

1 **Title: Benchmarking Short-Read ITS2 and Full-Length ITS Sequencing Reveals Pipeline-**  
2 **Dependent Biases in Indoor Fungal Community Profiling**

3 **Authors:** Mengyi Dong<sup>1,2,3</sup>, Denene Blackwood<sup>4</sup>, Megan Lott<sup>5</sup>, Sherlynette Pérez Castro<sup>4</sup>,  
4 Xavier Larkin<sup>4</sup>, Thomas J. Clerkin<sup>4</sup>, Heather Hemric<sup>1</sup>, Jake Nash<sup>6,7</sup>, Yeon Ji Kim<sup>1</sup>, Jason W.  
5 Arnold<sup>1</sup>, Lawrence A. David<sup>1</sup>, Rytas J. Vilgalys<sup>6</sup>, Anthony A. Fodor<sup>8</sup>, Rachel T. Noble<sup>4,5</sup>

6

7 **Corresponding authors:** Mengyi Dong; Rachel T. Noble

8 **Email:** [mengyidong@vt.edu](mailto:mengyidong@vt.edu); [rtnoble@email.unc.edu](mailto:rtnoble@email.unc.edu)

9

10 **Author Contributions:** R.N. and M.D. conceptualized the study. M.D., R.N., A.A.F., J.N.,  
11 R.J.V., J.W.A., Y.J.K. generated the methodology. D.B., M.L., X.L., T.C., S.P.C., H.H., and  
12 M.D. performed the investigation. M.D. conducted formal analysis. M.D., M.L., and S.P.C.  
13 wrote the original draft of the manuscript. M.D. performed the visualization. R.N., L.D., and  
14 M.D. acquired funding. J.N., R.J.V., A.A.F., J.W.A., and R.N. provided resources. R.N.  
15 supervised the work. All co-authors reviewed and edited the manuscript.

16

17 **Conflict Interest Statement:** The authors declare no competing interests.

18

19 **Keywords:** indoor mycobiome; ITS sequencing; PacBio HiFi; amplicon sequencing; fungal  
20 community profiling; taxonomic resolution; built environment; sequencing pipeline comparison;  
21 health-relevant fungi; ASV

22

---

**Affiliations:**

<sup>1</sup> Department of Molecular Genetics and Microbiology, Duke Microbiome Center, Duke University School of Medicine, Durham, North Carolina, USA

<sup>2</sup> Present address: Virginia Seafood Agricultural Research and Extension Center, Virginia Tech, Hampton, Virginia, USA

<sup>3</sup> Present address: Department of Food Science and Technology, Virginia Tech, Blacksburg, Virginia, USA

<sup>4</sup> Department of Earth, Marine, and Environmental Sciences, Institute of Marine Sciences (IMS), University of North Carolina at Chapel Hill, Morehead City, North Carolina, USA

<sup>5</sup> Department of Environmental Sciences and Engineering, the Gillings School of Global Public Health, University of North Carolina Chapel-Hill, Chapel Hill, North Carolina, USA

<sup>6</sup> Department of Biology, Duke University, Durham, North Carolina, USA

<sup>7</sup> Present address: Department of Biology, Boston University, Boston, Massachusetts, USA

<sup>8</sup> Department of Bioinformatics and Genomics, University of North Carolina at Charlotte, Charlotte, North Carolina, USA

## 23 Abstract

24 Short-read amplicon sequencing is widely used for fungal surveys but can limit taxonomic  
25 resolution. Long-read sequencing enables recovery of the full internal transcribed spacer (ITS)  
26 region and may improve ecological and taxonomic inference. Here, we conducted a paired  
27 comparison of Illumina ITS2 and PacBio HiFi full-length ITS sequencing using identical DNA  
28 extracts from built-environmental air and surface samples ( $n = 68$ ) collected across homes, a  
29 dormitory, and laboratories. Both datasets were taxonomically assigned using the same algorithm  
30 and reference database. We performed paired statistics, in-silico ITS2 trimming of long-read  
31 sequences, and cross-platform mapping at multiple identity thresholds. Full-length ITS provided  
32 higher taxonomic resolution, assigning a greater fraction of ASVs at the family (98% vs. 88%)  
33 and species (42% vs. 32%) ranks than ITS2 (paired Wilcoxon  $q = 0.002$ ). Alpha-diversity  
34 comparisons showed similar Shannon diversity across pipelines, whereas richness metrics were  
35 consistently higher for full-length ITS. Beta-diversity analyses indicated broadly comparable  
36 community-level patterns, although full-length ITS revealed stronger sample-type- and  
37 location-associated structure (PERMANOVA  $R^2 \geq 0.06$ ,  $p = 0.0001$ ). In-silico ITS2 trimming  
38 reduced these differences, indicating that amplicon length is a major contributor to enhanced  
39 taxonomic resolution and ecological inference. Cross-platform mapping further showed  
40 extensive one-to-many relationships between ITS2 and full-length ITS ASVs, consistent with  
41 increased sequence resolution in long-read data.

42 Together, these results show that ITS2 sequencing provides robust community-level profiling,  
43 while full-length ITS enables improved richness estimates and finer ecological and taxonomic  
44 resolution. This paired, bias-aware framework provides a practical template for selecting fungal  
45 amplicon sequencing strategies in built-environment mycobiome studies.

## 46 Importance

47 Fungal communities in built environments influence indoor air quality and human exposure, yet  
48 their characterization depends strongly on sequencing strategy. This study provides a controlled,  
49 paired comparison of short-read ITS2 and long-read full-length ITS sequencing, showing that  
50 differences in amplicon length substantially contribute to variation in taxonomic resolution and  
51 ecological inference. While both approaches yield comparable community-level patterns,  
52 full-length ITS improves richness estimates, species-level assignment, and environmental  
53 discrimination by resolving sequence variation collapsed in ITS2 surveys. By integrating paired  
54 diversity analyses, in-silico ITS2 trimming, and cross-platform ASV mapping, this work offers a  
55 bias-aware framework for evaluating fungal amplicon pipelines. Importantly, improved  
56 species-level resolution enables functional interpretation of indoor fungi, for example the  
57 identification of taxa associated with pathogenic traits, allergen production, or toxin synthesis,  
58 supporting the development of more informative exposure metrics and targeted assays relevant to  
59 human health in built environments.

## 60 Introduction

61 The built-environment microbiome influences indoor air quality, building integrity, and human  
62 health (1). While indoor bacterial communities have been extensively characterized (1, 2),  
63 growing attention has focused on fungal diversity, spatial variation, and exposure pathways in  
64 indoor environments (3–5). These studies show that indoor fungi originate from multiple sources  
65 and vary widely across buildings, yet accurate detection and classification of indoor fungal taxa  
66 remain challenging (1, 6).

67 A key challenge is the reliable identification of fungi at finer taxonomic ranks, especially for  
68 genera and species relevant to health. Indoor fungi span common saprophytes and opportunistic  
69 pathogens, with detection influenced by sampling approach, primer design, and sequencing  
70 platform (6, 7). Importantly, clinically relevant differences in pathogenicity and drug  
71 susceptibility are often observed among closely related fungal species or strains that are  
72 indistinguishable using coarse taxonomic markers. Well-documented examples include cryptic  
73 *Aspergillus* species within the *A. fumigatus* complex and members of the *Cryptococcus*  
74 *neoformans*/*C. gattii* species complexes (8). Indoor fungal exposure has been linked to  
75 respiratory health effects, particularly in children (9). Sequencing surveys routinely detect  
76 hundreds of fungal taxa within single homes, including medically important species such as  
77 *Aspergillus fumigatus*, *Stachybotrys chartarum*, and *Cryptococcus neoformans* (10–14).  
78 Sensitive and accurate molecular tools are therefore needed to interpret exposure risks and  
79 support environmental assessment (9).

80 Indoor fungal community composition varies across building types, ventilation, moisture, and  
81 occupant activity (15, 16). These factors shape both dominant indoor genera and  
82 lower-abundance taxa of potential concern. Such variability highlights the need for reproducible  
83 and comparable indoor fungal datasets (7).

84 Most environmental molecular fungal surveys rely on amplification of the internal transcribed  
85 spacer (ITS) region, with Illumina ITS1 or ITS2 sequencing widely used due to cost and  
86 throughput (17). ITS2 sequencing supports broad community profiling but often limits  
87 species-level resolution for closely related taxa (18). Long-read high-fidelity platforms such as  
88 PacBio HiFi can generate full-length ITS reads that retain additional phylogenetic information  
89 (19). Previous studies report improved taxonomic resolution with long-read ITS sequencing,  
90 alongside challenges related to primer design and performance across taxa (19–21). Recent  
91 improvements to PacBio sequencing throughput coupled with lower costs have made full-length  
92 ITS profiling more feasible for larger studies (22).

93 Despite the availability of both short- and long-read approaches, most indoor mycobiome studies  
94 rely on a single amplicon region (ITS1 or ITS2) (3, 5, 7). Few paired studies directly compare  
95 ITS2 and full-length ITS using the same indoor samples, and the relative contributions of  
96 amplicon length, primer choice, sequencing platform, and taxonomic classification remain  
97 unclear.

98 Here, we perform a paired comparison of Illumina ITS2 and PacBio HiFi full-length ITS  
99 sequencing to evaluate method-dependent differences in indoor fungal community profiling.

100 Specifically, we compared taxonomic assignment across ranks under classifier parity, assessed  
101 agreement in community structure and genus-level abundance, evaluated patterns among  
102 health-relevant taxa, and used in-silico ITS2 trimming and cross-platform mapping to isolate the  
103 influence of amplicon length and primer targeting. Together, this work provides a  
104 method-focused benchmark for selecting fungal amplicon strategies in built-environment studies.

## 105 **Materials and Methods**

### 106 **Sample collection and preprocessing**

107 Indoor air and surface samples (n = 68) were collected between 2023 and 2024 from multiple  
108 built environments in North Carolina, including two unoccupied residences, a student dormitory,  
109 and laboratory spaces at the University of North Carolina at Chapel Hill and Duke University  
110 (Supplementary Data 1, Fig. S1). Sampling locations represented a range of occupancy levels,  
111 moisture conditions, and visible fungal growth. Surface swabs and bioaerosol samples were  
112 collected using standardized protocols appropriate to each site.

113 Surface samples were collected using sterile rayon swabs pre-moistened with buffer and  
114 swabbed over standardized surface areas. Bioaerosol samples were collected using a BobCat  
115 (AC-200) sampler with dry electret filters and eluted following manufacturer recommendations.  
116 Field blanks, method blanks, negative controls, and positive controls were included at each site  
117 to monitor contamination and processing consistency. All eluates were aliquoted and stored at  
118  $-80^{\circ}\text{C}$  prior to DNA extraction. Detailed sampling locations, surface types, and control  
119 descriptions are provided in Supplementary Methods.

### 120 **DNA extraction and quality control**

121 Total nucleic acids were extracted using magnetic-bead-based extraction on a KingFisher™ Flex  
122 system (bioMérieux NucliSENS kit). Extractions were eluted in 100  $\mu\text{L}$  of elution buffer.  
123 Negative extraction controls and positive controls spiked with *Aspergillus niger* were included to  
124 assess contamination and extraction efficiency. All extracts were stored at  $-80^{\circ}\text{C}$  until library  
125 preparation. Extraction scripts and additional quality-control details are provided in the  
126 Supplementary Methods.

### 127 **Illumina ITS2 Library Preparation and Sequencing**

128 The fungal ITS2 region was amplified using a three-step PCR protocol adapted from prior work  
129 (23) with ITS3NGS and ITS4NGR primers (24–26). All primer sequences are listed in S. Table  
130 2. Amplicon sizes ( $\sim 400$  bp) were verified by agarose gel electrophoresis and purified using  
131 AMPure XP Beads (Beckman Coulter, Inc., Brea, CA, USA) at a  $1.8\times$  bead-to-sample volume  
132 ratio. DNA concentrations were quantified using the Qubit HS dsDNA kit, and samples were  
133 pooled in equimolar amounts to  $\sim 2\text{--}3$  ng/ $\mu\text{L}$ . Final libraries were sequenced on an Illumina  
134 MiniSeq platform (Illumina, Inc., San Diego, CA, USA) targeting a read depth of at least 20,000  
135 reads per sample using  $2\times 250$  bp paired-end reads in the Duke Sequencing and Genomic  
136 Technologies Shared Resource (Duke University, Durham, NC, USA). Detailed procedure and  
137 PCR cycling conditions are described in Supplementary Methods.

## 138 **PacBio Full-Length ITS Library Preparation and Sequencing**

139 Full-length ITS amplicons were generated from the same DNA extracts using Phusion Plus  
140 polymerase with ITS1catta and ITS4ngsUni primers (20, 27). Samples were barcoded, pooled,  
141 and prepared following the PacBio Kinnex protocol without modification (28). Size selected and  
142 cleaned libraries were loaded onto a PacBio SMRT® Cell and sequenced on the Revio system  
143 (Pacific Biosciences of California, Inc.) at the Duke Sequencing and Genomic Technologies  
144 Shared Resource. Amplicon verification, pooling volumes, and circularization steps are  
145 described in the Supplementary Methods.

## 146 **Illumina Data Processing**

147 Illumina paired-end reads were merged using PEAR (29), and ITS2 regions were extracted using  
148 ITSxpress (30). Sequences were processed in QIIME 2 (v2023.9) (31) using DADA2  
149 denoise-single for error correction, chimera removal, and ASV inference (32). Taxonomy was  
150 assigned using a naïve Bayes classifier trained on full-length ITS in UNITE v9.0 (33). QIIME 2  
151 artifacts were exported for downstream analysis in R.

## 152 **PacBio Data Processing**

153 PacBio circular consensus sequences were processed using DADA2 v1.28 with PacBio-specific  
154 error modeling (32). Primers were removed prior to filtering, and ASVs were inferred with  
155 chimera removal optimized for long amplicons. Taxonomy was assigned using the same UNITE  
156 v9.0 classifier applied to the ITS2 data, ensuring classifier parity.

## 157 **Mapping of Illumina ITS2 ASVs Against PacBio Full-Length ITS References**

158 Illumina ITS2 ASVs were mapped to PacBio full-length ITS ASVs using the official NCBI  
159 BLAST + Docker image (v2.11.0; <https://hub.docker.com/r/ncbi/blast>) at >97% identity to assess  
160 sequence overlap and resolution differences. Detailed parameters are provided in the  
161 Supplementary Methods. Maximum likelihood trees were constructed for the BLAST mapped  
162 sequences. For downstream mapping analysis, identity thresholds of  $\geq 99\%$  or 100% were used to  
163 quantify one-to-one and one-to-many relationships between pipelines.

## 164 **In-silico ITS2 Extraction from PacBio Full-Length ITS Reads**

165 An in-silico ITS2 dataset was generated from PacBio full-length ITS ASVs using ITSx (34).  
166 Extracted ITS2 sequences were assigned using the same UNITE v9.0 classifier applied to the  
167 full-length ITS and ITS2 data. This dataset was used for sensitivity analyses of taxonomic  
168 assignment and community structure.

