumina sequences were generated from the same DNA extracted with the
23 optimized protocol developed in this study, which emphasizes the importance of DNA
24 library preparation protocols in quantitative metagenomic profiles in our experimental
25 design.

1 The long reads obtained with ONT sequencing allowed us to perform "de novo"
2 assemblies for the most abundant species in the Microbaria cohort (*B. vulgatus*) as
3 well as potential metabolic disease-relevant bacteria such as *Akkermansia*
4 *muciniphila*, whose abundance in the sequenced samples is lower in comparison with
5 other less obese or overweight individuals, and is associated with metabolic health[38]
6 [48]. These assemblies showed a high level of completion when compared with closest
7 relatives, highlighting the strong potential of ONT sequencing in the assembly of
8 bacterial genomes directly from shotgun metagenomics data. In comparison, the
9 assembly of the *Akkermansia muciniphila* genome in a previous study required
10 "SOLiD" and shotgun sequencing to cover the genome with a scaffold and 56 gaps[49].
11 However, the structural annotation of the assembled genomes revealed a caveat
12 associated with this sequencing technology, which is the poor quality of the
13 assemblies, that we observed, with highly fragmented genes as a consequence of the
14 presence of frameshifts in the assembled sequences (as shown in **Additional file S8**).
15 This is explained by the low read accuracy of ONT reads in comparison with short-
16 read technologies, which makes hybrid strategies combining long reads product of
17 ONT sequencing with small reads product of short-read sequencing the most prevalent
18 approach for accurate genome assemblies [50][51].

19 It is also important to mention that longer reads generated from ONT did not
20 compensate for the overall sequence coverage in comparison with other sequencing
21 technologies. The low throughput of ONT sequencing is one of its major drawbacks for
22 quantitative metagenomic studies of complex microbial ecosystems like the human gut
23 microbiome. However, the results obtained showed a high similarity in bacterial
24 diversity estimation between the two sequencing methods in our test samples. Despite
25 improvements in experimental protocol such as DNA extraction and library

1 preparation, the demultiplexing step needs to be improved. For instance, the ratio of
2 unclassified reads was about 25%, knowing that the sequencing depth of ONT was
3 low, and the error rate elevated. However, the low classification rates (25%) did not
4 seem to impact the bacterial diversity and the bacterial compositions estimates. In this
5 context, we could consider ONT in the context of quantitative metagenomic studies as
6 a “shallow-sequencing” method in the line of proposed low-sequencing depth
7 approaches to characterize microbial ecosystems more accurately than 16S barcoding
8 approaches and with lower costs than deep shotgun sequencing[52].

9 Nanopore-based technology such as that used by ONT, is proposed as easily
10 accessible due to relatively low costs and a small benchtop footprint, providing an
11 avenue to perform NGS in clinical settings. Admittedly, this technology has some
12 drawbacks such as relatively modest sequencing depth, and error rates that remain
13 high (2 - 5%)[53] compared to Illumina (0,1%)[54]. However, in our current work, ONT
14 consistently replicated results for intestinal microbiome diversity and the composition
15 of main phyla in patients with severe obesity.

16

17

18

19

20

21

22

1 **Conclusion**

2 ONT should be considered as a suitable sequencing method in a clinical setting, and
3 we propose to pair this technology with AI-based approaches in the future, which may
4 allow the identification of interpretable and accurate signatures for disease diagnosis
5 and prognosis in an effort to adapt future personalized microbiome based therapies
6 [55].The present work indeed provides a useable and complete experimental and
7 analytical workflow for the use of ONT sequencing in quantitative metagenomic
8 studies. Both computational and wet-lab steps has been benchmarked to improve the
9 precision of quantitative metagenomic profiles and microbial diversity estimates. The
10 developed bioinformatic pipeline allows the quantification of metagenomic features in
11 different reference spaces (non-redundant gene catalogs, reference genomes with
12 Centrifuge) from raw ONT FAST5 files, including functional characteristics of
13 metagenomic species. This is a key dimension of metagenomic analyses to
14 understand the impact of variations in the composition of the microbiome from a
15 functional point of view. Wet lab experiments revealed that it is DNA extraction
16 technique rather than the sequencing technology what predominantly drives the
17 observed variation in microbial diversity and microbiome composition. Moreover,
18 despite variations in library preparation metagenomic profiles inferred from ONT
19 sequencing, data appears close to those obtained with Illumina sequencing in terms
20 of microbial diversity and taxonomic and functional profiles. Through simulations and
21 wet lab experiments read length emerged as a critical parameter to optimize in order
22 to maximize the classification of ONT reads. The experimental protocol developed
23 allowed this optimization providing not only a correct coverage of the taxonomic
24 profiles with the reduction of unclassified reads, but also in the assembly of near
25 complete microbial genomes even from limited number of reads. This could contribute

1 to improve the knowledge and completeness of microbial genome assemblies of
2 complex ecosystems like the human gut.

3 Overall, despite caveats and limitations of ONT sequencing regarding throughput and
4 sequencing errors, this proposed workflow paves the way to taxonomic and functional
5 profiling of microbial communities with this sequencing technology at competitive costs
6 and fast data, which corresponds to a great need in the microbiome community.

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

1 **Ethics approval and consent to participate**

2 All subjects from the MetaCardis study or the Microbaria study provided written
3 informed consent and the study was conducted in accordance with the Helsinki
4 Declaration. MetaCardis study is registered in clinical trial
5 <https://clinicaltrials.gov/show/NCT02059538> and Microbaria study was registered in
6 clinicaltrial.gov (NCT01454232).

7 **Consent for publication**

8 All authors gave consent for publication.

9 **Availability of data and material**

10 Sequences have been deposited in the European Bioinformatics Institute (EBI)
11 European Nucleotide Archive (ENA) under accession **number XXXX** (Private access
12 until paper acceptance). The computational pipeline is freely available in
13 <https://git.ummisco.fr/ebelda/nanopore>. Other data are available on request.

14 **Competing interests**

15 KC is a consultant for Danone Research, LNC therapeutics and CONFO therapeutics
16 for work unassociated with the present study. J-D.Z. is consultant for Quinten for work
17 unassociated with the present study.

18 **Funding**

19 This study was supported by the European Union FP7 Metacardis (grant agreement
20 HEALTH- F4-2012-305312), H2020 EPoS (H2020- PHC-2014-634413), the Innovative
21 Medicines Initiative 2 (IMI-2) Joint Undertaking under grant agreement No. 777377.
22 This Joint Undertaking receives support from the European Union's Horizon 2020
23 research and innovation programme and EFPIA. This study was also supported by

1 Transatlantic Networks of Excellence Award from the Leducq Foundation (17CVD01)
2 and JPI (A healthy diet for a healthy life; 2017-01996_3). The clinical study (Microbaria)
3 was sponsored by Assistance Publique Hopitaux de Paris.

4 **Author contribution**

5 K.C., E.P, E.B and J-D.Z conceived and designed the project. K.C. is the primary
6 investigator of the clinical study, from which, samples were used in this work. R.A.
7 performed all the wet-lab experiments and established the protocol (from stool sample
8 processing to sequencing). E.B developed the databases, analytical pipelines and
9 performed metagenomics and functional analysis. E.P. conceived the simulation and
10 validation protocol and supervised the development of the bioinformatics workflows.
11 K.C, J-D.Z., E.P., R.A., P.L. T.W. and E.B. contributed to results' discussion. R.A.,
12 E.B., E.P. and K.C wrote the paper and all authors commented and edited the
13 manuscript.

14 **Acknowledgements**

15 We thank Tim Swartz (Integrative Phenomics) for critical reading of the manuscript and
16 English correction.

17 **Authors' information**

18 Eugeni Belda is employee of Integrative phenomics but this activity is not related to
19 the topic of this publication