## 169 **Data Analysis**

170 All statistical analyses and data visualizations were performed in R Studio (v2024.04.2). ITS2  
171 sequencing depth sufficiency was evaluated using rarefaction curves and Good's coverage ( $1 -$   
172 singletons/total reads). To account for compositionality and differences in sequencing depth,

173 ASV count data were centered log-ratio (CLR) transformed using the `microbiome::transform`  
174 function (35). Taxa were agglomerated to higher taxonomic ranks using the `tax_glom` function in  
175 `phyloseq` (36) as needed. Alpha-diversity metrics, including Shannon diversity and Chao1  
176 richness, were calculated using `vegan` (37). Beta diversity was assessed using Bray-Curtis  
177 dissimilarity on relative abundance data and visualized with principal coordinate analysis  
178 (PCoA). Differences in community composition were tested using permutational multivariate  
179 analysis of variance (PERMANOVA) via `adonis2` with 10,000 permutations, using Euclidean  
180 distances on CLR-transformed data and accounting for paired sample structure where applicable.  
181 Ordination concordance between datasets was evaluated using symmetric Procrustes analysis  
182 (PROTEST) and Mantel tests (Spearman) comparing Bray-Curtis distance matrices.  
183 Taxonomic assignment success between Illumina ITS2 and PacBio full-length ITS pipelines was  
184 quantified as the percentage of ASVs assigned at each taxonomic rank per paired sample.  
185 Genus-level abundance agreement between pipelines was assessed using Spearman correlations  
186 on CLR-transformed data. Differences in abundance of health-relevant genera between pipelines  
187 and sample types were tested using paired Wilcoxon tests. Differences in per-sample correlation  
188 coefficients across sample types were evaluated using analysis of variance (ANOVA) followed  
189 by Tukey's HSD where appropriate. Assumptions of normality and homogeneity were assessed  
190 using Shapiro-Wilk and Levene's tests.

191

## 192 **Results**

### 193 *Comparison of ITS2 and Full-Length ITS Sequencing Pipeline Performance*

194 PacBio full-length ITS sequencing yielded substantially higher post-QC read counts and ASV  
195 richness than ITS2 (mean =  $739,601 \pm 1,565,353$  reads per sample; 6,419 ASVs), whereas  
196 Illumina ITS2 produced  $10,395 \pm 6,817$  reads per sample and 1,814 ASVs. Illumina ITS2  
197 libraries showed sufficient sequencing depth for downstream analyses, as indicated by  
198 rarefaction curves that approached saturation at  $\sim 2k$  reads for nearly all samples (Fig. S2).  
199 Good's coverage values were similarly high (median = 1.00; IQR = 1.00-1.00; mean = 0.97),  
200 confirming that most ITS2 samples were well covered. Three samples had low coverage ( $< 0.90$ )  
201 and were removed from paired ITS2-full-length ITS comparisons, taxonomic assignment  
202 analyses, and diversity tests; all descriptive summaries retained the full dataset. Sequence length  
203 profiles are shown in Fig. S3. Full-length ITS reads spanned  $\sim 400$ -800 bp with a peak around  
204 550 bp; ITS2 reads from Illumina pipeline or trimmed from full-length ITS sequences were  
205 similar ( $\sim 100$ -300 bp, peak  $\sim 150$ -200 bp), reflecting the targeted subregion and potential  
206 trimming effects. We further compared both per-sample percentage of ASVs assigned and total  
207 number of ASVs classified at each taxonomic rank of two pipelines (Fig. 1a). Both pipelines  
208 showed similar performance across the higher ranks, with nearly complete assignment at  
209 kingdom (100%), phylum, class, and order, and comparable results at the genus level (75%).  
210 Differences emerged at finer ranks. Full-length ITS assigned 5,742 ASVs (98% of paired-sample  
211 mean) at the family level and 2,473 ASVs (42%) at the species level, whereas ITS2 assigned  
212 1,354 ASVs (87%) at the family level and 679 ASVs (32%) at the species level.

213 At the phylum level (Fig. S4), both pipelines were dominated by Ascomycota and  
214 Basidiomycota. Chytridiomycota was detected only in the ITS2 dataset, reflecting our inclusion  
215 of chytrid-specific primer variants in the ITS2 pipeline improved amplification of these lineages.  
216 The top 20 most abundant genera based on mean CLR abundance showed broadly similar

217 patterns between the two sequencing pipelines, particularly among the highest-abundance taxa  
218 (Fig. 1b). Both methods recovered *Cladosporium*, *Penicillium*, *Aspergillus*, and  
219 *Toxicocladosporium* as dominant genera. Several taxa, including *Komagataella* and *Malassezia*,  
220 displayed higher CLR abundance in the ITS2 data, whereas *Didymella*, *Fusarium*, *Talaromyces*,  
221 and *Trametes* were more abundant in the full-length ITS dataset. Divergence increased among  
222 lower-abundance genera, where full-length ITS resolved several taxa that appeared reduced or  
223 absent in ITS2 profiles (Fig. S5). These differences likely reflect primer-specific amplification  
224 patterns.

### 225 *Community-Level Diversity and Compositional Differences Between Sequencing Pipelines*

226 Paired comparisons of alpha diversity showed that the two pipelines produced similar overall  
227 diversity profiles across sample types (Fig. 2a). Shannon diversity did not differ between ITS2  
228 and full-length ITS in any sample type (paired Wilcoxon, all  $q \geq 0.05$ ), and values ranged from 0  
229 to 6.5. In contrast, richness metrics showed consistent increases for the full-length ITS dataset.  
230 Chao1 richness was significantly higher for full-length ITS in air, swab, and negative control  
231 samples (paired Wilcoxon  $q < 0.05$ ). Paired scatter plots further showed weak sample-level  
232 correlation in Shannon diversity and poor concordance in richness estimates between pipelines  
233 (Fig. S6), despite similar group-level Shannon distributions and consistently higher richness  
234 recovered by full-length ITS.

235 Paired and unpaired beta-diversity analyses showed that the two sequencing pipelines captured  
236 broadly similar community-level patterns, although full-length ITS revealed stronger structure  
237 associated with environmental variables than ITS2 (Fig. 2b; Fig. S8). In paired PERMANOVA  
238 restricted to matched samples, sequencing workflow explained a small but significant fraction of  
239 Bray-Curtis dissimilarity ( $R^2 \approx 0.06$ ,  $p < 0.0001$ ), and paired ordination plots showed a consistent  
240 shift between ITS2 and full-length ITS profiles derived from the same samples (Fig. 2b).  
241 Unpaired principal coordinates analysis showed clearer separation by sample type and location  
242 for the full-length ITS dataset than for ITS2 (Fig. S8). PERMANOVA indicated stronger effects  
243 of location in the full-length ITS dataset ( $R^2 = 0.12$ ,  $p < 10^{-4}$ ) than in the ITS2 dataset ( $R^2 = 0.10$ ,  
244  $p < 10^{-4}$ ), with sample type contributing a smaller but significant proportion of explained  
245 variance in both cases. In contrast, the PacBio-derived in-silico ITS2 dataset displayed ordination  
246 patterns in between full-length ITS and Illumina ITS2 (location:  $R^2 = 0.11$ ,  $p < 10^{-4}$ ), and similar  
247 sample type separation ( $R^2 = 0.06$ ,  $p = 0.0001$ ) to full-length ITS. Measures of cross-pipeline  
248 agreement supported this interpretation (Fig. S8). Mantel correlations between Bray-Curtis  
249 distance matrices were weak and not significant for both ITS2 versus full-length ITS ( $r = 0.06$ ,  
250  $p = 0.098$ ) and ITS2 versus trimmed ITS2 ( $r = 0.03$ ,  $p = 0.273$ ). In contrast, symmetric Procrustes  
251 rotation of full ordination configurations indicated moderate and significant alignment in both  
252 comparisons (ITS2 vs full-length ITS:  $r = 0.67$ ,  $p = 0.001$ ; ITS2 vs trimmed ITS2:  $r = 0.63$ ,  
253  $p = 0.001$ ). Restricting the analysis to the first 15 PCoA axes reduced alignment strength but  
254 retained statistical significance ( $r \approx 0.40$ ,  $p = 0.001$ ), indicating that shared community structure  
255 is distributed across multiple ordination dimensions rather than confined to the dominant axes.  
256 These results implicate amplicon length as the primary driver of differences in community-level  
257 inference.

## 258 *Genus-Level Agreement and Taxon-Specific Differences Between Sequencing Pipelines*

259 Agreement between sequencing pipelines was evaluated at the genus level (Fig. 3). Spearman  
260 correlations of CLR-transformed abundances for overlapping genera ranged from approximately  
261 0 to 0.7 (mean  $\pm$  SD:  $0.27 \pm 0.15$ ) across sample types and locations (Fig. 3a). Correlation  
262 strength differed significantly by sample type. Abundance was significantly higher in air samples  
263 compared with positive control samples (ANOVA  $p = 0.001$ ), indicating that the degree of  
264 cross-pipeline similarity depends on sample characteristics rather than solely on sequencing  
265 pipelines.

266 To evaluate concordance in statistical inference across pipelines, we compared signed- $\log_{10}$   
267 (p-values) from differential abundance tests for all shared genera and environmental contrasts  
268 (Fig. 3b, left). The plot indicated a moderate but highly significant correlation between pipelines  
269 (Spearman's  $\rho = 0.40$ ,  $p = 1.5 \times 10^{-60}$ ). Most genera clustered near the origin, reflecting weak or  
270 nonsignificant differences between pipelines, whereas a subset of genera exhibited strong signals  
271 in one or both datasets. Notably, a few genera (e.g. *Komagataella*, *Zasmidium*, *Aspergillus*)  
272 deviated far from the 1:1 line, indicating lineage-specific sensitivity due to primer targeting  
273 differences. To summarize how concordance varies across environmental contrasts, we  
274 computed Spearman correlation coefficients for each pairwise location comparison (Fig. 3b,  
275 right). Correlation values ranged from weak or negative to moderately strong ( $\rho \approx -0.2$  to 0.53),  
276 with most comparisons showing statistically significant agreement after false-discovery-rate  
277 correction. Together, these analyses demonstrate that while the two pipelines produce broadly  
278 similar genus-level inferences, pipeline-specific differences for genera or environmental  
279 contrasts may influence downstream ecological interpretation.

280 Health-relevant fungal genera exhibited distinct abundance patterns across sequencing pipelines,  
281 sampling types, and locations (Fig. 3c and Table S3). Ubiquitous indoor genera such as  
282 *Cladosporium*, *Penicillium*, and *Aspergillus* were consistently detected with high CLR  
283 abundances that were significantly different between pipelines (Wilcoxon adjusted  $p < 0.01$ ).  
284 Several other genera, including *Malassezia*, *Fusarium*, *Candida*, and *Saccharomyces*, also  
285 showed significant differences in CLR abundance between sequencing pipelines (adjusted  
286  $p < 0.05$ ). Several genera exhibited sample-type-dependent shifts between pipelines (Fig. 3c). In  
287 swab samples, the full-length ITS pipeline detected higher CLR abundances of *Aspergillus* and  
288 *Fusarium* (adjusted  $p < 0.05$ ). In air samples, full-length ITS also yielded higher CLR  
289 abundances of *Aspergillus* and *Penicillium*, whereas ITS2 showed higher CLR abundance of  
290 *Malassezia* (adjusted  $p < 0.05$ ). These results indicate that pipeline influences the apparent  
291 abundance of selected health-relevant taxa in a manner that depends on sample type.

## 292 *Cross-Platform Taxonomic Mapping: Unmapped ASVs*

293 To assess the degree of complementarity between sequencing pipelines, ITS2 ASVs were  
294 mapped against the full-length ITS ASVs using BLAST, treating the full-length ITS dataset as a  
295 reference. A total of 1,136 ITS2 ASVs (63%) mapped to at least one full-length ITS ASV at  
296  $\geq 99\%$  identity, corresponding to 2,731 full-length ITS ASVs (45%) that received at least one  
297 ITS2 match (Fig. 4a). The remaining 677 ITS2 ASVs (37%) and 3,383 full-length ITS ASVs  
298 (55%) lacked detectable counterparts under these criteria. Analysis at the 100% identity further

299 illustrated this pattern (Figure 4b). Most ITS2 ASVs (n = 491) showed a one-to-one relationship  
300 with full-length ITS sequences, while a subset (n = 384) mapped to multiple (2-10) full-length  
301 ITS variants. Conversely, the majority of full-length ITS ASVs had one single ITS2 ASV  
302 mapping (n = 2,046), with very few having two ITS2 sequence mapping (n = 7). Together, these  
303 results demonstrate that the two pipelines recover overlapping but non-identical representations  
304 of fungal diversity. Differences in mapping patterns, including one-to-many relationships at  
305 100% identity, are consistent with amplicon-length-driven resolution differences and primer  
306 targeting effects, rather than the absence of taxa in either dataset.

307 Unmapped ASVs from both pipelines were predominantly low in read abundance (Fig. 4 c-d).  
308 Most unmapped ITS2 ASVs had fewer than 500 total reads, with the majority below 100 reads  
309 (Fig. 4c). Although most unmapped full-length ITS ASVs also occurred at low abundance, a  
310 subset displayed moderate to high read counts when visualized on a log<sub>10</sub> scale (Fig. 4c, right).  
311 This increase likely reflects both the substantially higher sequencing depth of the full-length ITS  
312 dataset and the greater resolving power of long-read amplicons. As a result, biologically  
313 abundant lineages that are represented by a small number of ITS2 ASVs can appear as multiple  
314 unmapped full-length ITS ASVs under strict 100% identity mapping criteria. Unmapped ASVs  
315 from both pipelines were primarily affiliated with Ascomycota and Basidiomycota, with smaller  
316 contributions (<100 ASVs) from under-classified fungi, Chytridiomycota (ITS2), Mucoromycota  
317 (full-length ITS), and Mortierellomycota (both pipelines) (Fig. 4 d). A subset of unmapped ASVs  
318 belonged to genera with relevance to indoor environmental health (Fig. 4e). Unmapped ITS2  
319 ASVs included low-abundance assignments to genera *Malassezia*, *Aspergillus*, *Candida*,  
320 *Penicillium*, and *Fusarium* (Fig. 4e, left). Unmapped full-length ITS ASVs also included genera  
321 *Cladosporium*, *Aspergillus*, *Penicillium*, and *Fusarium*, often with higher cumulative read counts  
322 than their ITS2 counterparts (Fig. 4e, right). These unmapped health-associated ASVs likely  
323 reflect the influences of primer targeting and amplicon length on pipeline-specific detection,  
324 which can vary even within clinically relevant lineages. Together, these results demonstrate that  
325 ITS2 and full-length ITS pipelines provide overlapping yet distinct representations of fungal  
326 community composition, with primer targeting, amplicon length, and sequencing depth  
327 contributing to pipeline-specific detection of low-abundance taxa.

### 328 *Cross-Platform Taxonomic Mapping: mapped ASVs*

329 Comparison of phylogenetic trees constructed from mapped ASVs demonstrates the enhanced  
330 resolving power of long-read sequencing (Fig. 5a). Although both phylogenies were largely  
331 bifurcating, the ITS2-derived tree showed significantly shorter branch lengths than the  
332 full-length ITS tree (median 0.0266 vs. 0.0367; Wilcoxon test,  $p = 1.6 \times 10^{-8}$ , Fig. S9), indicating  
333 reduced phylogenetic signal in the ITS2 region. Longer branches in the full-length ITS tree  
334 reflect increased sequence variation across ITS1, 5.8S, and ITS2. Among the ITS2 ASVs  
335 mapped at 100% identity to full-length ITS sequences, taxonomic assignment discrepancies were  
336 observed at multiple ranks (Fig. 5b). A total of 17 ASVs had taxonomy assignment differed at  
337 the species level and 6 at the genus level. These discrepancies likely reflect the increased  
338 discriminatory power of full-length ITS sequences. Cross-platform comparison also revealed  
339 shifts in taxonomic resolution (Fig. 5c). Among the mapped ITS2 ASVs, the majority (599  
340 ASVs) remained unchanged, suggesting broad agreement of taxa with robust representation  
341 across databases. Around 15% ASVs (n = 131) gained higher-resolution classification after

342 mapping to full-length ITS reference, highlighting the added discriminatory information obtained  
343 from longer reads. Around 14% ASVs (n = 121) were reassigned to broader taxonomic ranks,  
344 likely due to the primer target differences and the limited coverage of full-length ITS reference  
345 we constructed. These patterns illustrate that cross-platform mapping exerts an uneven influence  
346 on taxonomic resolution, depending on both sequence informativeness and reference database  
347 completeness.

348 The 24 ITS2 ASVs that had different assignment after mapping to the full-length ITS dataset  
349 were predominantly low in read abundance (< 1000 reads, Fig. 5d), except one ASV (ASV1382 ,  
350 assigned to *Komagataella kurtzmanii*) had over 8000 total reads. It was mapped at 100% identity  
351 to two full length ITS ASVs (ASV1411 and ASV6353, both assigned to *Komagataella ulmi*) at  
352 low abundance of 1061 total reads. This suggests that classification disagreement mainly affects  
353 rare taxa rather than dominant community members. A Sankey diagram visualizing taxonomic  
354 transitions between ITS2 and full-length ITS classifications further illustrate these dynamics  
355 (Fig. 5e). Most ASVs retained stable higher-level taxonomic placement, with a few ASVs gained  
356 different genus or species assignments, demonstrating how short-read and long-read sequences  
357 can yield divergent interpretations for certain taxa. These shifts emphasize the need for cautious  
358 cross-platform comparisons and highlight the importance of comprehensive, high-quality  
359 reference databases for achieving accurate taxonomic resolution in mycobiome studies.

## 360 Discussion

361 This study provides a systematic evaluation of Illumina ITS2 and PacBio HiFi full-length ITS  
362 sequencing for indoor fungal community profiling, with emphasis on how sequencing strategy  
363 influences taxonomic resolution, community-level inference, and detection of health-relevant  
364 taxa. By integrating paired analyses, in-silico ITS2 trimming, and cross-platform sequence  
365 mapping, we demonstrate that apparent pipeline differences arise primarily from amplicon length  
366 and primer targets (19, 20, 38, 39).

### 367 *Amplicon length as the dominant driver of taxonomic resolution*

368 Full-length ITS assigned over 40% of ASVs to species compared with approximately 32% for  
369 ITS2, reflecting the additional sequence variation captured across ITS1, 5.8S, and ITS2. This  
370 improvement is consistent with prior work showing that short-amplicon barcoding limits  
371 discrimination among closely related taxa, including clinically relevant species that differ by few  
372 substitutions within ITS2 alone or have identical ITS2 regions (18, 20, 21), and parallels gains  
373 reported when full-length 16S rRNA sequencing is applied to bacterial communities (40, 41).  
374 Reference database completeness imposes an additional constraint on both pipelines: the UNITE  
375 database remains incomplete for many indoor-relevant fungal lineages, potentially  
376 underestimating the resolution advantage of full-length ITS and introducing differential  
377 classification biases (18, 33). Critically, in-silico trimming of PacBio reads to ITS2 caused  
378 convergence toward Illumina ITS2 in both assignment rates (Fig. S7) and community structure  
379 (Fig. S8), confirming that amplicon length rather than platform-specific error models or  
380 sequencing depth underlies the observed differences (19, 38).

### 381 *Community-level diversity and structure are broadly conserved across pipelines*

382 Community-level diversity metrics were broadly conserved despite differences in resolution.  
383 Shannon diversity did not differ between pipelines in any sample type, indicating that dominant

384 taxa and overall evenness are robust to sequencing strategy, a pattern consistent with platform  
385 comparisons in bacterial microbiome research (40, 41). In contrast, richness metrics (Chao1 and  
386 observed ASVs) were consistently higher for full-length ITS, reflecting increased resolution of  
387 intra-taxon sequence variation that accumulates as elevated richness in long-read datasets (19,  
388 20). At the beta-diversity level, both pipelines recovered broadly similar community patterns, but  
389 full-length ITS revealed stronger structure associated with sampling location and sample type.  
390 PERMANOVA indicated that workflow explained only a small fraction of compositional  
391 variance ( $R^2 \approx 0.06$ ) relative to environmental factors, Procrustes analyses showed moderate  
392 ordination alignment, and weak Mantel correlations indicated that cross-pipeline agreement is  
393 distributed across many dimensions rather than concentrated in primary axes. The intermediate  
394 ordination position of in-silico trimmed PacBio-ITS2 data between full-length ITS and Illumina  
395 ITS2 further corroborates amplicon length as the primary driver (38), and confirms that ITS2  
396 captures major community gradients in indoor environments (7, 16) while full-length ITS  
397 enhances discrimination along secondary axes of environmental heterogeneity.

#### 398 *Genus-level agreement and context-dependent discrepancies*

399 Genus-level abundance patterns were generally concordant between pipelines, particularly for  
400 dominant indoor genera such as *Cladosporium*, *Penicillium*, and *Aspergillus*, consistent with  
401 their prevalence across sequencing-based indoor surveys worldwide (5–7, 42). However,  
402 agreement varied by sample type and environment, with laboratory samples showing higher  
403 cross-platform correlations than residential or dormitory environments. This context dependence  
404 likely reflects interactions between community complexity and primer performance. In occupied  
405 buildings, higher  $\alpha$ -diversity and greater representation of amplification-resistant lineages may  
406 amplify primer-targeting biases documented for ITS3/ITS4-type primers, particularly for  
407 Basidiomycota (27, 39). Differential abundance testing showed moderate global concordance  
408 (Spearman  $\rho = 0.40$ ) alongside lineage-specific discrepancies, with certain genera (including  
409 *Komagataella*, *Zasmidium*, and *Aspergillus*) deviating markedly from the 1:1 line. For example,  
410 a high-abundance ITS2 ASV assigned to *K. kurtzmanii* mapped at 100% identity to two full-  
411 length ITS ASVs assigned to *K. ulmi*, illustrating that species-level inference can diverge  
412 between pipelines even for abundant taxa. Researchers should therefore be cautious when cross-  
413 comparing differential abundance results across sequencing strategies (18, 43). Although both  
414 datasets were classified using the same reference database and algorithm, concordance remains  
415 limited by incomplete fungal ITS reference resources and the insufficient discriminatory power  
416 of ITS1 or ITS2 alone for many taxa (43). Consequently, long-read ITS can resolve sequence  
417 variation that remains ambiguous in shorter ITS2 queries, producing marker-dependent  
418 differences that reflect reference limitations rather than pipeline error.

#### 419 *Health-relevant taxa are differentially affected by sample type and marker choice*

420 Differences in health-relevant genera between pipelines were driven primarily by sample type  
421 rather than location. Full-length ITS detected higher CLR abundances of *Aspergillus* and  
422 *Fusarium* in swab samples, and of *Aspergillus* and *Penicillium* in air samples, while ITS2  
423 showed higher apparent abundance of *Malassezia* in air samples. The enrichment of *Malassezia*  
424 in ITS2 data reflects its unusually short ITS regions, which are preferentially amplified during  
425 PCR, combined with strong UNITE database representation (39, 44). Conversely, the greater  
426 sensitivity of full-length ITS for *Aspergillus* and *Penicillium* likely reflects improved resolution

427 of species complexes collapsed to single ITS2 ASVs (10, 11); these genera encompass clinically  
428 distinct species with direct implications for indoor exposure and risk assessment (6, 9).

#### 429 *Cross-platform mapping reveals resolution, not detection, differences*

430 Cross-platform mapping analyses reinforce these interpretations: while most ITS2 ASVs mapped  
431 to full-length ITS at  $\geq 99\%$  identity, full-length ITS frequently partitioned abundant taxa into  
432 multiple closely related ASVs that collapsed into single ITS2 sequences, consistent with ITS2  
433 underestimating within-taxon sequence diversity at the ASV level (18, 19, 38). Unmapped ASVs  
434 from both pipelines were predominantly low-abundance, though a subset of unmapped full-  
435 length ITS ASVs had moderate read counts, reflecting higher PacBio sequencing depth and  
436 greater resolving power rather than detection failure. Notably, health-relevant genera including  
437 *Aspergillus*, *Penicillium*, and *Fusarium* appeared among unmapped ASVs in both datasets (9–  
438 11), highlighting the value of inspecting unmapped fractions in pipeline comparisons.  
439 Phylogenetic trees from mapped ASVs further illustrated resolution differences: full-length ITS  
440 trees exhibited longer branches and more bifurcating nodes than ITS2 trees at the same  
441 evolutionary scale (19, 21), and taxonomic discrepancies among 100%-identity-mapped ASVs at  
442 species and genus levels reinforce calls for comprehensive reference databases and transparent  
443 reporting of taxonomic resolution (18, 33, 43).

#### 444 *Implications for indoor mycobiome studies*

445 Taken together, these results support a complementary, use-case-driven approach to indoor  
446 mycobiome sequencing. ITS2 remains efficient and cost-effective for large-scale surveys and  
447 community-level comparisons where patterns among dominant taxa are of primary interest, and  
448 is well-suited to epidemiological studies linking indoor fungi to health outcomes (3, 5, 9). Full-  
449 length ITS is particularly valuable when species-level resolution informs exposure or clinical risk  
450 assessment, especially for genera such as *Aspergillus* and *Fusarium* (6, 10). Recent  
451 improvements in PacBio HiFi throughput via the Revio system and Kinnex library preparation  
452 have reduced per-sample costs and expanded multiplexing feasibility (22, 28). With these  
453 improvements, full-length ITS will become increasingly tractable for large built-environmental  
454 studies. Meanwhile, a hybrid strategy combining ITS2 for community-wide profiling with  
455 targeted full-length ITS validation of taxa of health concern provides a practical middle ground.  
456 Improved species-level identification also enables downstream interpretation relevant to human  
457 health. In particular, genomic and phenotypic markers associated with mycotoxin production,  
458  $\beta$ -glucan content, immune activation, and other pathogenic traits can inform the development of  
459 targeted assays and more informative exposure and health-risk assessments (45–47). More  
460 broadly, we found that community complexity in occupied settings amplifies discrepancies  
461 between pipelines, highlighting the importance of transparent reporting of sequencing strategy,  
462 reference database version, and bioinformatic pipeline choices for meaningful cross-study  
463 comparisons (18, 27, 43), particularly given the documented geographic and environmental  
464 structuring of indoor fungal communities (7, 15, 16).

#### 465 *Limitations*

466 This study has several limitations. The number of sites was limited, which may restrict  
467 generalizability of environment-specific conclusions (7, 16). Taxonomic assignments depend on  
468 the completeness of database, which remains incomplete for many lineages and may  
469 differentially affect short- and long-read annotations (18, 33). In-silico ITS2 trimming does not  
470 fully recapitulate PCR amplification with ITS2-specific primers, including primer binding  
471 efficiency, amplicon length bias, and chimera formation (19, 38). In addition, this analysis relied

472 on DNA-based profiling and did not assess fungal viability, activity, or mycotoxin production;  
473 future integration of RNA-based or culture-based approaches would provide a more complete  
474 picture of active indoor fungal communities (6, 9, 43, 48). Finally, the higher sequencing depth  
475 of the PacBio dataset may have contributed to differences in low-abundance taxon detection.  
476 Future studies should target more comparable sequencing depths across platforms to better  
477 disentangle depth effects from those of amplicon length and primer targeting. As whole-genome  
478 sequencing of environmental and built-environment fungal isolates expands, reference databases  
479 are expected to better support long-read ITS classification, further improving species-level  
480 resolution beyond what is currently achievable.

481  
482 In conclusion, amplicon length is the primary determinant of differences between Illumina ITS2  
483 and PacBio full-length ITS pipelines for indoor mycobiome analysis. Both pipelines recover  
484 broadly similar community-level patterns for dominant taxa, but full-length ITS provides  
485 enhanced taxonomic resolution, improved richness estimates, and stronger ecological  
486 discrimination at the family and species levels. Platform-specific differences in the apparent  
487 abundance of health-relevant genera including *Aspergillus*, *Penicillium*, *Fusarium*, and  
488 *Malassezia* highlight how sequencing strategy and sampling matrix jointly shape the  
489 interpretation of indoor fungal exposure. Method selection should be guided by study scale,  
490 required taxonomic resolution, and research or exposure-assessment objectives in built-  
491 environment mycobiome studies.

#### 492 **Data Availability**

493 Sequences have been deposited at NCBI SRA (National Center for Biotechnology Information  
494 Sequence Read Archive, <https://www.ncbi.nlm.nih.gov/sra>) under BioProject ID  
495 PRJNA1294855.

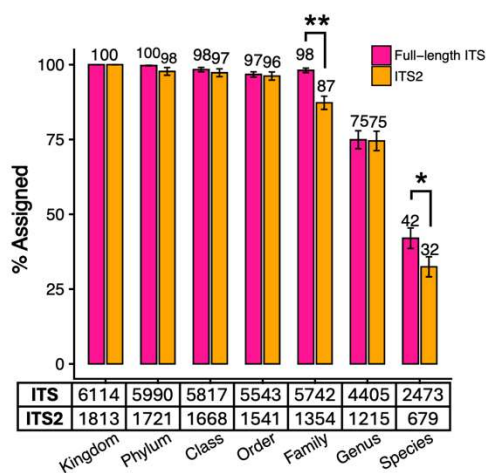
#### 496 **Acknowledgement**

497 This work was supported by the National Science Foundation Engineering Research Center  
498 Precision Microbiome Engineering (NSF PreMiEr ERC) through Award No. 2133504.  
499 Sequencing services were provided by the Duke Microbiome Sequencing Core, and we thank the  
500 core staff for their technical support and guidance.