20

21

1 References

- 2 1. Qin J, Li R, Raes J, Arumugam M, Burgdorf KS, Manichanh C, et al. A human gut
3 microbial gene catalogue established by metagenomic sequencing. *Nature*.
4 2010;464:59–65.
- 5 2. Vieira-Silva S, Falony G, Belda E, Nielsen T, Aron-Wisnewsky J, Chakaroun R, et
6 al. Statin therapy is associated with lower prevalence of gut microbiota dysbiosis.
7 *Nature*. 2020;581:310–5.
- 8 3. Aron-Wisnewsky J, Gaborit B, Dutour A, Clement K. Gut microbiota and non-
9 alcoholic fatty liver disease: new insights. *Clin Microbiol Infect*. 2013;19:338–48.
- 10 4. Aron-Wisnewsky J, Vigliotti C, Witjes J, Le P, Holleboom AG, Verheij J, et al. Gut
11 microbiota and human NAFLD: disentangling microbial signatures from metabolic
12 disorders. *Nat Rev Gastroenterol Hepatol*. 2020;17:279–97.
- 13 5. Turnbaugh PJ, Ley RE, Mahowald MA, Magrini V, Mardis ER, Gordon JI. An obesity-
14 associated gut microbiome with increased capacity for energy harvest. *Nature*.
15 2006;444:1027–31.
- 16 6. Qin N, Yang F, Li A, Prifti E, Chen Y, Shao L, et al. Alterations of the human gut
17 microbiome in liver cirrhosis. *Nature*. 2014;513:59–64.
- 18 7. Cotillard A, Kennedy SP, Kong LC, Prifti E, Pons N, Le Chatelier E, et al. Dietary
19 intervention impact on gut microbial gene richness. *Nature*. 2013;500:585–8.
- 20 8. Aron-Wisnewsky J, Prifti E, Belda E, Ichou F, Kayser BD, Dao MC, et al. Major
21 microbiota dysbiosis in severe obesity: fate after bariatric surgery. *Gut*. 2019;68:70–
22 82.
- 23 9. Raes J, Bork P. Molecular eco-systems biology: towards an understanding of
24 community function. *Nat Rev Microbiol*. 2008;6:693–9.
- 25 10. Jones MB, Highlander SK, Anderson EL, Li W, Dayrit M, Klitgord N, et al. Library
26 preparation methodology can influence genomic and functional predictions in human
27 microbiome research. *Proc Natl Acad Sci USA*. 2015;112:14024–9.
- 28 11. Henderson G, Cox F, Kittelmann S, Miri VH, Zethof M, Noel SJ, et al. Effect of DNA
29 Extraction Methods and Sampling Techniques on the Apparent Structure of Cow and
30 Sheep Rumen Microbial Communities. *PLOS ONE*. 2013;8:e74787.
31 doi:10.1371/journal.pone.0074787.
- 32 12. Santiago A, Panda S, Mengels G, Martinez X, Azpiroz F, Dore J, et al. Processing
33 faecal samples: a step forward for standards in microbial community analysis. *BMC*
34 *Microbiology*. 2014;14:112. doi:10.1186/1471-2180-14-112.
- 35 13. Kennedy NA, Walker AW, Berry SH, Duncan SH, Farquarson FM, Louis P, et al.
36 The Impact of Different DNA Extraction Kits and Laboratories upon the Assessment of
37 Human Gut Microbiota Composition by 16S rRNA Gene Sequencing. *PLOS ONE*.
38 2014;9:e88982. doi:10.1371/journal.pone.0088982.
- 39 14. Voigt AY, Costea PI, Kultima JR, Li SS, Zeller G, Sunagawa S, et al. Temporal and
40 technical variability of human gut metagenomes. *Genome Biol*. 2015;16:73.
- 41 15. Turnbaugh PJ, Ley RE, Hamady M, Fraser-Liggett CM, Knight R, Gordon JI. The
42 human microbiome project. *Nature*. 2007;449:804–10.
- 43 16. Costea PI, Zeller G, Sunagawa S, Pelletier E, Alberti A, Levenez F, et al. Towards
44 standards for human fecal sample processing in metagenomic studies. *Nat Biotechnol*.
45 2017;35:1069–76.
- 46 17. Harstad H, Ahmad R, Bredberg A. Nanopore-based DNA sequencing in clinical
47 microbiology: preliminary assessment of basic requirements. *bioRxiv*. 2018;:382580.
48 doi:10.1101/382580.

- 1 18. Nayfach S, Pollard KS. Toward Accurate and Quantitative Comparative
2 Metagenomics. *Cell*. 2016;166:1103–16. doi:10.1016/j.cell.2016.08.007.
- 3 19. Cranmer K, Brehmer J, Louppe G. The frontier of simulation-based inference. *Proc*
4 *Natl Acad Sci U S A*. 2020.
- 5 20. Nielsen HB, Almeida M, Juncker AS, Rasmussen S, Li J, Sunagawa S, et al.
6 Identification and assembly of genomes and genetic elements in complex
7 metagenomic samples without using reference genomes. *Nat Biotechnol*.
8 2014;32:822–8.
- 9 21. Li J, Jia H, Cai X, Zhong H, Feng Q, Sunagawa S, et al. An integrated catalog of
10 reference genes in the human gut microbiome. *Nat Biotechnol*. 2014;32:834–41.
- 11 22. Fritz A, Hofmann P, Majda S, Dahms E, Dröge J, Fiedler J, et al. CAMISIM:
12 simulating metagenomes and microbial communities. *Microbiome*. 2019;7:17.
- 13 23. <https://nanoporetech.com/community>.
- 14 24. Wick RR, Judd LM, Holt KE. Performance of neural network basecalling tools for
15 Oxford Nanopore sequencing. *Genome Biology*. 2019;20:129. doi:10.1186/s13059-
16 019-1727-y.
- 17 25. Kim D, Song L, Breitwieser FP, Salzberg SL. Centrifuge: rapid and sensitive
18 classification of metagenomic sequences. *Genome Res*. 2016;26:1721–9.
- 19 26. Li H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*.
20 2018;34:3094–100.
- 21 27. Chamberlain SA, Szöcs E. taxize: taxonomic search and retrieval in R. *F1000Res*.
22 2013;2:191. doi:10.12688/f1000research.2-191.v2.
- 23 28. McMurdie PJ, Holmes S. phyloseq: an R package for reproducible interactive
24 analysis and graphics of microbiome census data. *PLoS ONE*. 2013;8:e61217.
- 25 29. Lawrence M, Huber W, Pagès H, Aboyoun P, Carlson M, Gentleman R, et al.
26 Software for Computing and Annotating Genomic Ranges. *PLOS Computational*
27 *Biology*. 2013;9:e1003118. doi:10.1371/journal.pcbi.1003118.
- 28 30. Lee S, Cook D, Lawrence M. plyranges: a grammar of genomic data
29 transformation. *Genome Biology*. 2019;20:4. doi:10.1186/s13059-018-1597-8.
- 30 31. <http://meta.genomics.cn/meta/home>.
- 31 32. Kanehisa M, Goto S, Sato Y, Furumichi M, Tanabe M. KEGG for integration and
32 interpretation of large-scale molecular data sets. *Nucleic Acids Res*. 2012;40 Database
33 issue:D109-114.
- 34 33. <https://CRAN.R-project.org/package=vegan>.
- 35 34. <https://community.nanoporetech.com/protocols/native-barcoding-genomic-dna/>.
- 36 35. James@cancer. CoreGenomics: How do SPRI beads work? CoreGenomics. 2012.
37 <http://core-genomics.blogspot.com/2012/04/how-do-spri-beads-work.html>. Accessed
38 18 Nov 2020.
- 39 36. Hp B, Sc F, Bo A, N K, Ba N, Md S, et al. Culturing of “unculturable” human
40 microbiota reveals novel taxa and extensive sporulation. *Nature*. 2016;533:543–6.
41 doi:10.1038/nature17645.
- 42 37. Krajmalnik-Brown R, Ilhan Z-E, Kang D-W, DiBaise JK. Effects of gut microbes on
43 nutrient absorption and energy regulation. *Nutr Clin Pract*. 2012;27:201–14.
- 44 38. Dao MC, Belda E, Prifti E, Everard A, Kayser BD, Bouillot J-L, et al. Akkermansia
45 muciniphila abundance is lower in severe obesity, but its increased level after bariatric
46 surgery is not associated with metabolic health improvement. *Am J Physiol Endocrinol*
47 *Metab*. 2019;317:E446–59.
- 48 39. Ondov BD, Treangen TJ, Melsted P, Mallonee AB, Bergman NH, Koren S, et al.
49 Mash: fast genome and metagenome distance estimation using MinHash. *Genome*
50 *Biology*. 2016;17:132. doi:10.1186/s13059-016-0997-x.

- 1 40. Bäckhed F, Roswall J, Peng Y, Feng Q, Jia H, Kovatcheva-Datchary P, et al.
2 Dynamics and Stabilization of the Human Gut Microbiome during the First Year of Life.
3 *Cell Host Microbe*. 2015;17:690–703.
- 4 41. Yachida S, Mizutani S, Shiroma H, Shiba S, Nakajima T, Sakamoto T, et al.
5 Metagenomic and metabolomic analyses reveal distinct stage-specific phenotypes of
6 the gut microbiota in colorectal cancer. *Nat Med*. 2019;25:968–76.
- 7 42. Wesolowska-Andersen A, Bahl MI, Carvalho V, Kristiansen K, Sicheritz-Pontén T,
8 Gupta R, et al. Choice of bacterial DNA extraction method from fecal material
9 influences community structure as evaluated by metagenomic analysis. *Microbiome*.
10 2014;2:19.
- 11 43. McOrist AL, Jackson M, Bird AR. A comparison of five methods for extraction of
12 bacterial DNA from human faecal samples. *J Microbiol Methods*. 2002;50:131–9.
- 13 44. Ariefdjohan MW, Savaiano DA, Nakatsu CH. Comparison of DNA extraction kits
14 for PCR-DGGE analysis of human intestinal microbial communities from fecal
15 specimens. *Nutr J*. 2010;9:23.
- 16 45. <http://www.microbiome-standards.org/>.
- 17 46. Ott SJ, Musfeldt M, Timmis KN, Hampe J, Wenderoth DF, Schreiber S. In vitro
18 alterations of intestinal bacterial microbiota in fecal samples during storage. *Diagn*
19 *Microbiol Infect Dis*. 2004;50:237–45.
- 20 47. Cardona S, Eck A, Cassellas M, Gallart M, Alastrue C, Dore J, et al. Storage
21 conditions of intestinal microbiota matter in metagenomic analysis. *BMC Microbiol*.
22 2012;12:158.
- 23 48. Dao MC, Everard A, Aron-Wisnewsky J, Sokolovska N, Prifti E, Verger EO, et al.
24 *Akkermansia muciniphila* and improved metabolic health during a dietary intervention
25 in obesity: relationship with gut microbiome richness and ecology. *Gut*. 2016;65:426–
26 36.
- 27 49. Caputo A, Dubourg G, Croce O, Gupta S, Robert C, Papazian L, et al. Whole-
28 genome assembly of *Akkermansia muciniphila* sequenced directly from human stool.
29 *Biol Direct*. 2015;10:5.
- 30 50. Rang FJ, Kloosterman WP, de Ridder J. From squiggle to basepair: computational
31 approaches for improving nanopore sequencing read accuracy. *Genome Biology*.
32 2018;19:90. doi:10.1186/s13059-018-1462-9.
- 33 51. Moss EL, Maghini DG, Bhatt AS. Complete, closed bacterial genomes from
34 microbiomes using nanopore sequencing. *Nature Biotechnology*. 2020;38:701–7.
35 doi:10.1038/s41587-020-0422-6.
- 36 52. Hillmann B, Al-Ghalith GA, Shields-Cutler RR, Zhu Q, Gohl DM, Beckman KB, et
37 al. Evaluating the Information Content of Shallow Shotgun Metagenomics. *mSystems*.
38 2018;3. doi:10.1128/mSystems.00069-18.
- 39 53. Tedersoo L, Drenkhan R, Anslan S, Morales-Rodriguez C, Cleary M. High-
40 throughput identification and diagnostics of pathogens and pests: Overview and
41 practical recommendations. *Mol Ecol Resour*. 2019;19:47–76.
- 42 54. Fox EJ, Reid-Bayliss KS, Emond MJ, Loeb LA. Accuracy of Next Generation
43 Sequencing Platforms. *Next Gener Seq Appl*. 2014;1.
- 44 55. Prifti E, Chevalleyre Y, Hanczar B, Belda E, Danchin A, Clément K, et al.
45 Interpretable and accurate prediction models for metagenomics data. *Gigascience*.
46 2020;9.
- 47

1 **Figures :**

2 **Figure 1: Summary of the workflow** : (A) Simulated data processing. (B) Wet-lab
3 optimization. (C) Summary of ONT sequencing comparison with Illumina and SOLiD
4 technologies.

5 **Figure 2: Metagenomic profiles from simulated samples between minimap2**
6 **results and minimap2 results filtered from secondary alignments.** (A) Boxplots of
7 mean lengths of Nanopore reads of 250 simulated samples (y axis) between those
8 aligned and unaligned over the 506 reference genomes from minimap2 results (x-axis).
9 (B) Boxplots of recall values of species richness estimates in 250 simulated samples
10 (y-axis) between metagenomic profiles inferred from all minimap2 alignments (mmap2
11 raw) and from minimap2 primary alignments only (mmap2APfilt, x-axis). (C) Boxplots
12 of precision values of species richness estimates in 250 simulated samples (y-axis)
13 between metagenomic profiles inferred from all minimap2 alignments (mmap2 raw)
14 and from minimap2 primary alignments only (mmap2APfilt, x-axis). (D) Principal
15 Coordinates Analysis (PCoA) of metagenomic profiles from the reference and 250
16 simulated samples inferred from all minimap2 alignments (mmap2raw) and from
17 minimap2 primary alignments only (mmap2APfilt, x-axis). Dashed lines connect points
18 coming from the same sample (reference, simulated ones; 3 points per sample).

19

20 **Figure 3: Impact of filtering minimap2 primary alignments of Nanopore reads at different**
21 **thresholds of sequence identity.** Boxplots of recall (A) and precision (B) values of species
22 richness estimates in 250 simulated samples (y-axis) between metagenomic profiles inferred
23 from primary alignments of Nanopore reads filtered by different thresholds of sequence identity
24 (from 0 to 90%; x-axis) stratified by the number of species in reference metagenomic profiles.
25 (C) PCoA of reference metagenomic profiles and metagenomic profiles of the 250 simulated

1 samples inferred from minimap2 primary alignments filtered by different thresholds of
2 sequence identity (from 0 to 90%; x-axis). Dashed lines connect points coming from the same
3 sample (reference, simulated ones; 11 points per sample) with different shapes assigned to
4 samples from different reference species richness. Points corresponding to reference samples
5 and simulated samples with no filtering by sequence identity (id_0) are highlighted with larger
6 point sizes. (D) Boxplots of Spearman's Rho coefficients in correlation analyses between
7 taxonomic profiles of reference and simulated samples (y-axis) at different thresholds of
8 sequence identity (from 0 to 90%; x-axis) stratified by the number of species in reference
9 metagenomic profiles. Points are colored according with the sequencing depth of simulated
10 samples.

11

12 **Figure 4: Impact of filtering minimap2 primary alignments of Nanopore reads at different**
13 **thresholds of mapQ score.** Boxplots of recall (A) and precision (B) values of species richness
14 estimates in 250 simulated samples (y-axis) between metagenomic profiles inferred from
15 primary alignments of Nanopore reads filtered by different thresholds of mapQ score (from 0
16 to 50; x-axis) stratified by the number of species in reference metagenomic profiles. (C) PCoA
17 of reference metagenomic profiles and metagenomic profiles of the 250 simulated samples
18 inferred from minimap2 primary alignments filtered by different thresholds of mapQ score (from
19 0 to 50; x-axis). Dashed lines connect points coming from the same sample (reference,
20 simulated ones) with different shapes assigned to samples from different reference species
21 richness. Points corresponding to reference samples and simulated samples with no filtering
22 by mapQ (0) and filtered by mapQ>5 (5) are highlighted with larger point sizes. (D) Boxplots
23 of Spearman's Rho coefficients in correlation analyses between taxonomic profiles of
24 reference and simulated samples (y-axis) at different thresholds of sequence identity (from 0
25 to 90%; x-axis) stratified by the number of species in reference metagenomic profiles. Points
26 are colored according with the sequencing depth of simulated samples.

27

1 **Figure 5: DNA extraction kits, fragmentation and end repair impact over human**
2 **stool metagenomic composition from Nanopore sequencing data.** Read length
3 distributions of nanopore reads across different DNA extraction kits (A, n=29) and
4 between DNA fragmentation (B, n=6 paired samples Fragmented/non-fragmented)
5 and DNA end repair (C, n=6 paired samples end vs. no end repair) steps for Invitrogen
6 samples. Blue dashed lines correspond to the median value of log₂-transformed read
7 lengths used to stratify reads as long or short. (D) Differences between the fraction of
8 classified reads by Centrifuge approach between long and short reads for 29 samples
9 in panel A. (E) Differences in microbial diversity (observed species) between extraction
10 kits (n=29). (F) Differences in microbial diversity by DNA fragmentation (n=4 paired
11 samples). (G) Differences in microbial diversity by DNA end-repair step (n=4 paired
12 samples). (H) PCoA ordination of 29 samples in panel A colored by extraction kit. (I)
13 PCoA ordination of 8 samples in panel E. (J) PCoA of 8 samples of panel G colored
14 by DNA end-repair step. ns (panel F, G) =Non-significant differences in paired
15 Wilcoxon rank-sum tests. ***=Pvalue<0.0001 in paired Wilcoxon rank-sum test

16

17 **Figure 6: Optimization of DNA extraction and library preparation protocols.** (A)
18 Steps of the bacterial DNA extraction protocol. In black, the steps include in the
19 protocol of the Invitrogen kit, in red the improvement steps recommended by the IHMS
20 consortium. B) Improvement of library preparation by application of NEB
21 recommendation and decrease the SPRI/DNA ratio.

22

23 **Figure 7: Impact of Invitrogen optimized protocol over human stool**
24 **metagenomic composition from Nanopore sequencing data.** (A) Read length
25 distributions of ONT reads across different DNA extraction kits including the optimized

1 Invitrogen protocol. Blue dashed lines correspond to the median value of log₂-
2 transformed read lengths used to stratify reads as long or short. (B) Differences in the
3 fraction of classified reads between Invitrogen optimized kit and original Invitrogen kit
4 (n=6 paired samples; *=P-value<0.05, Paired Wilcoxon rank-sum test). (C) Differences
5 in microbial diversity between Invitrogen optimized kit and original Invitrogen kit (n=6
6 paired samples; ns=P-value>0.05, Paired Wilcoxon rank-sum test). (D) PCoA
7 ordination of 12 samples extracted with Invitrogen optimized kit and original Invitrogen
8 kit. Dashed lines connect samples coming from the same fecal stool sample collected
9 at different dates.

10

11 **Figure 8: Collection kits and storage conditions Impact on metagenomic human**
12 **stool composition from Nanopore sequencing data.** (A) log₂-read length
13 distribution of ONT reads across collection kits and temperature storage conditions.
14 For comparison, the log₂-read length distribution of initial Invitrogen optimized reads
15 is included. Dashed blue line represents the median log₂-read length from the entire
16 dataset. (B) Difference in the fraction of classified reads by Centrifuge strategy
17 between collection kits stratified by storage condition. (C) Differences in microbial
18 diversity between collection kits stratified by storage condition. (D) Differences in
19 microbial diversity between donors of fecal samples in this experiment. (E) Impact of
20 difference experimental variables (donor, temperature, collection kit) over microbiome
21 composition of studied samples. The barplot represents the effect sizes (R²) from
22 PERMANOVA tests of variables in Y-axis over a beta-diversity distance matrix
23 computed from Centrifuge-based genus abundance data (*=P-value<0.05,
24 PERMANOVA test). (F) PCoA ordination of samples from collection kits experiments

1 coloured by donor. Dashed lines connect samples collected with same collection kits
2 (Omnigen, Ozyme, Norgen).