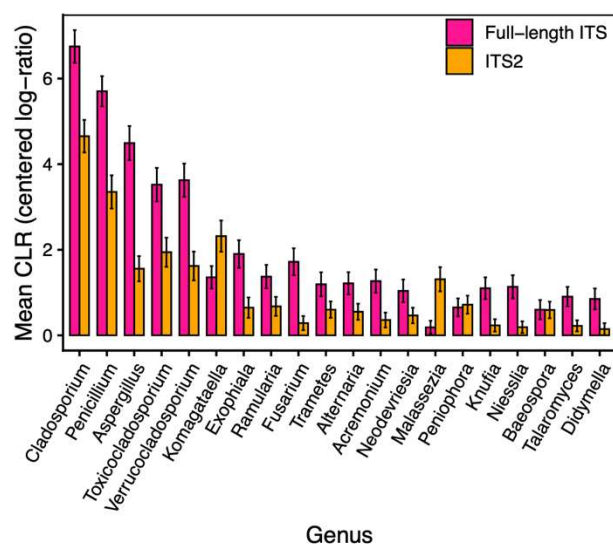
501

502

**a. Paired Taxonomy Resolution Comparison**

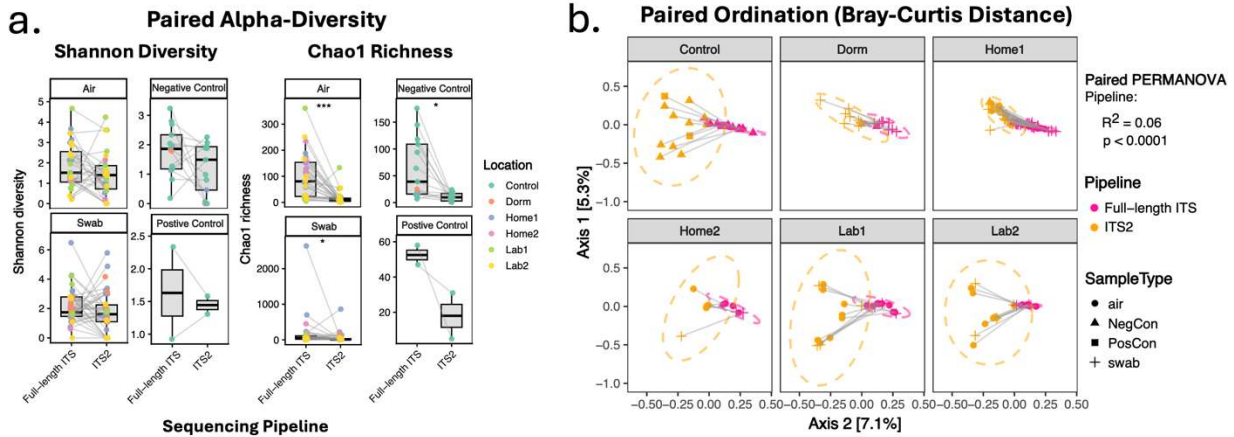


**b. Top 20 most abundant genera (Mean CLR ± SE)**



504

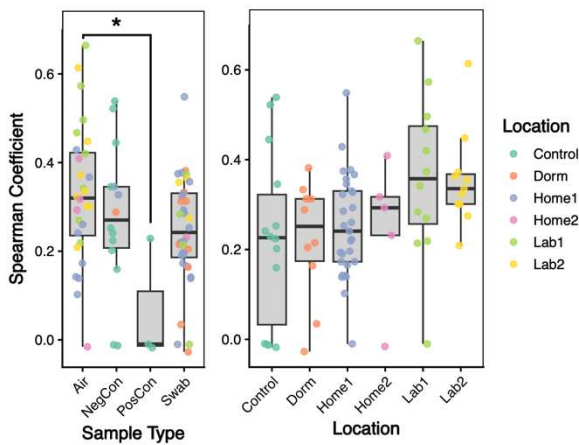
505 Fig 1. Comparative performance of Illumina ITS2 and PacBio full-length ITS sequencing  
 506 platforms. (a) Taxonomic assignment across all ranks comparing PacBio full-length (deep pink)  
 507 and Illumina ITS2 (orange) pipelines by paired sample comparison, numbers above the bars  
 508 indicate percent of ASV classified ± standard error (SE) for each rank, numbers in the table  
 509 under the plot showed the total number of ASVs classified at each rank from each pipeline. (b)  
 510 Top 20 genera based on mean CLR abundance across both sequencing pipelines. Bars show  
 511 mean centered log-ratio (CLR) abundance ± standard error (SE) for each genus, with paired  
 512 values plotted for full-length ITS and ITS2.



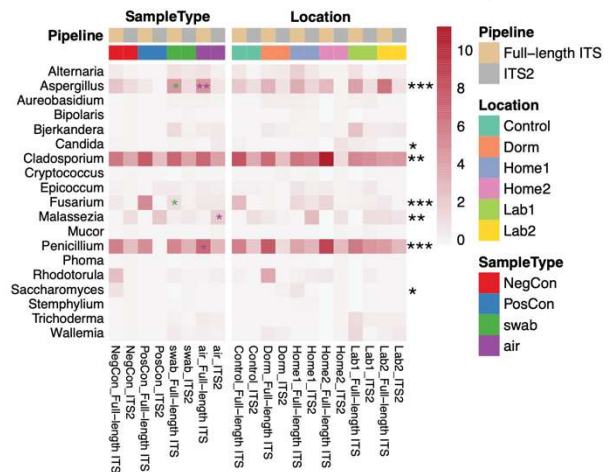
513

514 Fig 2. Community-Level Diversity and Compositional Differences Between Sequencing  
 515 Pipelines. (a) Paired alpha-diversity comparisons for ITS2 and full-length ITS across sample  
 516 types. Side-by-side boxplots show Shannon diversity (left) and Chao1 richness (right) for paired  
 517 samples, with connecting lines linking the same sample across pipelines. Points are colored by  
 518 sampling location. Statistical significance labeled for paired Wilcoxon test comparing alpha-  
 519 diversity index values between pipelines. (b). Paired ordination comparing ITS2 and full-length  
 520 ITS profiles using Bray-Curtis dissimilarity. Principal coordinates analysis (PCoA) shows each  
 521 sample pair connected by a line, with ITS2 and full-length ITS placed in proximity but  
 522 occupying distinct positions. Ellipses represent 95% confidence intervals for location groups.  
 523 Paired PERMANOVA indicated a significant pipeline effect ( $R^2 \approx 0.06$ ,  $p < 0.001$ ), consistent  
 524 with the systematic shifts observed between ITS2 and full-length ITS profiles. Asterisks indicate  
 525 statistically significant differences between short vs long reads or sample types (\*\* for p-value  
 526 between 0-0.001, \* for p-value between 0.001-0.01, \* for p-value between 0.01 - 0.05).

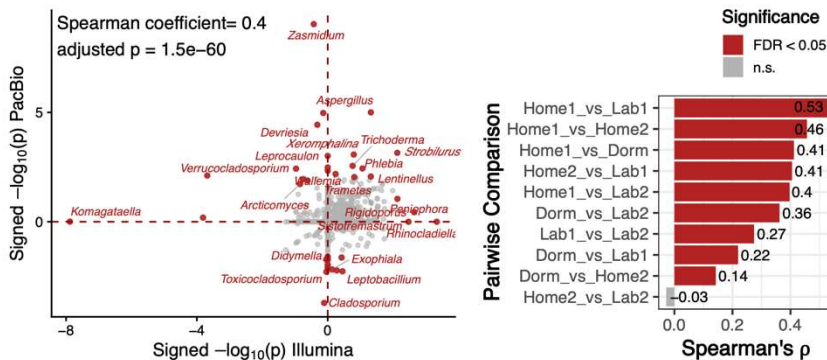
**a. Abundance correlation of overlapped genera (per sample, CLR transformed)**



**c. CLR abundance of environmental health genera**



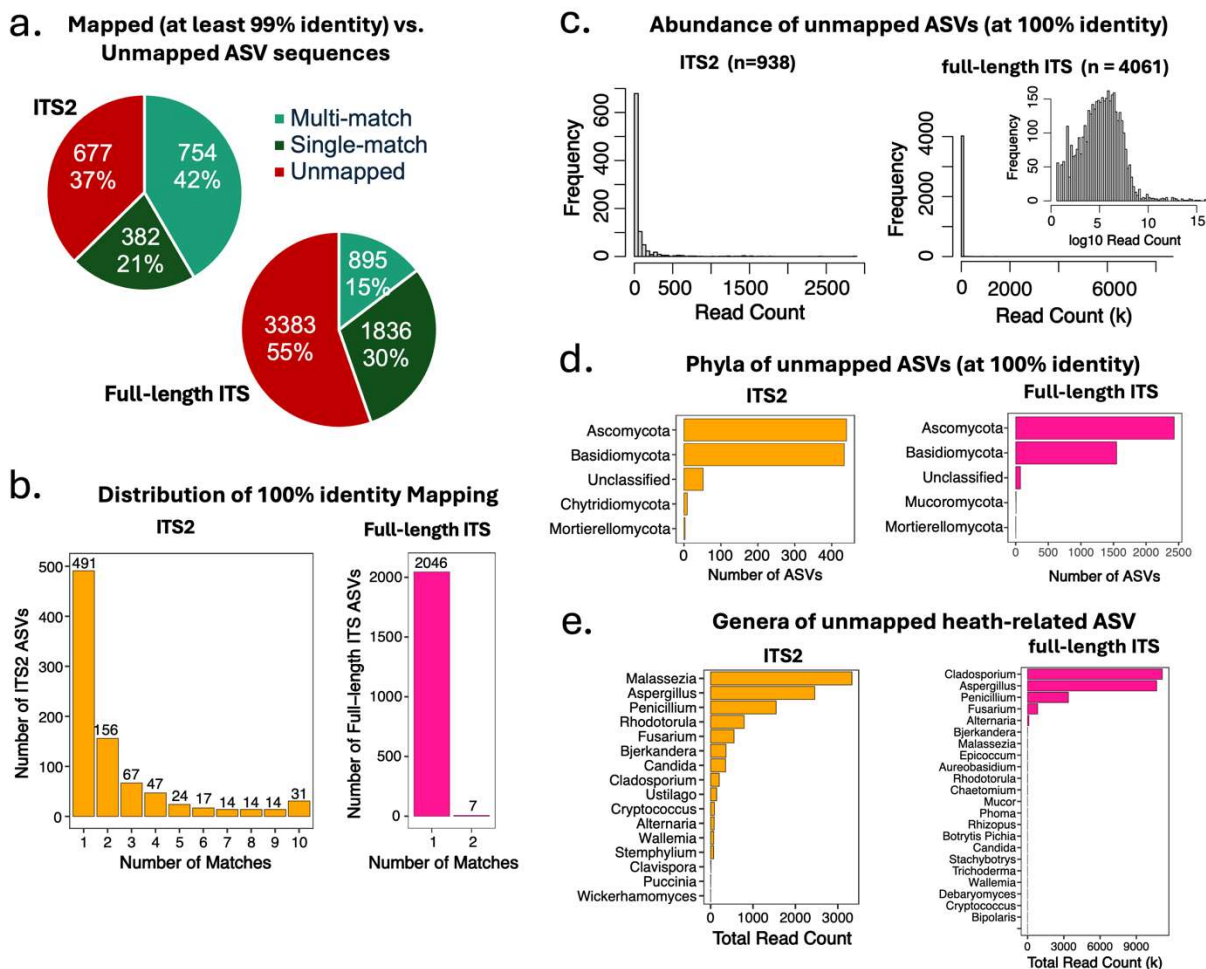
**b. Concordance of Overlapping Genera between pipelines**



527

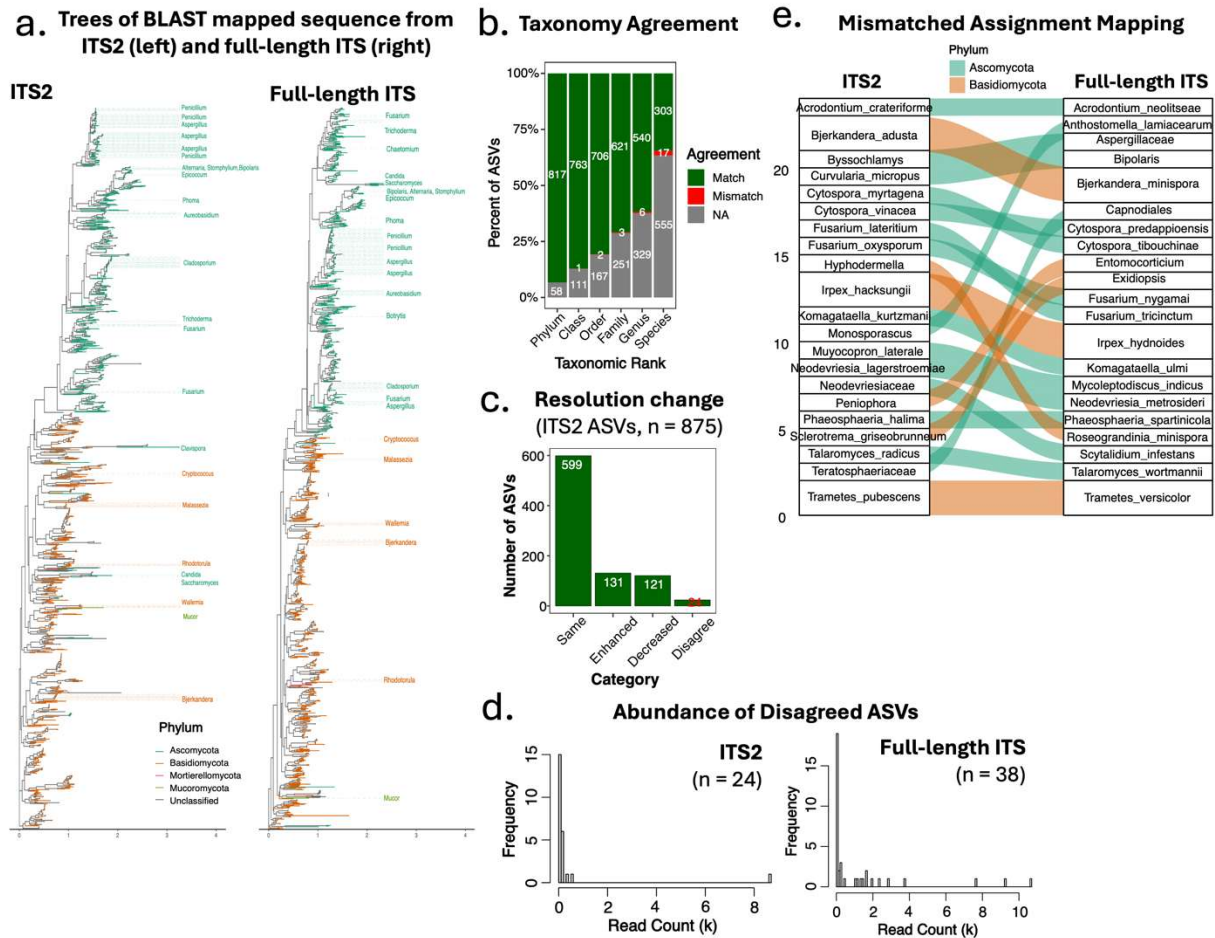
528 Fig 3. Genus-level agreement and taxon-specific differences between sequencing pipelines. (a)  
 529 Abundance correlation of overlapping genera. Spearman correlations (per sample, CLR-  
 530 transformed abundances) between ITS2 and full-length ITS for overlapping genera, grouped by  
 531 sample type (left) and location (right). Correlation strength varied significantly across sample  
 532 types and locations, with laboratory samples showing the highest cross-pipeline agreement. Each  
 533 point represents one sample, colored by sampling location. (b) Concordance of taxon-level  
 534 statistical inference across pipelines. Top: Scatter plot of signed  $-\log_{10}(p)$ -values from  
 535 differential abundance analyses for shared genera across all pairwise environmental contrasts,  
 536 comparing ITS2 (x-axis) and full-length ITS (y-axis). Red points indicate genera with  
 537  $FDR < 0.05$  in either pipeline. Dashed red lines mark direction reversals. Annotated genera  
 538 highlight lineages with strong or pipeline-specific signals. The overall Spearman correlation  
 539 ( $\rho = 0.40$ , adjusted  $p = 1.5 \times 10^{-60}$ ) indicates moderate global concordance. Right: Spearman  
 540 correlations for each pairwise location comparison. Bars show correlation strength, with  
 541 significant comparisons ( $FDR < 0.05$ ) highlighted. c. Abundance patterns of health-relevant  
 542 genera. Heatmap of mean CLR abundance for selected health-relevant fungal genera across  
 543 sample types (air, swab, controls) and indoor locations (Home1, Home2, Dorm, Lab). Asterisks  
 544 indicate significant differences between pipelines or sample types ( $*** p < 0.001$ ,  $** p < 0.01$ ,  
 545  $* p < 0.05$ ). Several genera displayed strong pipeline-dependent or location-specific patterns,

546 illustrating how sequencing strategy and environmental context shape the detected indoor  
 547 mycobiome.