3

4 **Figure 9: Comparison of quantitative metagenomic profiles of Microbaria**
5 **samples between sequencing technologies.** Correlation between gene richness
6 from SOLiD sequencing (x-axis) and Observed Species inferred from Nanopore(ONT)
7 sequencing data using Centrifuge approach (A, n=33), gene richness inferred from
8 ONT sequencing data (B, n=33) and gene richness inferred from Illumina sequencing
9 data (C, n=21). The strength of the similarities was evaluated with Spearman
10 correlation test (Spearman's Rho and P-value included in the scatter plots). (D)
11 Lineplots representing the scaled diversity (from zero to 1) of Microbaria samples from
12 different diversity metrics based on SOLiD, ONT and Illumina sequencing data.
13 Samples in x-axis are ordered based on the scaled diversity of the gene richness from
14 the original Microbaria study (GeneRichness3.9SOLiD). (E) Heatmap of Spearman's
15 Rho representing similarities in abundance vectors of taxonomic features in x-axis
16 between ONT quantifications based on Centrifuge data and Illumina and SOLiD
17 quantifications based on metagenomic species of the IGC gene catalog (y-axis; #=P-
18 valueadj<0.05, BH method; *=P-value<0.05). On the bottom of the heatmap is
19 represented the prevalence of taxonomic features in x-axis based on ONT sequencing
20 data.

21

1 **Additional files**

2 **Additional file S1:** Density distributions of mapQ scores in primary alignments of 250
3 simulated samples stratified by the number of species in reference samples (50
4 samples per reference species richness).

5 **Additional file S2: Statistical comparison of differences between reference and**
6 **simulated samples at different thresholds of mapQ scores.** At each level of
7 reference species richness (from 50 to 450 species, 50 samples per level), we
8 compare the distributions of the similarities in species abundances between reference
9 and simulated samples (Spearman's Rho coefficients of correlations between
10 reference and simulated species abundance vectors) for all possible pairs of mapQ
11 thresholds evaluated with Tukey's post-hoc pairwise tests. The 95% family-wise
12 confidence level in the difference between pairs of mapQ threshold is represented
13 colored by the significance of the difference according to adjusted P-values in Tukey's
14 tests. If we focus on the mapQ=5, we observe that higher mapQ values leads to higher
15 similarities between reference and simulated species abundance vectors (positive
16 values in the confidence levels of the differences) in R50 and R150 simulated samples,
17 whereas this is not the case for more complex/rich simulated samples (R250-R450),
18 where we observe that the similarities with the reference decrease as we increase the
19 stringency of the mapQ filtering (negative values in the confidence levels of the
20 differences, being significant for R450 samples).

21 **Additional file S3: Taxonomic profile of ZymoBIOMICS mock community inferred**
22 **from Nanopore sequencing.** The reference composition of ZymoBIOMICS mock
23 community is compared with the taxonomic profile obtained from Nanopore
24 sequencing data with Centrifuge only and with Centrifuge combined with filtering of
25 read bins by minimap2 mapping against the corresponding reference genomes with

1 parameters derived from simulation experiments (primary alignments only, min.
2 mapQ=5)

3 **Additional file S4: Comparison of microbial composition of Microbaria samples**
4 **between Nanopore, Illumina and SOLiD sequencing data.** (A) PCoA of samples
5 from Microbaria study based on genus-level MGS abundance data from three different
6 sequencing technologies (n=34 for Nanopore (ONT) and SOLiD; n=21 for Illumina).
7 Significant effect of sequencing technology in microbiome composition is observed in
8 PERMANOVA test (P-value=0.001; R2=0.11), with sample points from Illumina and
9 ONT sequencing data (both generated with Invitrogen optimized protocol) closer than
10 sample points from SOLiD sequencing (different DNA extraction method) (B)
11 Hierarchical clustering of Microbaria samples product of different sequencing methods
12 based on same genus-level MGS abundance data as PCoA in panel A. Sample points
13 from Illumina and ONT sequencing over same biological sample tends to cluster
14 together in the dendrogram.

15 **Additional file S5: Comparison of similarities in the abundance of taxonomic**
16 **features between Nanopore and SOLiD-Illumina sequencing data.** (A)
17 Correlations of taxonomic feature abundances at different levels of taxonomic
18 hierarchy between Nanopore(ONT) abundance data based on Centrifuge approach
19 and Illumina and SOLiD abundance data (based on MGS from IGC catalog). (B)
20 Correlations of taxonomic feature abundances at different levels of taxonomic
21 hierarchy between ONT abundance data and Illumina and SOLiD abundance data
22 based on MGS from IGC catalog. Dashed lines connect the same taxonomic feature
23 across comparisons. ** P-value<0.01, Paired Wilcoxon rank-sum test.

24 **Additional file S6: Comparison of similarities in KEGG functional modules**
25 **abundance between Nanopore (ONT) and SOLiD-Illumina sequencing data.** (A)

1 Vulcano plots comparing the results of Spearman correlations of individual KEGG
2 functional modules between ONT and Illumina-SOLiD sequencing data. (B)
3 Comparison of similarities in module abundance data (Spearman's Rho) between ONT
4 abundance data (from Centrifuge and from MGS abundance data) and Illumina and
5 SOLiD abundance data (based on MGS abundance data). P-values from pairwise
6 Wilcoxon rank-sum tests of Spearman's rho distributions between comparisons in x-
7 axis are shown above the violin plots.

8 **Additional file S7:** Scatterplots of KEGG Sporulation module M00485 abundance and
9 microbial diversity across different quantifications of diversity and module abundance
10 based on Nanopore (ONT), SOLiD and Illumina sequencing data. Results of Spearman
11 correlation tests are shown for each comparison.

12 **Additional file S8:** Scatterplots of abundances of KEGG LPS biosynthesis modules
13 (M00060, M00063) and microbial diversity across different quantifications of diversity
14 and module abundance based on Nanopore (ONT), SOLiD and Illumina sequencing
15 data. Results of Spearman correlation tests are shown for each comparison.

16 **Additional file S9: Summary of *Bacteroides vulgatus* and *Akkermansia***
17 ***muciniphila* genomes the assemblies from Microbaria Nanopore sequencing**
18 **data.** (A) Mean \pm standard error of the top 26 bacterial species with the highest
19 abundance in 33 Microbaria samples based on Centrifuge workflow. (B) MUMMER
20 alignment dot plot between *Bacteroides vulgatus* assembly from Nanopore (ONT)
21 reads (y axis) and the reference genome of *Bacteroides vulgatus* ATCC 8482 (x axis).
22 Each point represents a contig in y-axis matching the reference genome on x-axis, with
23 red representing the same orientation and blue representing inverse orientation of
24 contig in y axis vs. reference genome. Points across the diagonal represents
25 assembled contigs collinear with the reference genome on x-axis. (C) MUMMER

1 alignment dot plot between *Akkermansia muciniphila* assembly from Nanopore (ONT)
2 reads (y axis) and the reference genome of *Akkermansia muciniphila* isolate *Urmite* (x
3 axis). (D) Comparison of log₂-gene length distributions between genes from
4 *Bacteroides vulgatus* Nanopore (ONT) assembly and genes from *Bacteroides vulgatus*
5 *ATCC 8482* reference genome. (E) Comparison of log₂-gene length distributions
6 between genes from *Akkermansia muciniphila* Nanopore (ONT) assembly and genes
7 from *Akkermansia muciniphila* isolate *Urmite* reference genome.

Figures

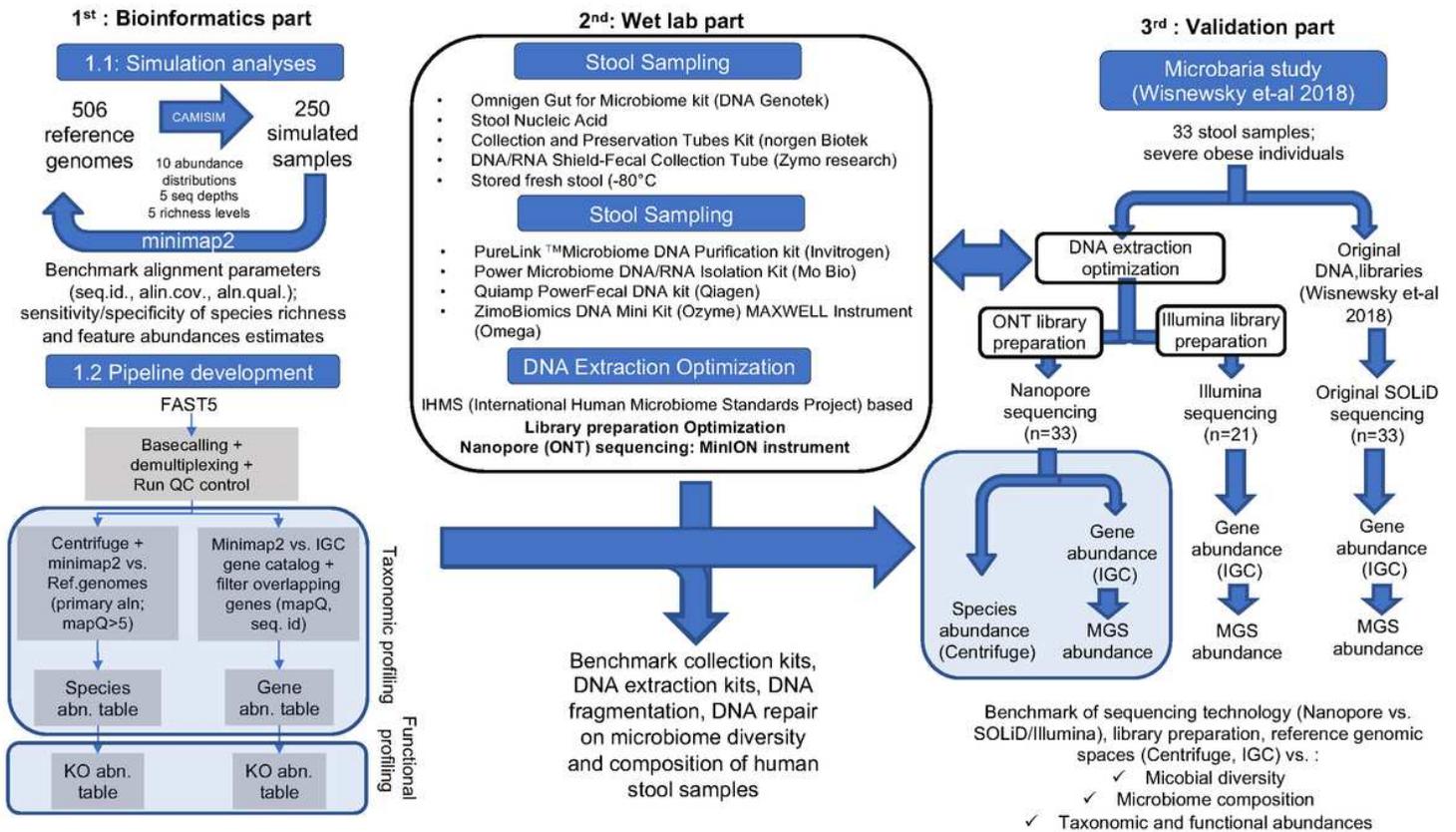


Figure 1

Summary of the workflow : (A) Simulated data processing. (B) Wet-lab optimization. (C) Summary of ONT sequencing comparison with Illumina and SOLiD technologies.

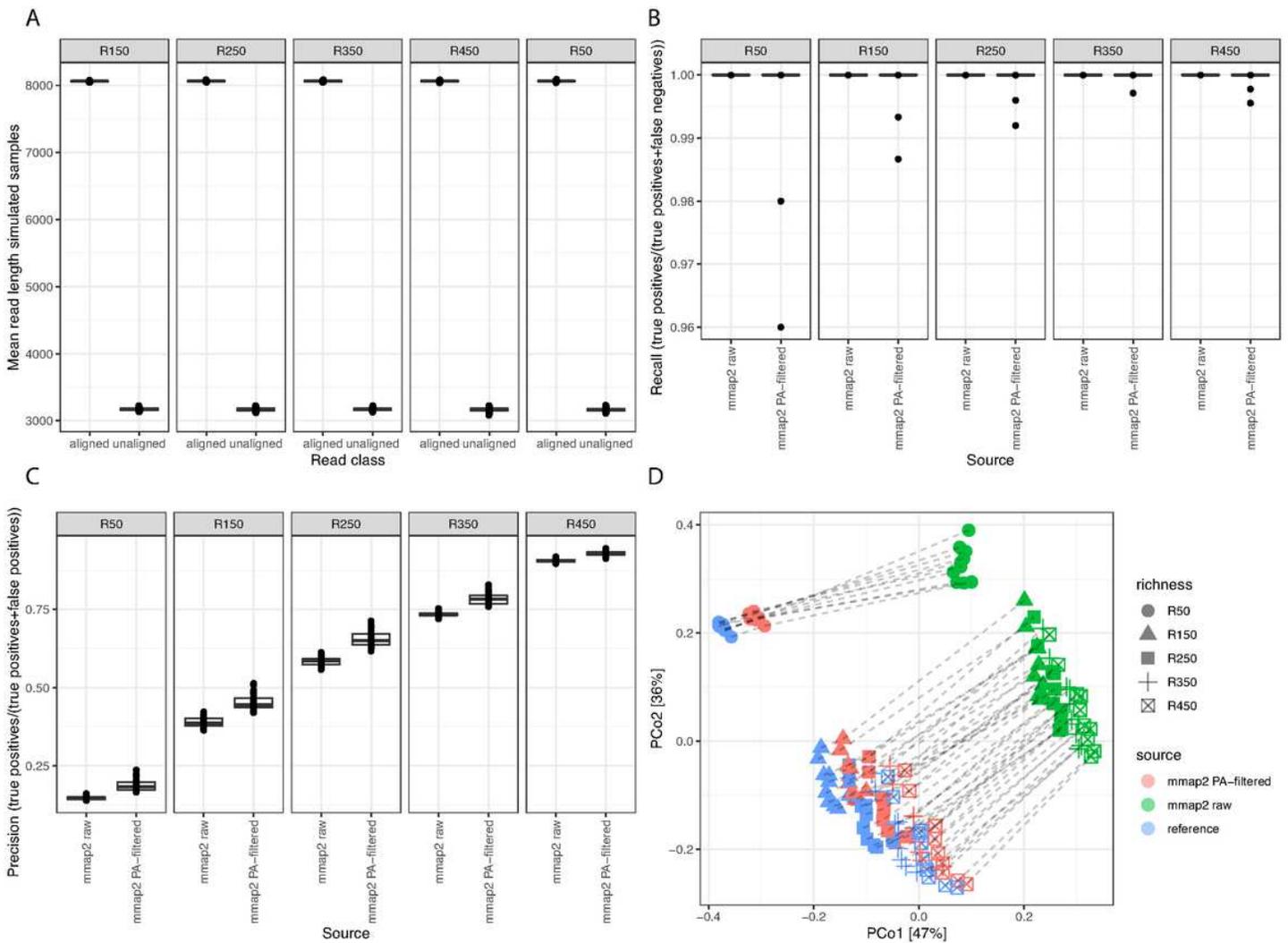


Figure 2

Metagenomic profiles from simulated samples between minimap2 results and minimap2 results filtered from secondary alignments. (A) Boxplots of mean lengths of Nanopore reads of 250 simulated samples (y axis) between those aligned and unaligned over the 506 reference genomes from minimap2 results (x-axis). (B) Boxplots of recall values of species richness estimates in 250 simulated samples (y-axis) between metagenomic profiles inferred from all minimap2 alignments (mmap2 raw) and from minimap2 primary alignments only (mmap2APfilt, x-axis). (C) Boxplots of precision values of species richness estimates in 250 simulated samples (y-axis) between metagenomic profiles inferred from all minimap2 alignments (mmap2 raw) and from minimap2 primary alignments only (mmap2APfilt, x-axis). (D) Principal Coordinates Analysis (PCoA) of metagenomic profiles from the reference and 250 simulated samples inferred from all minimap2 alignments (mmap2raw) and from minimap2 primary alignments only (mmap2APfilt, x-axis). Dashed lines connect points coming from the same sample (reference, simulated ones; 3 points per sample).