548

549 Fig 4. Cross-platform taxonomic mapping analysis of ITS2 pipeline against full-length ITS  
 550 pipeline, using full-length ITS data as the reference dataset. (a) Pie chart showing mapped and  
 551 unmapped ASV counts and percentages in ITS2 (left) and full-length ITS (right) pipelines.  
 552 Unmapped indicates ITS2 ASV sequence was not mapped to any of full-length ITS sequence,  
 553 single match indicates ITS2 sequence was mapped to one full-length ITS sequence, multi-match  
 554 indicates ITS2 sequence mapped to 2-10 full length ITS sequence. (b) Distribution of 100%  
 555 identity mapping. (c) Read abundance distribution of unmapped ITS2 (left) and full-length ITS  
 556 (right) ASVs demonstrating predominance of low-abundance taxa. (d) Phyla distributions of  
 557 unmapped ITS2 (left) and full-length ITS (right) ASVs. (e) Read abundance of unmapped ASVs  
 558 assigned to health-related genera from the ITS2 (left) and full-length ITS (right) pipelines.



559

560 Fig 5. Cross-platform taxonomic mapping analysis of ITS2 pipeline against full-length ITS  
 561 pipeline, using full-length ITS data as the reference dataset. (a) Rooted maximum likelihood  
 562 trees constructed using BLAST mapped (> 97% identity) sequence from ITS2 (left) and full-  
 563 length ITS (right) pipelines. The branches are colored by phylum, and health-related taxa are  
 564 labeled at the genus level. (b) Taxonomic agreement analysis for 100% mapped ASVs, showing  
 565 number of mismatched assignments at each taxonomic rank, NA(gray) indicates missing data in  
 566 either or both pipelines. (c) Changes in ITS2 taxonomic resolution following mapping (100%  
 567 identity), categorized as decreased, enhanced, or same resolution compared to original  
 568 assignments. (d) Abundance distribution of 100% mapped but taxonomically mismatched ITS2  
 569 and full-length ITS ASVs. (e) Sankey diagram of taxonomic assignment mapping of  
 570 taxonomically mismatched ASVs: ITS2 (left) and full-length ITS (right), with flow colors  
 571 representing major fungal phyla.

572 **References**

- 573 1. Adams RI, Bateman AC, Bik HM, Meadow JF. 2015. Microbiota of the indoor environment:  
574 a meta-analysis. *Microbiome* 3.
- 575 2. Bosch TCG, Wigley M, Colomina B, Bohannon B, Meggers F, Amato KR, Azad MB, Blaser  
576 MJ, Brown K, Dominguez-Bello MG, Ehrlich SD, Elinav E, Finlay BB, Geddie K, Geva-  
577 Zatorsky N, Giles-Vernick T, Gros P, Guillemin K, Haraoui L-P, Johnson E, Keck F, Lorimer  
578 J, Mcfall-Ngai MJ, Nichter M, Pettersson S, Poinar H, Rees T, Tropini C, Undurraga EA,  
579 Zhao L, Melby MK. 2024. The potential importance of the built-environment microbiome and  
580 its impact on human health. *Proceedings of the National Academy of Sciences* 121.
- 581 3. Hegarty B, Haverinen-Shaughnessy U, Shaughnessy RJ, Peccia J. 2019. Spatial Gradients of  
582 Fungal Abundance and Ecology throughout a Damp Building. *Environ Sci Technol Lett*  
583 6:329–333.
- 584 4. Gilbert JA, Hartmann EM. 2024. The indoors microbiome and human health. *Nat Rev*  
585 *Microbiol* 22:742–755.
- 586 5. Hickman B, Kirjavainen PV, Täubel M, de Vos WM, Salonen A, Korpela K. 2022.  
587 Determinants of bacterial and fungal microbiota in Finnish home dust: Impact of  
588 environmental biodiversity, pets, and occupants. *Front Microbiol* 13.
- 589 6. Nevalainen A, Täubel M, Hyvärinen A. 2015. Indoor fungi: companions and contaminants.  
590 *Indoor Air* 25:125–156.
- 591 7. Amend AS, Seifert KA, Samson R, Bruns TD. 2010. Indoor fungal composition is  
592 geographically patterned and more diverse in temperate zones than in the tropics. *Proceedings*  
593 *of the National Academy of Sciences* 107:13748–13753.
- 594 8. Borman AM, Johnson EM. 2023. Changes in fungal taxonomy: mycological rationale and  
595 clinical implications. *Clinical Microbiology Reviews* 36:e00099-22.
- 596 9. Dannemiller KC, Mendell MJ, Macher JM, Kumagai K, Bradman A, Holland N, Harley K,  
597 Eskenazi B, Peccia J. 2014. Next-generation DNA sequencing reveals that low fungal  
598 diversity in house dust is associated with childhood asthma development. *Indoor Air* 24:236–  
599 247.
- 600 10. Latgé J-P, Chamilos G. 2019. *Aspergillus fumigatus* and Aspergillosis in 2019. *Clinical*  
601 *Microbiology Reviews* 33:10.1128/cmr.00140-18.
- 602 11. Bennett JW, Klich M. 2003. Mycotoxins. *Clinical Microbiology Reviews* 16:497–516.
- 603 12. Casadevall A, Perfect JR. 1998. *Cryptococcus neoformans*. Citeseer.
- 604 13. Elkady EA, Torky H. 2016. *Cryptococcus Neoformans* Isolated from Pigeon, Chicken  
605 Dropping Samples and Dust Collected From Their House on Urease Hydrolysis. *Alexandria*  
606 *Journal of Veterinary Sciences* 51.
- 607 14. Rajasingham R, Smith RM, Park BJ, Jarvis JN, Govender NP, Chiller TM, Denning DW,  
608 Loyse A, Boulware DR. 2017. Global burden of disease of HIV-associated cryptococcal  
609 meningitis: an updated analysis. *Lancet Infect Dis* 17:873–881.

- 610 15. Chase J, Fouquier J, Zare M, Sonderegger DL, Knight R, Kelley ST, Siegel J, Caporaso JG.  
611 2016. Geography and Location Are the Primary Drivers of Office Microbiome Composition.  
612 *mSystems* 1:10.1128/msystems.00022-16.
- 613 16. Lax S, Smith DP, Hampton-Marcell J, Owens SM, Handley KM, Scott NM, Gibbons SM,  
614 Larsen P, Shogan BD, Weiss S, Metcalf JL, Ursell LK, Vázquez-Baeza Y, Van Treuren W,  
615 Hasan NA, Gibson MK, Colwell R, Dantas G, Knight R, Gilbert JA. 2014. Longitudinal  
616 analysis of microbial interaction between humans and the indoor environment. *Science*  
617 345:1048–1052.
- 618 17. Schoch CL, Seifert KA, Huhndorf S, Robert V, Spouge JL, Levesque CA, Chen W, Fungal  
619 Barcoding Consortium, Fungal Barcoding Consortium Author List, Bolchacova E, Voigt K,  
620 Crous PW, Miller AN, Wingfield MJ, Aime MC, An K-D, Bai F-Y, Barreto RW, Begerow D,  
621 Bergeron M-J, Blackwell M, Boekhout T, Bogale M, Boonyuen N, Burgaz AR, Buyck B, Cai  
622 L, Cai Q, Cardinali G, Chaverri P, Coppins BJ, Crespo A, Cubas P, Cummings C, Damm U,  
623 de Beer ZW, de Hoog GS, Del-Prado R, Dentinger B, Diéguez-Uribeondo J, Divakar PK,  
624 Douglas B, Dueñas M, Duong TA, Eberhardt U, Edwards JE, Elshahed MS, Fliiegerova K,  
625 Furtado M, García MA, Ge Z-W, Griffith GW, Griffiths K, Groenewald JZ, Groenewald M,  
626 Grube M, Gryzenhout M, Guo L-D, Hagen F, Hambleton S, Hamelin RC, Hansen K, Harrold  
627 P, Heller G, Herrera C, Hirayama K, Hirooka Y, Ho H-M, Hoffmann K, Hofstetter V,  
628 Högnabba F, Hollingsworth PM, Hong S-B, Hosaka K, Houbraken J, Hughes K, Huhtinen S,  
629 Hyde KD, James T, Johnson EM, Johnson JE, Johnston PR, Jones EBG, Kelly LJ, Kirk PM,  
630 Knapp DG, Kõljalg U, Kovács GM, Kurtzman CP, Landvik S, Leavitt SD, Ligenstoffer AS,  
631 Liimatainen K, Lombard L, Luangsa-ard JJ, Lumbsch HT, Maganti H, Maharachchikumbura  
632 SSN, Martin MP, May TW, McTaggart AR, Methven AS, Meyer W, Moncalvo J-M,  
633 Mongkolsamrit S, Nagy LG, Nilsson RH, Niskanen T, Nyilasi I, Okada G, Okane I, Olariaga  
634 I, Otte J, Papp T, Park D, Petkovits T, Pino-Bodas R, Quaedvlieg W, Raja HA, Redecker D,  
635 Rintoul TL, Ruibal C, Sarmiento-Ramírez JM, Schmitt I, Schüßler A, Shearer C, Sotome K,  
636 Stefani FOP, Stenroos S, Stielow B, Stockinger H, Suetrong S, Suh S-O, Sung G-H, Suzuki  
637 M, Tanaka K, Tedersoo L, Telleria MT, Tretter E, Untereiner WA, Urbina H, Vágvolgyi C,  
638 Vialle A, Vu TD, Walther G, Wang Q-M, Wang Y, Weir BS, Weiß M, White MM, Xu J,  
639 Yahr R, Yang ZL, Yurkov A, Zamora J-C, Zhang N, Zhuang W-Y, Schindel D. 2012. Nuclear  
640 ribosomal internal transcribed spacer (ITS) region as a universal DNA barcode marker for  
641 Fungi. *Proceedings of the National Academy of Sciences* 109:6241–6246.
- 642 18. Nilsson RH, Anslan S, Bahram M, Wurzbacher C, Baldrian P, Tedersoo L. 2019.  
643 Mycobiome diversity: high-throughput sequencing and identification of fungi. *Nat Rev*  
644 *Microbiol* 17:95–109.
- 645 19. Tedersoo L, Tooming-Klunderud A, Anslan S. 2018. PacBio metabarcoding of Fungi and  
646 other eukaryotes: errors, biases and perspectives. *New Phytologist* 217:1370–1385.
- 647 20. Tedersoo L, Anslan S. 2019. Towards PacBio-based pan-eukaryote metabarcoding using  
648 full-length ITS sequences. *Environmental microbiology reports* 11:659–668.
- 649 21. Tedersoo L, Anslan S, Bahram M, Kõljalg U, Abarenkov K. 2020. Identifying the  
650 ‘unidentified’ fungi: a global-scale long-read third-generation sequencing approach. *Fungal*  
651 *Diversity* 103:273–293.

- 652 22. Wenger AM, Peluso P, Rowell WJ, Chang P-C, Hall RJ, Concepcion GT, Ebler J,  
653 Functammasan A, Kolesnikov A, Olson ND, Töpfer A, Alonge M, Mahmoud M, Qian Y,  
654 Chin C-S, Phillippy AM, Schatz MC, Myers G, DePristo MA, Ruan J, Marschall T,  
655 Sedlazeck FJ, Zook JM, Li H, Koren S, Carroll A, Rank DR, Hunkapiller MW. 2019.  
656 Accurate circular consensus long-read sequencing improves variant detection and assembly of  
657 a human genome. *Nat Biotechnol* 37:1155–1162.
- 658 23. Lundberg DS, Yourstone S, Mieczkowski P, Jones CD, Dangl JL. 2013. Practical  
659 innovations for high-throughput amplicon sequencing. *Nat Methods* 10:999–1002.
- 660 24. Tedersoo L, Bahram M, Põlme S, Kõljalg U, Yorou NS, Wijesundera R, Ruiz LV, Vasco-  
661 Palacios AM, Thu PQ, Suija A, others. 2014. Global diversity and geography of soil fungi.  
662 *science* 346:1256688.
- 663 25. White TJ, Bruns T, Lee S, Taylor J, others. 1990. Amplification and direct sequencing of  
664 fungal ribosomal RNA genes for phylogenetics. *PCR protocols: a guide to methods and*  
665 *applications* 18:315–322.
- 666 26. Cregger M, Veach A, Yang Z, Crouch M, Vilgalys R, Tuskan G, Schadt C. 2018. The  
667 *Populus holobiont*: dissecting the effects of plant niches and genotype on the microbiome.  
668 *Microbiome* 6:1–14.
- 669 27. Tedersoo L, Lindahl B. 2016. Fungal identification biases in microbiome projects.  
670 *Environmental microbiology reports* 8:774–779.
- 671 28. Pacific Bioscience. 2025. Preparing Kinnex libraries from 16S rRNA amplicons: Procedure  
672 and Checklist.pdf. PacBio.
- 673 29. Zhang J, Kobert K, Flouri T, Stamatakis A. 2014. PEAR: a fast and accurate Illumina Paired-  
674 End reAd mergeR. *Bioinformatics* 30:614–620.
- 675 30. Rivers AR, Weber KC, Gardner TG, Liu S, Armstrong SD. 2018. ITSxpress: Software to  
676 rapidly trim internally transcribed spacer sequences with quality scores for marker gene  
677 analysis. *F1000Res* 7:1418.
- 678 31. Bolyen E, Rideout JR, Dillon MR, Bokulich NA, Abnet CC, Al-Ghalith GA, Alexander H,  
679 Alm EJ, Arumugam M, Asnicar F, Bai Y, Bisanz JE, Bittinger K, Brejnrod A, Brislawn CJ,  
680 Brown CT, Callahan BJ, Caraballo-Rodríguez AM, Chase J, Cope EK, Da Silva R, Diener C,  
681 Dorrestein PC, Douglas GM, Durall DM, Duvallet C, Edwardson CF, Ernst M, Estaki M,  
682 Fouquier J, Gauglitz JM, Gibbons SM, Gibson DL, Gonzalez A, Gorlick K, Guo J, Hillmann  
683 B, Holmes S, Holste H, Huttenhower C, Huttley GA, Janssen S, Jarmusch AK, Jiang L,  
684 Kaehler BD, Kang KB, Keefe CR, Keim P, Kelley ST, Knights D, Koester I, Kosciolk T,  
685 Kreps J, Langille MGI, Lee J, Ley R, Liu Y-X, Loftfield E, Lozupone C, Maher M, Marotz C,  
686 Martin BD, McDonald D, McIver LJ, Melnik AV, Metcalf JL, Morgan SC, Morton JT,  
687 Naimey AT, Navas-Molina JA, Nothias LF, Orchanian SB, Pearson T, Peoples SL, Petras D,  
688 Preuss ML, Pruesse E, Rasmussen LB, Rivers A, Robeson MS, Rosenthal P, Segata N,  
689 Shaffer M, Shiffer A, Sinha R, Song SJ, Spear JR, Swofford AD, Thompson LR, Torres PJ,  
690 Trinh P, Tripathi A, Turnbaugh PJ, Ul-Hasan S, van der Hooft JJJ, Vargas F, Vázquez-Baeza  
691 Y, Vogtmann E, von Hippel M, Walters W, Wan Y, Wang M, Warren J, Weber KC,  
692 Williamson CHD, Willis AD, Xu ZZ, Zaneveld JR, Zhang Y, Zhu Q, Knight R, Caporaso JG.  
693 2019. Reproducible, interactive, scalable and extensible microbiome data science using  
694 QIIME 2. *Nat Biotechnol* 37:852–857.