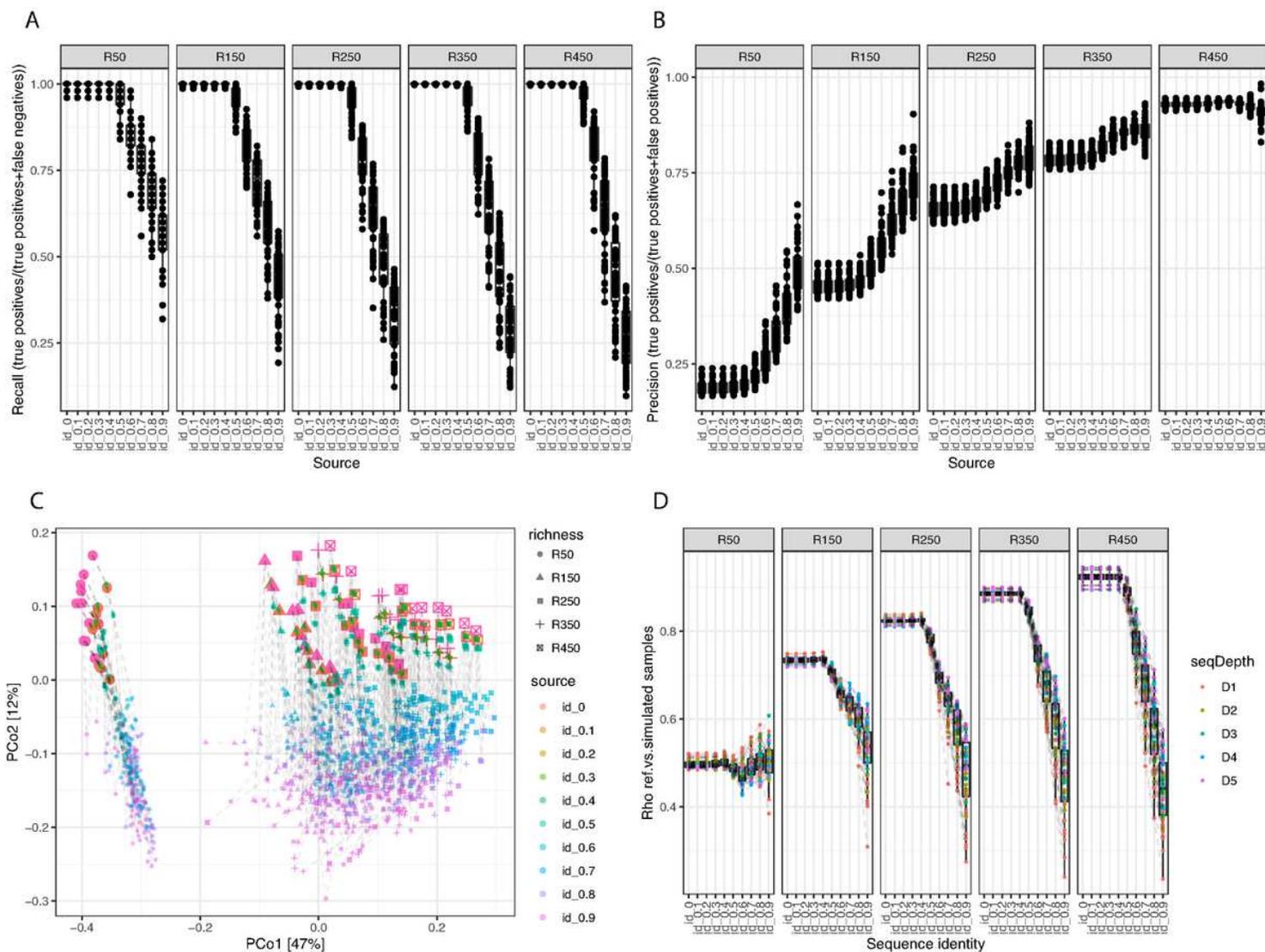


Figure 3

Impact of filtering minimap2 primary alignments of Nanopore reads at different thresholds of sequence identity. Boxplots of recall (A) and precision (B) values of species richness estimates in 250 simulated samples (y-axis) between metagenomic profiles inferred from primary alignments of Nanopore reads filtered by different thresholds of sequence identity (from 0 to 90%; x-axis) stratified by the number of species in reference metagenomic profiles. (C) PCoA of reference metagenomic profiles and metagenomic profiles of the 250 simulated samples inferred from minimap2 primary alignments filtered by different thresholds of sequence identity (from 0 to 90%; x-axis). Dashed lines connect points coming from the same sample (reference, simulated ones; 11 points per sample) with different shapes assigned to samples from different reference species richness. Points corresponding to reference samples and simulated samples with no filtering by sequence identity (id_0) are highlighted with larger point sizes. (D) Boxplots of Spearman's Rho coefficients in correlation analyses between taxonomic profiles of reference and simulated samples (y-axis) at different thresholds of sequence identity (from 0 to 90%; x-axis) stratified by the number of species in reference metagenomic profiles. Points are colored according with the sequencing depth of simulated samples.

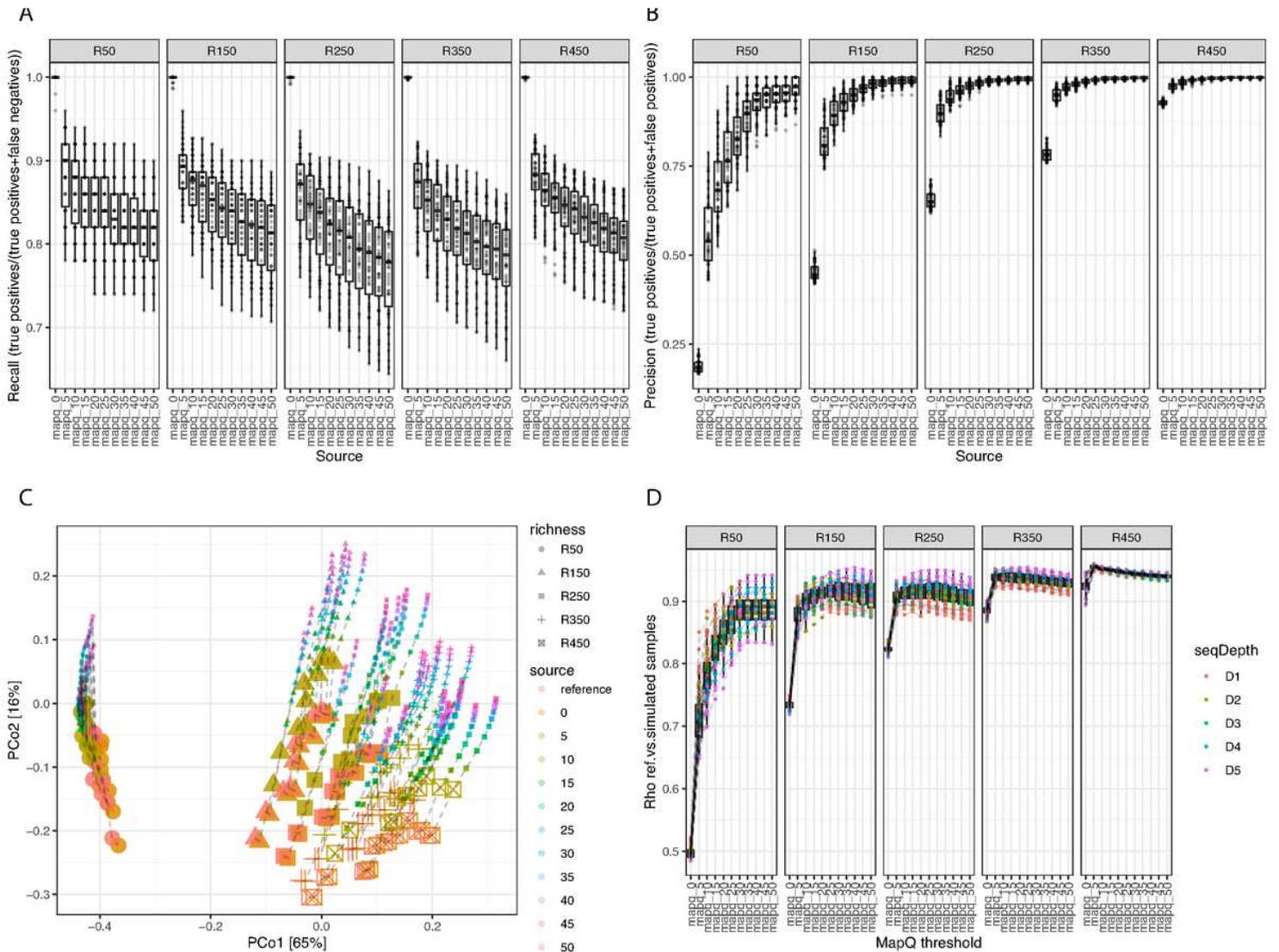


Figure 4

Impact of filtering minimap2 primary alignments of Nanopore reads at different thresholds of mapQ score. Boxplots of recall (A) and precision (B) values of species richness estimates in 250 simulated samples (y-axis) between metagenomic profiles inferred from primary alignments of Nanopore reads filtered by different thresholds of mapQ score (from 0 to 50; x-axis) stratified by the number of species in reference metagenomic profiles. (C) PCoA of reference metagenomic profiles and metagenomic profiles of the 250 simulated samples inferred from minimap2 primary alignments filtered by different thresholds of mapQ score (from 0 to 50; x-axis). Dashed lines connect points coming from the same sample (reference, simulated ones) with different shapes assigned to samples from different reference species richness. Points corresponding to reference samples and simulated samples with no filtering by mapQ (0) and filtered by mapQ>5 (5) are highlighted with larger point sizes. (D) Boxplots of Spearman's Rho coefficients in correlation analyses between taxonomic profiles of reference and simulated samples (y-axis) at different thresholds of sequence identity (from 0 to 90%; x-axis) stratified by the number of species in reference metagenomic profiles. Points are colored according with the sequencing depth of simulated samples.

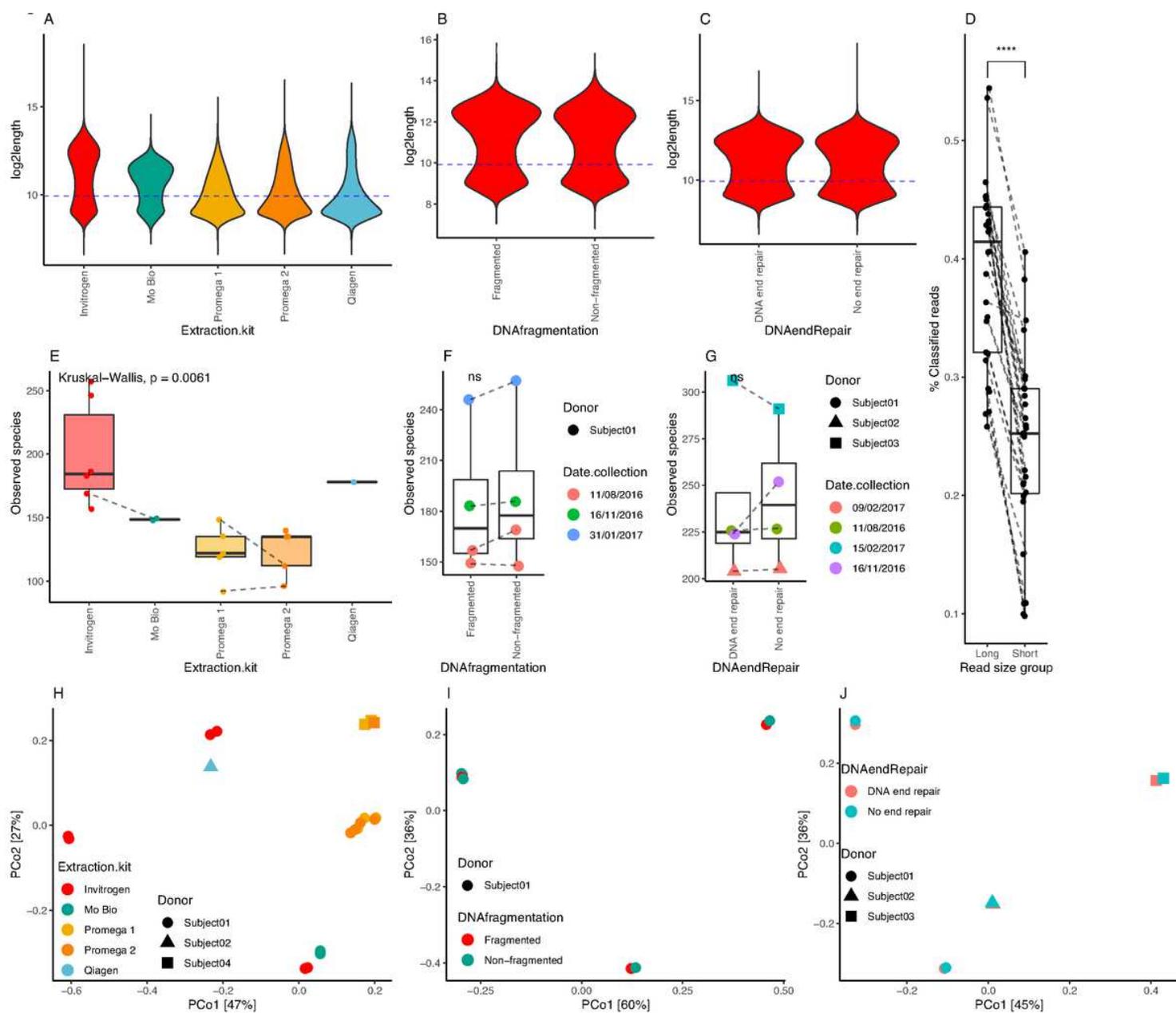


Figure 5

DNA extraction kits, fragmentation and end repair impact over human stool metagenomic composition from Nanopore sequencing data. Read length distributions of nanopore reads across different DNA extraction kits (A, $n=29$) and between DNA fragmentation (B, $n=6$ paired samples Fragmented/non-fragmented) and DNA end repair (C, $n=6$ paired samples end vs. no end repair) steps for Invitrogen samples. Blue dashed lines correspond to the median value of \log_2 -transformed read lengths used to stratify reads as long or short. (D) Differences between the fraction of classified reads by Centrifuge approach between long and short reads for 29 samples in panel A. (E) Differences in microbial diversity (observed species) between extraction kits ($n=29$). (F) Differences in microbial diversity by DNA fragmentation ($n=4$ paired samples). (G) Differences in microbial diversity by DNA end-repair step ($n=4$ paired samples). (H) PCoA ordination of 29 samples in panel A colored by extraction kit. (I) PCoA ordination of 8 samples in panel E. (J) PCoA of 8 samples of panel G colored by DNA end-repair step. ns

(panel F, G) = Non-significant differences in paired Wilcoxon rank-sum tests. ***=Pvalue<0.0001 in paired Wilcoxon rank-sum test

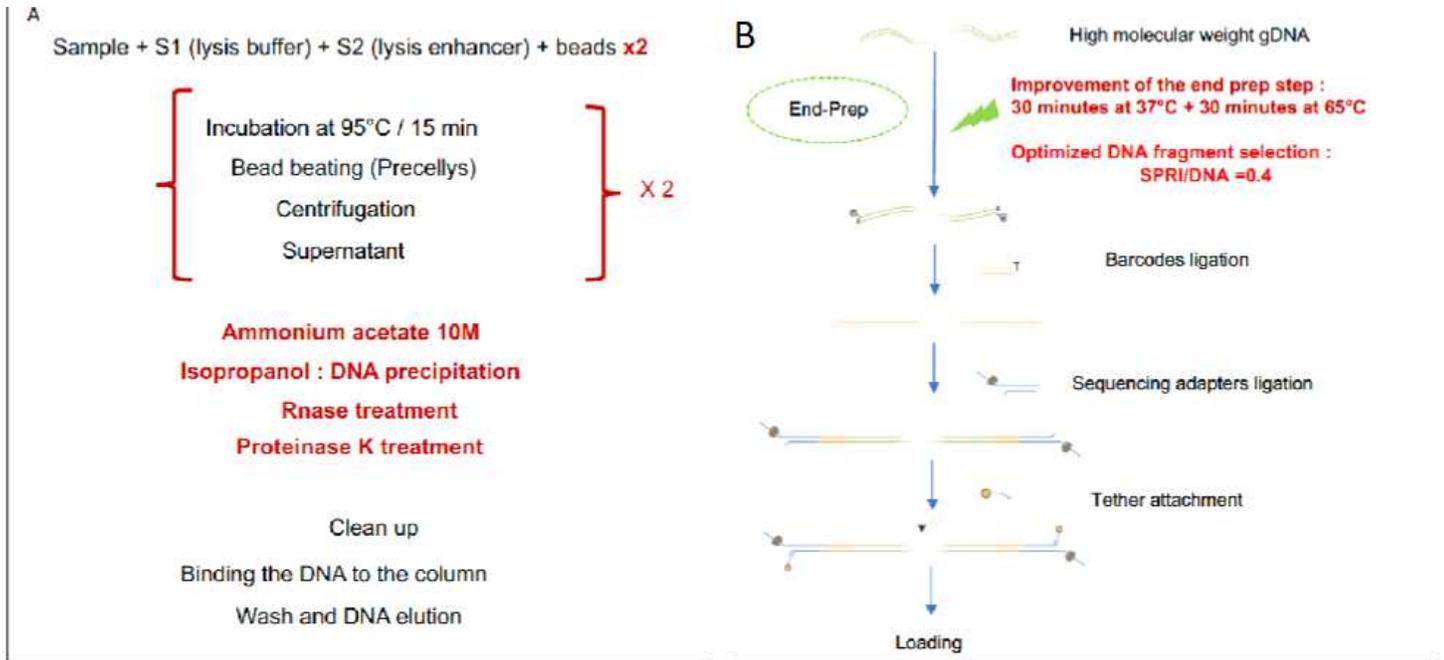


Figure 6

Optimization of DNA extraction and library preparation protocols. (A) Steps of the bacterial DNA extraction protocol. In black, the steps include in the protocol of the Invitrogen kit, in red the improvement steps recommended by the IHMS consortium. B) Improvement of library preparation by application of NEB recommendation and decrease the SPRI/DNA ratio.

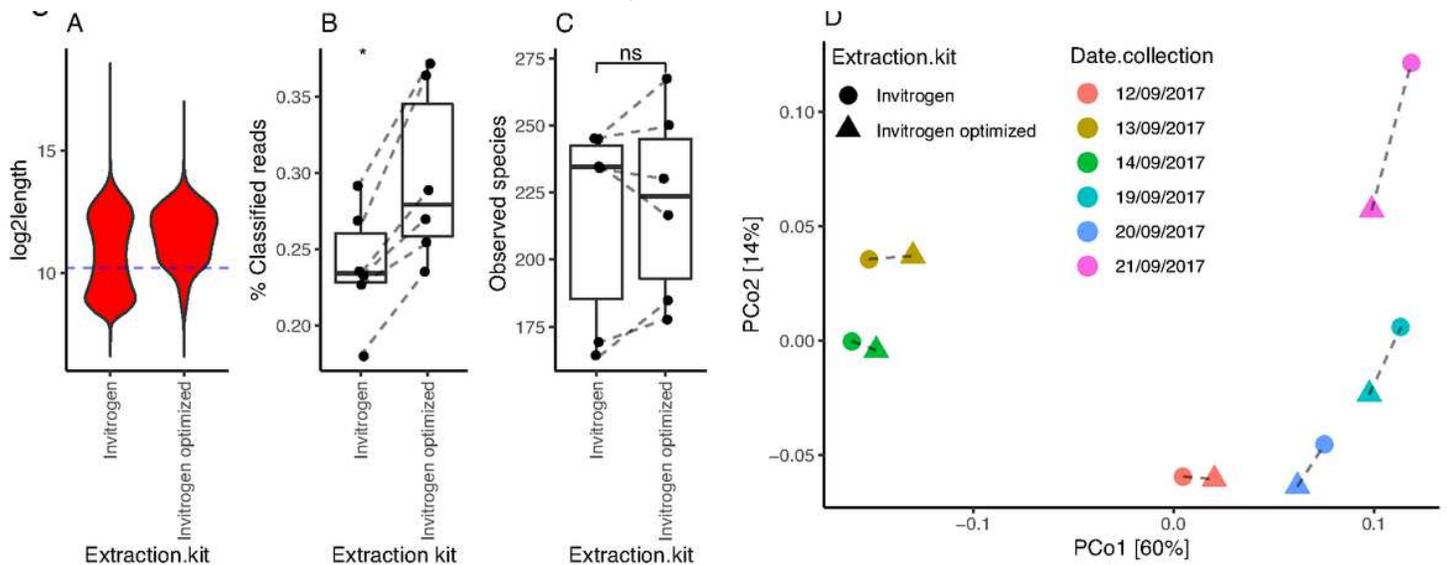


Figure 7

Impact of Invitrogen optimized protocol over human stool metagenomic composition from Nanopore sequencing data. (A) Read length distributions of ONT reads across different DNA extraction kits

including the optimized Invitrogen protocol. Blue dashed lines correspond to the median value of log₂-transformed read lengths used to stratify reads as long or short. (B) Differences in the fraction of classified reads between Invitrogen optimized kit and original Invitrogen kit (n=6 paired samples; * = P-value < 0.05, Paired Wilcoxon rank-sum test). (C) Differences in microbial diversity between Invitrogen optimized kit and original Invitrogen kit (n=6 paired samples; ns = P-value > 0.05, Paired Wilcoxon rank-sum test). (D) PCoA ordination of 12 samples extracted with Invitrogen optimized kit and original Invitrogen kit. Dashed lines connect samples coming from the same fecal stool sample collected at different dates.