- 695 32. Callahan BJ, McMurdie PJ, Rosen MJ, Han AW, Johnson AJA, Holmes SP. 2016. DADA2:  
696 High-resolution sample inference from Illumina amplicon data. *Nat Methods* 13:581–583.
- 697 33. Abarenkov K, Nilsson RH, Larsson K-H, Taylor AFS, May TW, Frøslev TG, Pawlowska J,  
698 Lindahl B, Pöldmaa K, Truong C, Vu D, Hosoya T, Niskanen T, Piirmann T, Ivanov F, Zirk  
699 A, Peterson M, Cheeke TE, Ishigami Y, Jansson AT, Jeppesen TS, Kristiansson E,  
700 Mikryukov V, Miller JT, Oono R, Ossandon FJ, Paupério J, Saar I, Schigel D, Suija A,  
701 Tedersoo L, Kõljalg U. 2024. The UNITE database for molecular identification and  
702 taxonomic communication of fungi and other eukaryotes: sequences, taxa and classifications  
703 reconsidered. *Nucleic Acids Research* 52:D791–D797.
- 704 34. Bengtsson-Palme J, Ryberg M, Hartmann M, Branco S, Wang Z, Godhe A, De Wit P,  
705 Sánchez-García M, Ebersberger I, de Sousa F, Amend A, Jumpponen A, Unterseher M,  
706 Kristiansson E, Abarenkov K, Bertrand YJK, Sanli K, Eriksson KM, Vik U, Veldre V,  
707 Nilsson RH. 2013. Improved software detection and extraction of ITS1 and ITS2 from  
708 ribosomal ITS sequences of fungi and other eukaryotes for analysis of environmental  
709 sequencing data. *Methods in Ecology and Evolution* 4:914–919.
- 710 35. Lahti L, Shetty SA. 2017. Introduction to the microbiome R package.
- 711 36. Mcmurdie PJ, Holmes S. 2013. phyloseq: An R Package for Reproducible Interactive  
712 Analysis and Graphics of Microbiome Census Data. *PLoS ONE* 8:e61217.
- 713 37. Oksanen J, Simpson GL, Blanchet FG, Kindt R, Legendre P, Minchin PR, O’Hara RB,  
714 Solymos P, Stevens MHH, Szoecs E, Wagner H, Barbour M, Bedward M, Bolker B, Borcard  
715 D, Borman T, Carvalho G, Chirico M, Caceres MD, Durand S, Evangelista HBA, FitzJohn R,  
716 Friendly M, Furneaux B, Hannigan G, Hill MO, Lahti L, Martino C, McGlenn D, Ouellette  
717 M-H, Cunha ER, Smith T, Stier A, Braak CJFT, Weedon J. 2025. *vegan: Community  
718 Ecology Package (2.7-1)*.
- 719 38. Castaño C, Berlin A, Brandström Durling M, Ihrmark K, Lindahl BD, Stenlid J,  
720 Clemmensen KE, Olson Å. 2020. Optimized metabarcoding with Pacific biosciences enables  
721 semi-quantitative analysis of fungal communities. *New Phytologist* 228:1149–1158.
- 722 39. Mbareche H, Veillette M, Bilodeau G, Duchaine C. 2020. Comparison of the performance of  
723 ITS1 and ITS2 as barcodes in amplicon-based sequencing of bioaerosols. *PeerJ* 8:e8523.
- 724 40. Biada I, Santacreu MA, González-Recio O, Ibáñez-Escriche N. 2025. Comparative analysis  
725 of Illumina, PacBio, and nanopore for 16S rRNA gene sequencing of rabbit’s gut microbiota.  
726 *Front Microbiomes* 4.
- 727 41. Wagner J, Paul Coupland, Browne HP, Lawley TD, Parkhill J. 2016. Evaluation of PacBio  
728 sequencing for full-length bacterial 16S rRNA gene classification. *BMC Microbiol* 16:274.
- 729 42. Estensmo ELF, Morgado L, Maurice S, Martin-Sanchez PM, Engh IB, Mattsson J, Kauserud  
730 H, Skrede I. 2021. Spatiotemporal variation of the indoor mycobiome in daycare centers.  
731 *Microbiome* 9:220.
- 732 43. Tedersoo L, Bahram M, Zinger L, Nilsson RH, Kennedy PG, Yang T, Anslan S, Mikryukov  
733 V. 2022. Best practices in metabarcoding of fungi: From experimental design to results.  
734 *Molecular Ecology* 31:2769–2795.

- 735 44. Hoggard M, Vesty A, Wong G, Montgomery JM, Fourie C, Douglas RG, Biswas K, Taylor  
736 MW. 2018. Characterizing the Human Mycobiota: A Comparison of Small Subunit rRNA,  
737 ITS1, ITS2, and Large Subunit rRNA Genomic Targets. *Front Microbiol* 9.
- 738 45. Brown GD, Ballou ER, Bates S, Bignell EM, Borman AM, Brand AC, Brown AJP, Coelho  
739 C, Cook PC, Farrer RA, Govender NP, Gow NAR, Hope W, Hoving JC, Dangarembizi R,  
740 Harrison TS, Johnson EM, Mukaremera L, Ramsdale M, Thornton CR, Usher J, Warris A,  
741 Wilson D. 2024. The pathobiology of human fungal infections. *Nat Rev Microbiol* 22:687–  
742 704.
- 743 46. Camilli G, Tabouret G, Quintin J. 2018. The Complexity of Fungal  $\beta$ -Glucan in Health and  
744 Disease: Effects on the Mononuclear Phagocyte System. *Front Immunol* 9.
- 745 47. Ráduly Z, Szabó L, Madar A, Pócsi I, Csernoch L. 2020. Toxicological and Medical Aspects  
746 of *Aspergillus*-Derived Mycotoxins Entering the Feed and Food Chain. *Front Microbiol* 10.
- 747 48. Nash J, Tremble K, Schadt C, Cregger MA, Bryan C, Vilgalys R. 2025. Time-series RNA  
748 metabarcoding of the active *Populus tremuloides* root microbiome reveals hidden temporal  
749 dynamics and dormant core members. *mSystems* 10:e00285-25.
- 750 49. Nguyen NH, Song Z, Bates ST, Branco S, Tedersoo L, Menke J, Schilling JS, Kennedy PG.  
751 2016. FUNGuild: An open annotation tool for parsing fungal community datasets by  
752 ecological guild. *Fungal Ecology* 20:241–248.

753

754

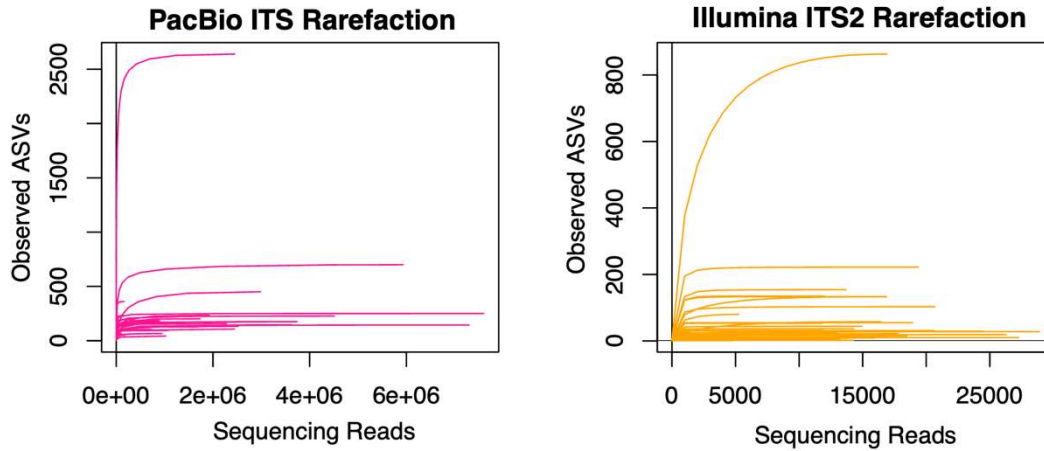
755 **Supplementary Figures**



756

757 Fig. S1 Images of sampling locations. (A) Home 1, an unoccupied residence in Carrboro, NC.  
758 (B) Home 2, a testbed residence for indoor environmental research at Duke University.  
759 (C) A student dormitory room at the UNC Institute of Marine Sciences (IMS). (D) Laboratory ceiling  
760 at the University of North Carolina at Chapel Hill (UNC), showing mold growth around an air  
761 vent. (E) Walk-in cold room at UNC with visible mold contamination on ceiling light diffuser  
762 panels.

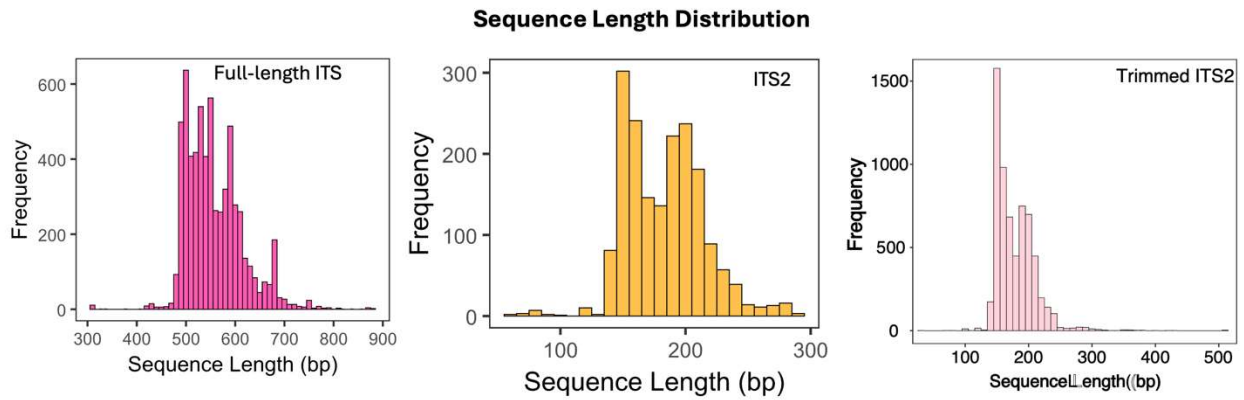
763



764

765 Fig S2. Rarefaction curves of full-length ITS (left) dataset and ITS2 (right) dataset.

766

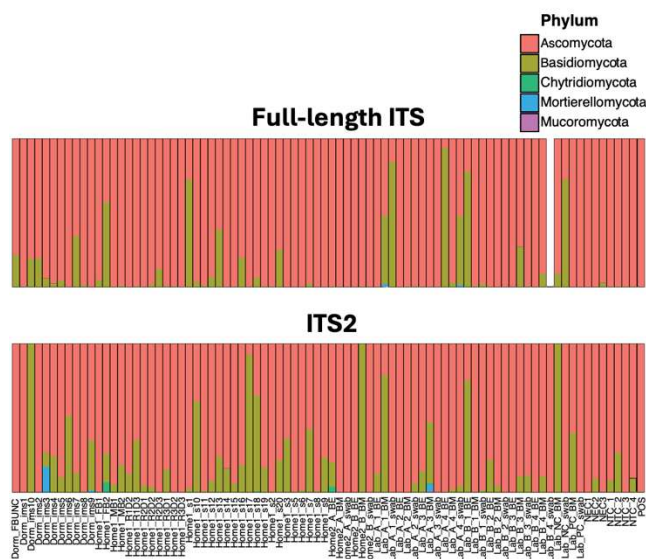


767

768 Fig S3. Sequence length distributions demonstrating the expected difference between full-length  
 769 ITS sequences (upper panel, 400-800 bp), ITS2 fragments (middle panel, 100-300 bp), and in-  
 770 silico trimmed ITS2 fragments (bottom panel, 100-300 bp).

771

### Relative Abundance at Phylum Level

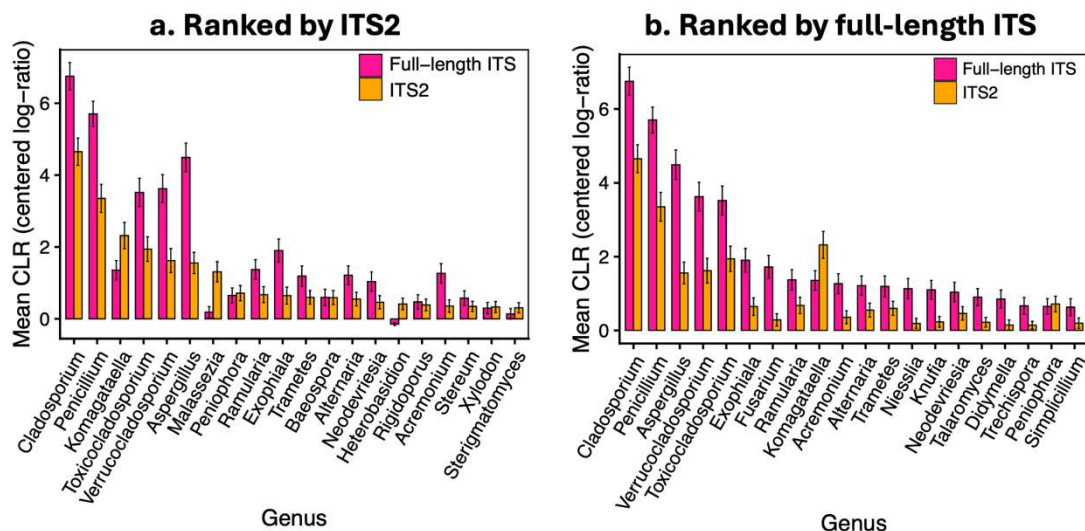


772  
773 Fig. S4. Relative abundance at phylum level derived from the full-length ITS pipeline (top panel)  
774 and ITS2 pipeline (bottom panel).

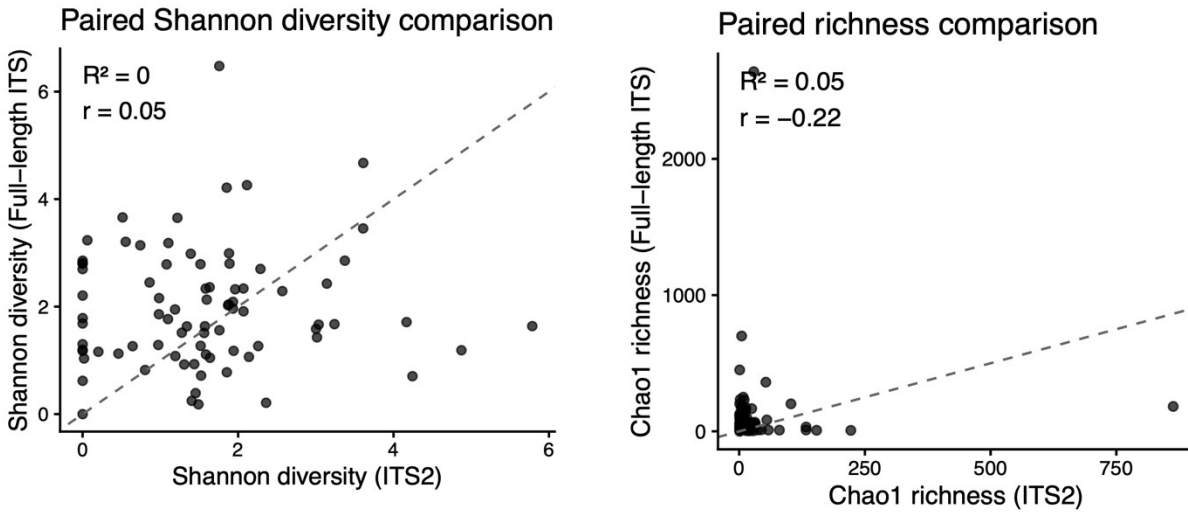
775

776

### Top 20 most abundant genera (Mean CLR $\pm$ SE)



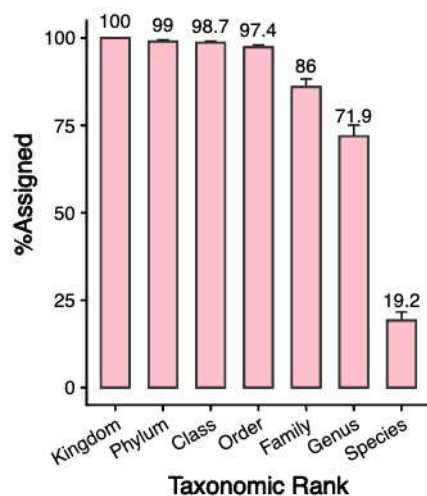
777  
778 Fig. S5. Top 20 genera ranked within each pipeline. (a) Top 20 by ITS2 mean CLR abundance;  
779 (b) Top 20 by Full-length ITS mean CLR abundance. Bars show mean CLR  $\pm$  SE across  
780 samples. These panels illustrate method-specific rank shifts that are partially obscured when  
781 genera are ranked by combined abundance.



782

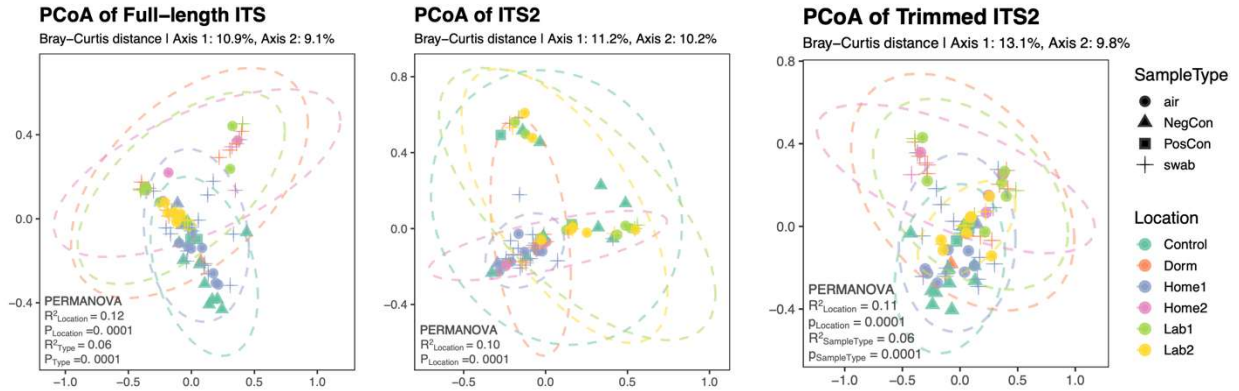
783 Fig. S6. **Paired comparisons of alpha-diversity metrics between sequencing pipelines.** Paired  
 784 scatter plots compare per-sample **Shannon diversity** (left) and **Chao1 richness** (right)  
 785 estimated from Illumina ITS2 and PacBio full-length ITS datasets. Each point represents a matched sample  
 786 processed by both pipelines, with dashed lines indicating the 1:1 relationship. Spearman  
 787 correlation was used for both indexes. Shannon diversity showed minimal sample-level  
 788 correlation despite similar group-level distributions, whereas richness estimates exhibited poor  
 789 concordance, reflecting systematic differences in rare-taxon resolution across markers. These  
 790 patterns indicate that alpha-diversity indices, particularly richness, are not directly  
 791 interchangeable between ITS2 and full-length ITS pipelines.

792



793

794 Fig. S7. Taxonomic assignment across all ranks comparing of in-silico trimmed ITS2 sequences  
 795 by paired samples. Numbers above the bars indicate percent of ASV classified  $\pm$  standard error  
 796 (SE) for each rank.



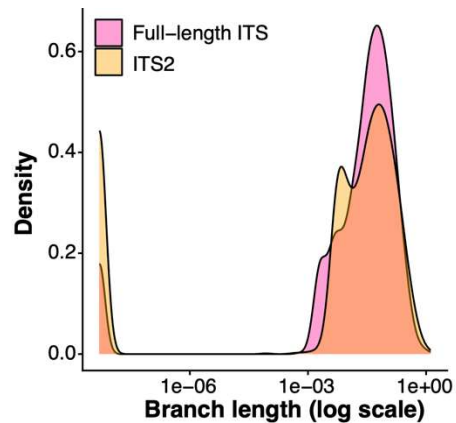
Metrics	ITS2 vs. Full-length ITS	ITS2 vs. Trimmed ITS2
Mantel on Bray-Curtis	$R = 0.06$ , $p = 0.098$	$R = 0.03$ , $p = 0.273$
Procrustes (symmetric)	$R = 0.67$ (ss = 0.55), $p = 0.001$	$R = 0.63$ (ss = 0.60), $p = 0.001$

797

798 **Fig. S8. Comparison of full-length ITS, ITS2, and trimmed ITS2 ordinations based on**  
 799 **Bray-Curtis dissimilarity.** Principal coordinates analysis (PCoA) plots show fungal community  
 800 composition derived from **full-length ITS** (left), **Illumina ITS2** (middle), and **PacBio-derived**  
 801 **in-silico ITS2** (right). Ellipses represent 95% confidence intervals for location groups.  
 802 Full-length ITS showed the strongest location-level separation (PERMANOVA  $R^2 = 0.12$ ,  
 803  $p = 10^{-4}$ ), whereas both ITS2 and trimmed ITS2 displayed weaker but significant location  
 804 structure ( $R^2 = 0.10$ ,  $p = 10^{-4}$ ; and  $R^2 = 0.11$ ,  $p = 10^{-4}$ , respectively). Mantel tests comparing  
 805 Bray-Curtis distance matrices indicated weak and non-significant correlations between ITS2 and  
 806 full-length ITS ( $r = 0.06$ ,  $p = 0.098$ ) as well as between ITS2 and trimmed ITS2 ( $r = 0.03$ ,  
 807  $p = 0.273$ ). In contrast, symmetric Procrustes analysis showed moderate and significant  
 808 alignment of ordination configurations for both comparisons (ITS2 vs full-length ITS:  $r = 0.67$ ,  
 809 sum of squares = 0.55,  $p = 0.001$ ; ITS2 vs trimmed ITS2:  $r = 0.63$ , sum of squares = 0.60,  
 810  $p = 0.001$ ). A robustness analysis restricted to the first 15 PCoA axes yielded lower but still  
 811 significant alignment ( $r \approx 0.40$ ,  $p = 0.001$ ), indicating that shared community structure is  
 812 distributed across multiple ordination dimensions.

813

814



Metrics	ITS2	Full-length ITS
True Polytomies	1	1
Percent internal nodes	0.73%	0.97%
Total tree length	187.86	139.02
median branch length	0.0266	0.0367
Wilcoxon test	$p = 1.64 \times 10^{-8}$	

815

816 **Fig. S9. Branch length distributions for ITS2 and full-length ITS phylogenies.** Kernel  
 817 density plots show the distribution of branch lengths inferred from phylogenies constructed using  
 818 ITS2 sequences (orange) and full-length ITS sequences (deep pink), plotted on a log scale. Both  
 819 phylogenies were largely bifurcating, with a similarly low proportion of multifurcating internal  
 820 nodes (polytomies; <0.1% of internal nodes in both trees). In contrast, branch lengths were  
 821 significantly shorter in the ITS2-derived tree than in the full-length ITS tree (median branch  
 822 length: 0.0266 vs. 0.0367; Wilcoxon rank-sum test,  $p = 1.64 \times 10^{-8}$ ), indicating reduced  
 823 phylogenetic signal in the ITS2 region. Longer branches in the full-length ITS phylogeny reflect  
 824 greater sequence variation captured across the combined ITS1, 5.8S, and ITS2 regions. Although  
 825 the ITS2 phylogeny exhibited greater total tree length, this reflected the accumulation of many  
 826 short branches rather than increased phylogenetic resolution. In contrast, longer median branch  
 827 lengths in the full-length ITS tree indicate greater information content per split.

828

829 **Supplementary Tables**830 **Table S1. Samples summary**

<b>Location</b>	<b>Sample (n)</b>	<b>Illumina ITS2 reads</b>	<b>PacBio full-length ITS reads</b>
Home1 (n = 31)	Surface (19)	195,294	15,429,812
	Air (8)	80,283	4,203,379
	NegCon-Field Blank (2)	43,224	3,613
	NegCon-Method Blank (2)	16,903	3,640
Home2 (n = 6)	Surface (2)	26,610	4,727,465
	Air (4)	45,776	8,336,048
Dorm (n = 11)	Surface (10)	114,889	1,937,575
	NegCon-Field Blank (1)	20,688	10,311
Lab 1 (n = 12)	Surface (4)	40,678	243,014
	Air (8)	78,524	209,857
	NegCon-Method Blank (2)	2,341	45,176
Lab 2 (n = 12)	Surface (4)	43,248	8,779
	Air (8)	111,713	9,077,577
Pipeline Control (n = 12)	Positive Control (3)	41,957	4,558,955
	NegCon -No Template (10)	42,207	14,810,490

831 \*Reads reflect total raw reads in corresponding categories and were not adjusted for polling factors.

832 NegCon stands for Negative Control.

833 Table S2. Primers and sequences used in Illumina ITS2 pipeline (5' to 3')(24–26).

Name	Sequence
ITS3NGS1	CATCGATGAAGAACGCAG
ITS3NGS2	CAACGATGAAGAACGCAG
ITS3NGS3	CACCGATGAAGAACGCAG
ITS3NGS4	CATCGATGAAGAACGTAG
ITS3NGS5	CATCGATGAAGAACGTGG
ITS3NGS10	CATCGATGAAGAACGCTG TCGTCGGCAGCGTCAGATGTGTATAAAGAGACAGNNNNNNNN TT
ITS3NGS1-F1	CATCGATGAAGAACGCAG TCGTCGGCAGCGTCAGATGTGTATAAAGAGACAGNNNNTNNNN TT
ITS3NGS1-F2	CATCGATGAAGAACGCAG TCGTCGGCAGCGTCAGATGTGTATAAAGAGACAGNNNNCTNNNN TT
ITS3NGS1-F3	CATCGATGAAGAACGCAG TCGTCGGCAGCGTCAGATGTGTATAAAGAGACAGNNNNACTNNNN TT
ITS3NGS1-F4	CATCGATGAAGAACGCAG TCGTCGGCAGCGTCAGATGTGTATAAAGAGACAGNNNNGACTNNNN TT
ITS3NGS1-F5	CATCGATGAAGAACGCAG TCGTCGGCAGCGTCAGATGTGTATAAAGAGACAGNNNNTGACTNNNN TT
ITS3NGS1-F6	CATCGATGAAGAACGCAG TCGTCGGCAGCGTCAGATGTGTATAAAGAGACAGNNNNNNNN TT
ITS3NGS2-F1	CAACGATGAAGAACGCAG TCGTCGGCAGCGTCAGATGTGTATAAAGAGACAGNNNNTNNNN TT
ITS3NGS3-F2	CACCGATGAAGAACGCAG TCGTCGGCAGCGTCAGATGTGTATAAAGAGACAGNNNNCTNNNN TT
ITS3NGS4-F3	CATCGATGAAGAACGTAG TCGTCGGCAGCGTCAGATGTGTATAAAGAGACAGNNNNACTNNNN TT
ITS3NGS5-F4	CATCGATGAAGAACGTGG TCGTCGGCAGCGTCAGATGTGTATAAAGAGACAGNNNNGACTNNNN TT
ITS3NGS10-F5	CATCGATGAAGAACGCTG GTCTCGTGGGCTCGGAGATGTGTATAAAGAGACAGNNNNN GA
ITS4NGR-F1	TCCTSCGCTTATTGATATGC GTCTCGTGGGCTCGGAGATGTGTATAAAGAGACAGNNTNNN GA
ITS4NGR-F2	TCCTSCGCTTATTGATATGC GTCTCGTGGGCTCGGAGATGTGTATAAAGAGACAGNNTNNN GA
ITS4NGR-F3	TCCTSCGCTTATTGATATGC GTCTCGTGGGCTCGGAGATGTGTATAAAGAGACAGNNTNNN GA
ITS4NGR-F4	TCCTSCGCTTATTGATATGC GTCTCGTGGGCTCGGAGATGTGTATAAAGAGACAGNNGACTNNN GA
ITS4NGR-F5	TCCTSCGCTTATTGATATGC GTCTCGTGGGCTCGGAGATGTGTATAAAGAGACAGNNTGACTNNN GA
ITS4NGR-F6	TCCTSCGCTTATTGATATGC GTCTCGTGGGCTCGGAGATGTGTATAAAGAGACAGNNNNN GA
ARCH-ITS4	TCCTCGCTTATTGATATGC

834

835

836 Table S3. Heatmap (Fig. 3C) environmental health-related genera statistics

Detected Targeted Genera	Class	Wilcoxon test adjusted p-value	
		Overall	Subsets
Alternaria	Airborne allergen	0.072	
Aspergillus	Opportunistic pathogen	0.00003	Air, 0.008 Swab, 0.04
Aureobasidium	Opportunistic pathogen	0.072	
Bipolaris	Plant pathogen	0.087	
Bjerkandera	Wood-decay	0.389	
Candida	Pathogen	0.039	
Cladosporium	Allergen	0.004	
Cryptococcus	Pathogen	0.056	
Epicoccum	Indoor fungus	0.917	
Fusarium	Plant pathogen	0.0004	Swab, 0.03
Malassezia	Pathogen	0.002	Air, 0.04
Mucor	Pathogen	0.572	
Penicillium	Pathogen	0.0008	Air, 0.04
Phoma	Plant pathogen	0.057	
Rhodotorula	Opportunistic pathogen	0.057	
Saccharomyces	Opportunistic pathogen	0.015	
Stemphylium	Plant pathogen	0.057	
Trichoderma	Opportunistic pathogen	0.146	
Wallemia	Allergen	0.052	

837 \*Subsets are comparing between pipelines within select subsets of sample type (Air, Swab,  
838 Negative Control, Positive Control) or location (Home1, Home2, Lab 1, Lab 2, Dorm). Subsets  
839 with significant differences were listed.