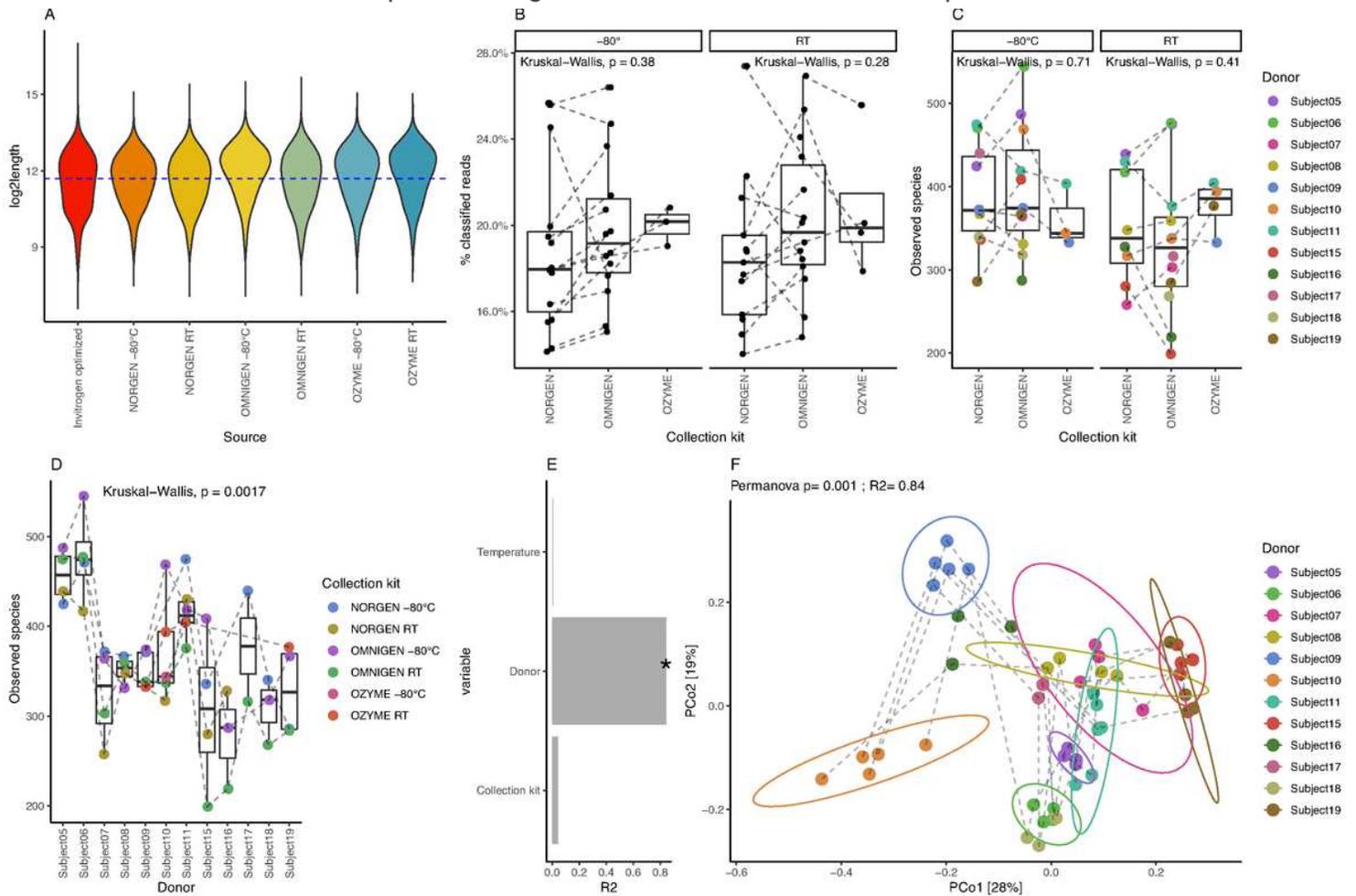


Figure 8

Collection kits and storage conditions Impact on metagenomic human stool composition from Nanopore sequencing data. (A) log₂-read length distribution of ONT reads across collection kits and temperature storage conditions. For comparison, the log₂-read length distribution of initial Invitrogen optimized reads is included. Dashed blue line represents the median log₂-read length from the entire dataset. (B) Difference in the fraction of classified reads by Centrifuge strategy between collection kits stratified by storage condition. (C) Differences in microbial diversity between collection kits stratified by storage condition. (D) Differences in microbial diversity between donors of fecal samples in this experiment. (E) Impact of difference experimental variables (donor, temperature, collection kit) over microbiome composition of studied samples. The barplot represents the effect sizes (R²) from PERMANOVA tests of variables in Y-axis over a beta-diversity distance matrix computed from Centrifuge-based genus

abundance data (*=P-value<0.05, PERMANOVA test). (F) PCoA ordination of samples from collection kits experiments coloured by donor. Dashed lines connect samples collected with same collection kits (Omnigen, Ozyme, Norgen).

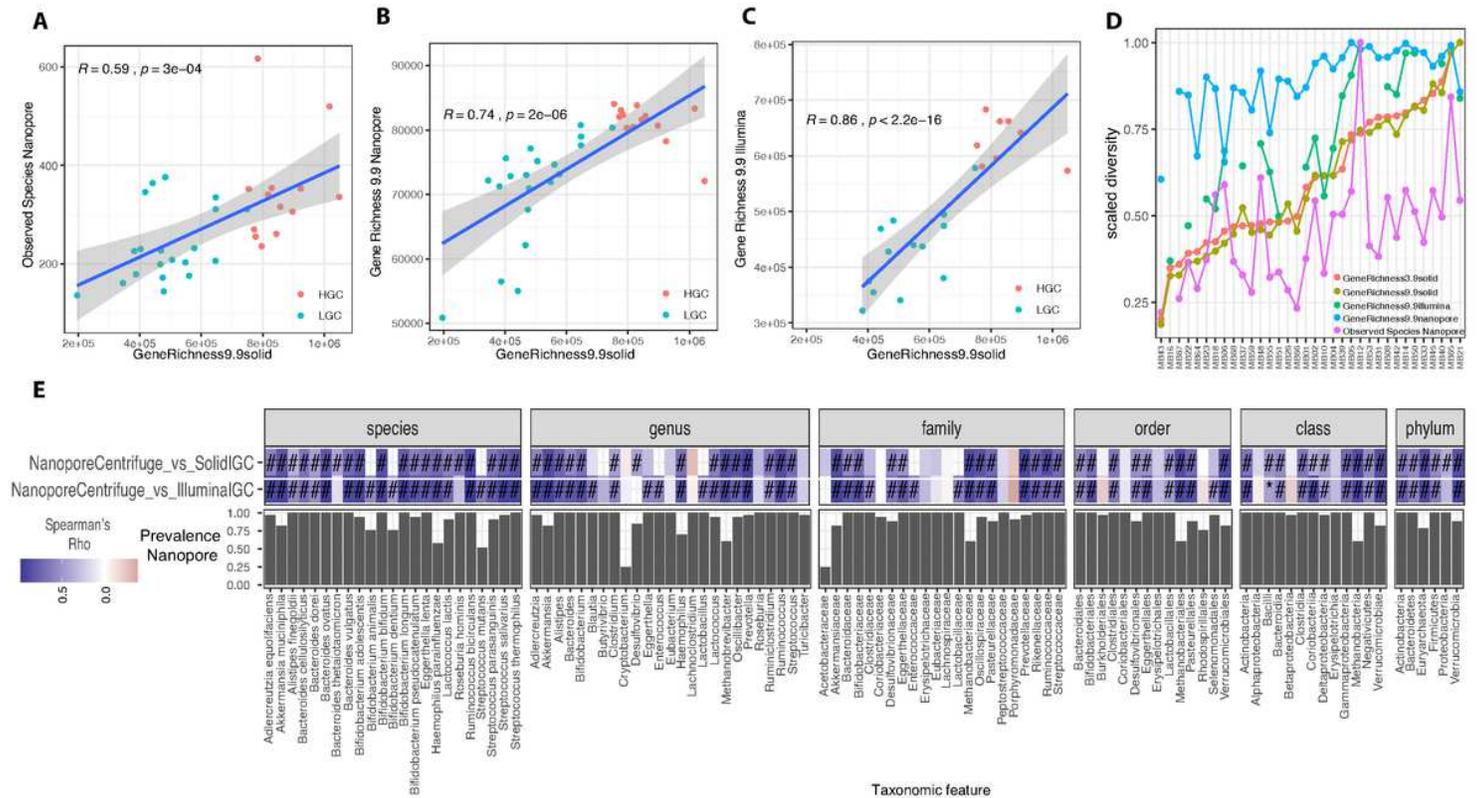


Figure 9

Comparison of quantitative metagenomic profiles of Microbaria samples between sequencing technologies. Correlation between gene richness from SOLiD sequencing (x-axis) and Observed Species inferred from Nanopore(ONT) sequencing data using Centrifuge approach (A, n=33), gene richness inferred from ONT sequencing data (B, n=33) and gene richness inferred from Illumina sequencing data (C, n=21). The strength of the similarities was evaluated with Spearman correlation test (Spearman's Rho and P-value included in the scatter plots). (D) Lineplots representing the scaled diversity (from zero to 1) of Microbaria samples from different diversity metrics based on SOLiD, ONT and Illumina sequencing data. Samples in x-axis are ordered based on the scaled diversity of the gene richness from the original Microbaria study (GeneRichness3.9SOLiD). (E) Heatmap of Spearman's Rho representing similarities in abundance vectors of taxonomic features in x-axis between ONT quantifications based on Centrifuge data and Illumina and SOLiD quantifications based on metagenomic species of the IGC gene catalog (y-axis; #=P-valueadj<0.05, BH method; *=P-value<0.05). On the bottom of the heatmap is represented the prevalence of taxonomic features in x-axis based on ONT sequencing data.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [SupplFigure1.pdf](#)
- [SupplFigure2.pdf](#)
- [SupplFigure3.pdf](#)
- [SupplFigure4.pdf](#)
- [SupplFigure5.pdf](#)
- [SupplFigure6.pdf](#)
- [SupplFigure7.pdf](#)
- [SupplFigure8.pdf](#)
- [SupplFigure9.pdf](#)
- [SupplementalTables.xlsx](#)