840

## 841 **Supplementary Method**

### 842 *Sample collection and preprocessing*

#### 843 *Home 1 and Dorm*

844 At Home 1 and Dorm, surface swabs were collected from a variety of surfaces using Isohelix™  
845 SK-3S rayon swabs pre-moistened with 50 mM Tris buffer. Swabbing was performed over a 110  
846 cm<sup>2</sup> area using standardized multidirectional S-strokes for 2 minutes. To account for potential  
847 contamination during field handling, field blanks were collected at both locations: two at Home 1  
848 and one at Dorm. For each field blank, a sterile swab was opened in the field, briefly exposed to  
849 ambient air for 5 seconds, and placed directly into elution buffer without contacting any surface.  
850 Following sample and field blank collection, all swabs were transferred into an elution buffer (20  
851 mM Tris, pH 8.0, 1% polyvinylpyrrolidone [PVP], 1% Tween 20) for transport. In the  
852 laboratory, swabs were briefly vortexed, and the resulting eluates were collected and aliquoted  
853 for downstream processing. All eluates were stored at -80 °C for no more than 6 months prior to  
854 analysis. At Home 1, bioaerosol samples (n = 8) were collected from three bedrooms using a  
855 BobCat (BC) AC-200 Sampler (InnovaPrep, Drexel, MO, USA), operated at 200 L/min for 2  
856 hours with a dry electret filter. After sampling, the filter was eluted with 8 mL of wet foam  
857 (AC08100T, InnovaPrep). To control for potential contamination introduced during the aerosol  
858 collection and processing workflow, two method blanks were prepared by passing sterile wet  
859 foam through clean dry electret filters using the same procedure as for bioaerosol samples, but  
860 without field deployment. All eluates, including those from method blanks, were aliquoted and  
861 stored at -80 °C.

#### 862 *Home 2 and Lab*

863 At Home 2 and Lab (1 and 2), surface swabs were collected using Isohelix™ SK-3S rayon swabs  
864 moistened with QIAGEN Buffer AE (10 mM Tris-Cl; 0.5 mM EDTA; pH 9.0). At Home 2,  
865 swabs were collected from a basement wall and an entryway wall. At UNC, swabs were taken  
866 from visibly mold-affected locations, including an air vent in a laboratory (Lab\_1) and a light  
867 diffuser panel in a walk-in cold room (Lab\_2). To monitor for potential contamination and assess  
868 processing consistency, one negative swab control (NC\_swab) was included by exposing a  
869 sterile swab to ambient air in the laboratory for 5 seconds before placing it directly into QIAGEN  
870 Buffer AE without surface contact. One positive swab control (PC\_swab) was prepared by  
871 spiking a sterile swab with *Aspergillus niger*. Following sample and control preparation, all  
872 swabs were transferred to QIAGEN Buffer AE for transport. In the laboratory, swabs were  
873 briefly vortexed, and the resulting eluates were collected and aliquoted for downstream  
874 processing. All eluates were stored at -80 °C for no more than 6 months prior to analysis.  
875 Bioaerosol samples were collected at both Home 2 and UNC using two complementary methods:  
876 (1) the BobCat (BC) as described above, operated at 200 L/min for 5 minutes; and (2) filter  
877 concentration, in which a 1 mL aliquot of the BC eluate was passed through a 47 mm, 0.45 µm  
878 mixed cellulose ester (MCE) membrane filter (Millipore HAWP 04700) to enhance particle  
879 retention. One negative bioaerosol control (NC\_BM) was prepared by passing sterile wet foam  
880 through a clean dry electret filter without field deployment. One positive bioaerosol control  
881 (PC\_BM) was prepared by spiking a sterile BobCat eluate with *Aspergillus niger* derived from  
882 ATCC 6275 (KwikStik). All eluates, including those from negative and positive controls, were  
883 aliquoted and stored at -80 °C until processing.

884

885 *Illumina ITS2 library preparation PCR conditions*

886 Fungal ITS2 regions were amplified using a three-step PCR protocol adapted from Lundberg et  
887 al. (23). Initial amplification of ITS2 was performed on 3 $\mu$ L of extracted total DNA in a 10- $\mu$ L  
888 volume containing 0.3  $\mu$ L of 10  $\mu$ M forward (ITS3NGS) and reverse(ITS4NGR) primers (24–  
889 26), 5  $\mu$ L of 2X KAPA HiFi HotStart ReadyMix (F. Hoffmann-La Roche AG, Basel,  
890 Switzerland), 0.1  $\mu$ L of 100X SYBR Green I (Thermo Fisher Scientific, Waltham, MA, USA),  
891 and 1.3  $\mu$ L nuclease-free water. Reaction mixes were denatured for 3 min at 95°C prior to 30  
892 cycles of denaturation at 98°C for 20 seconds, annealing at 58°C for 15 seconds, and extension at  
893 72°C for 1 min, followed by immediately cooling reactions to 4°C until transferred to ice. A  
894 second PCR introduced Illumina sequencing adaptors using frameshifted primers (ITS3NGS F1–  
895 F6, ITS4NGR F1–F6). 3 $\mu$ L of PCR product from first PCR was added to PCR master mix with  
896 same recipe as the previous PCR and ran at the same thermocycle conditions except for only  
897 repeat 10 cycles of denaturation-annealing-extension. Third PCR amplification to add Illumina  
898 adapters and dual 8-bp indices for sample multiplexing was performed in a 50  $\mu$ l volume  
899 containing 5  $\mu$ l of 2.5  $\mu$ M forward and reverse indexing primers, 25  $\mu$ l of 2X KAPA HiFi buffer,  
900 0.5  $\mu$ l of 100X SYBR Green I, 9.5  $\mu$ l nuclease free water, and 5  $\mu$ l of PCR product from previous  
901 step. Reaction mixes were denatured for 3 min at 95°C prior to 30 cycles of denaturation at 98°C  
902 for 20 seconds, annealing at 60°C for 15 seconds, and extension at 72°C for 1 min, followed by  
903 immediately cooling reactions to 4°C until transferred to ice. All primer sequences are listed in  
904 S. Table 2.

905 *PacBio full-length ITS library preparation*

906 Full-length ITS library was prepared using the same extracted DNA from ITS2 pipeline. Five  $\mu$ L  
907 of extracted total DNA was subject to PCR amplification of the ITS region using Phusion Plus  
908 PCR Master Mix (Thermo Fisher Scientific, Waltham, MA, USA) with forward primer  
909 ITS1catta (5'-ACCWGC GGARGGATCATTA-3') and reverse primer ITS4ngsUni (5'-  
910 CCTSCSCTTANTDATATGC-3') containing unique barcodes and Kinnex adaptors at a final  
911 concentration of 0.3 $\mu$ M(20, 27). Reaction mixes were denatured for 30 seconds at 98°C prior to  
912 35 cycles of denaturation at 98°C for 10 seconds, annealing at 57°C for 20 seconds, and  
913 extension at 72°C for 75 seconds. After the last cycle, a final extension at 72°C for 5 minutes  
914 was performed. Completed PCR reactions were visualized on an Invitrogen E-Gel™ EX Agarose  
915 Gels (Thermo Fisher Scientific, Waltham, MA, USA) to ensure that amplicon size was correct  
916 (~800bp), and that each sample amplified appropriately. Amplicon libraries were subsequently  
917 pooled based on gel band intensity (1.5, 5, 10, or 20 $\mu$ L per reaction). Library pool was cleaned  
918 and concentrated using 1.4x volume of SMRTbell® cleanup beads (Pacific Biosciences of  
919 California, Inc., Menlo Park, CA, USA) and eluted in 50 $\mu$ l of Low TE Elution Buffer (PacBio).  
920 Cleaned libraries were quantified using the Qubit HS dsDNA kit (Thermo Fisher Scientific,  
921 Waltham, MA, USA) and stored at -20°C prior to Kinnex PCR for concatenation and  
922 circularization. Kinnex PCR was performed as outlined in the PacBio Kinnex 16S kit's published  
923 protocol with no modifications (PacBio)(28). Size selected and cleaned libraries were loaded  
924 onto a PacBio SMRT® Cell and sequenced on the Revio system (PacBio) in the Duke  
925 Sequencing and Genomic Technologies Shared Resource.

## 926 *Illumina ITS2 Sequence Processing*

927 Raw Illumina paired-end reads were processed using a multi-step pipeline combining several  
928 bioinformatic tools. First, paired-end reads were merged using PEAR (29) with default  
929 parameters to reconstruct full-length ITS2 amplicons. The ITS2 region was then extracted from  
930 merged reads using ITSxpress (30) with the following parameters: --region ITS2, --taxa Fungi, --  
931 single\_end, and --threads 64. This step ensured that only the ITS2 region was retained for  
932 downstream analysis, removing flanking conserved regions.

933 Subsequent processing was performed using QIIME 2 (version 2023.9) (31) within a Singularity  
934 container environment. A manifest file was generated to import the ITSxpress-processed  
935 sequences into QIIME 2 as single-end data using the SingleEndFastqManifestPhred33V2 format.  
936 Quality control and amplicon sequence variant (ASV) inference were conducted using the  
937 DADA2 plugin with the denoise-single workflow (32). No length truncation was applied (--p-  
938 trunc-len 0) since sequences had already been quality processed through ITSxpress. The  
939 maximum expected error threshold was set to the default value of 2 (--p-max-ee 2.0). The  
940 DADA2 algorithm performed error correction, denoising, and chimera removal to generate high-  
941 resolution ASVs.

942 Taxonomic classification was performed using a pre-trained naive Bayes classifier based on the  
943 UNITE database (version 9.0, July 2023) obtained from the unite-train repository (33). The  
944 sklearn-based classifier was applied to representative ASV sequences using the classify-sklearn  
945 function with parallel processing enabled (--p-n-jobs -1). Final outputs including the feature  
946 table, representative sequences, and taxonomic assignments were exported from QIIME 2  
947 artifacts using qiime tools export and converted to standard formats (TSV, FASTA) for  
948 downstream analysis in R.

## 949 *PacBio full-length ITS Sequence Processing*

950 PacBio circular consensus sequences (CCS) were processed using the DADA2 pipeline (version  
951 1.28) with PacBio-specific modifications (32). Primer sequences (ITS1catta: 5'-  
952 ACCWGC GGARGGATCATTA-3' and ITS4ngsUni: 5'-CCTSCSCTTANTDATATGC-3') were  
953 removed using the removePrimers function with orientation correction enabled. Quality filtering  
954 and trimming were performed using filterAndTrim with the following parameters: minimum  
955 quality score of 3, minimum length of 250 bp, maximum length of 1200 bp, maximum expected  
956 errors of 3, and complete removal of sequences containing ambiguous nucleotides (maxN=0).

957 Error modeling was conducted using the PacBioErrfun function specifically designed for PacBio  
958 data, which accounts for the distinct error profile of long-read sequencing. Sequence variants  
959 were inferred using the dada function with PacBio-optimized error rates. Chimeric sequences  
960 were identified and removed using removeBimeraDenovo with the consensus method and a  
961 minimum fold-parent-over-abundance threshold of 3.5, which is optimized for longer amplicons.

962 Taxonomy was assigned using the same classifier based on the same UNITE database (version  
963 9.0, July 2023) and sklearn-based classifier as Illumina short-read ITS2 sequences (33). The final

964 amplicon sequence variant (ASV) table retained sequences ranging from 250-1200 bp,  
965 encompassing the full-length ITS region including ITS1, 5.8S rRNA gene, and ITS2.

### 966 *Cross-platform ASV Mapping*

967 The full-length ITS sequences and taxonomic assignments from the PacBio pipeline were  
968 extracted from the phyloseq object and exported to FASTA (`pacbio_refseq.fasta`) and tab-  
969 delimited taxonomy files (`pacbio_taxonomy.tsv`). To annotate Illumina-derived ITS2 amplicon  
970 sequence variants (ASVs), we used a direct sequence similarity search strategy. First, a  
971 nucleotide BLAST database was constructed from the PacBio reference sequences using  
972 `makeblastdb` from the NCBI BLAST+ suite. Illumina ITS2 representative sequences were then  
973 queried against this database using `blastn`, with a minimum sequence identity threshold of 97%  
974 and an e-value cutoff of  $1e^{-20}$ . The resulting alignment table was annotated with taxonomic  
975 information by joining the subject sequence identifiers (`sseqid`) with the PacBio taxonomy table  
976 using a custom R script. Final outputs included both the raw alignment metrics and the  
977 corresponding taxonomic assignments. We observed that several Illumina ASVs matched  
978 multiple full-length PacBio sequences, reflecting the limited resolving power of ITS2.

### 979 *In-silico ITS2 Extraction from PacBio Full-Length ITS Reads*

980 An in-silico ITS2 dataset was generated from the PacBio full-length ITS sequences  
981 (`pacbio_refseq.fasta`). The ITS2 region was extracted using `ITSx` with the `--only ITS2` option  
982 to restrict output to the ITS2 subregion, and `--taxa F` to specify fungal ITS boundaries (34). The `-`  
983 `-preserve` flag retained the original ASV identifiers, and the command generated the file  
984 `pacbio_ITSx.ITS2.fasta`. The resulting PacBio-ITS2 sequences were then processed and  
985 assigned with the same Illumina ITS2 pipeline to ensure classifier parity and identical  
986 QC/denoising settings. This PacBio-ITS2 dataset was used in sensitivity analyses to (i) compare  
987 rank-assignment rates to Illumina ITS2 under identical bioinformatic conditions, and (ii) quantify  
988 how much of the Full-length ITS vs ITS2 difference could be attributed to amplicon region and  
989 primer targeting rather than platform chemistry.

### 990 *Statistical analysis*

991 All statistical analyses and data visualizations were performed in R Studio Version 2024.04.2  
992 (Posit, PBC, Boston, MA, USA). ITS2 sequencing depth sufficiency was evaluated by (i)  
993 rarefaction curve using function `rarecurve` and (ii) computing Good's coverage per sample ( $1 -$   
994  $\text{singletons}/\text{total reads}$ ). Before formal data analysis, functional guilds were added to fungal taxa  
995 in phyloseq objects using FUNGuild (49) with a confidence ranking of "Probable" or "Highly  
996 Probable". Taxa were categorized into ecological guilds (saprotrophs, plant pathogens, animal  
997 pathogens, etc.) and trophic modes. To handle differences in sequencing depth, taxa abundance  
998 data underwent centered-log ratio (CLR) transformation using function `transform` from the  
999 `microbiome` package (35) to address the compositional nature of microbiome data. For  
1000 visualization purposes, relative abundance transformations were also performed using the same  
1001 package. Taxonomic data were agglomerated at various taxonomic ranks using the `tax_glom`  
1002 function in phyloseq when appropriate.

1003 Alpha diversity metrics including observed richness and Shannon diversity index were calculated  
1004 using the vegan (37) package. Beta diversity was assessed using Bray-Curtis distances on  
1005 relative abundance transformed data, ordination plot was using the Principal Coordinate Analysis  
1006 (PCoA) in phyloseq package. Differences in fungal community composition between sample  
1007 groups were tested using permutational multivariate analysis of variance (PERMANOVA) with  
1008 the adonis2 function from vegan (37), using 10,000 permutations and Euclidean distances on  
1009 CLR-transformed data. Ordination agreement was evaluated using symmetric  
1010 Procrustes/PROTEST (vegan), reporting the Procrustes correlation  $r = \sqrt{1 - ss}$  and  
1011 permutation p-value, and Mantel (Spearman) on Bray-Curtis distance matrices. Analyses used  
1012 matched samples and 999 permutations.

1013 The taxonomic assignment success between Illumina and PacBio platforms was compared by  
1014 calculating the percentage of ASVs assigned at each taxonomic level from kingdom through  
1015 species. Correlations between platforms were assessed using Spearman's rank correlation  
1016 coefficients calculated on CLR-transformed abundance data at the genus level. Paired Wilcoxon  
1017 tests -tests were performed to test differences between sequencing methods and between sample  
1018 types in detecting environmental health related genera. Analysis of variance (ANOVA) was  
1019 performed to test differences in per-sample correlation coefficients between sequencing  
1020 platforms, followed by Tukey's HSD post-hoc tests when significant differences were detected.  
1021 Normality and homogeneity of variance assumptions were assessed using Shapiro-Wilk and  
1022 Levene's tests, respectively.