

# UCSF

## UC San Francisco Previously Published Works

### Title

Fast and accurate metagenotyping of the human gut microbiome with GT-Pro

### Permalink

<https://escholarship.org/uc/item/76p7w1wx>

### Journal

Nature Biotechnology, 40(4)

### ISSN

1087-0156

### Authors

Shi, Zhou Jason  
Dimitrov, Boris  
Zhao, Chunyu  
[et al.](#)

### Publication Date

2022-04-01

### DOI

10.1038/s41587-021-01102-3

Peer reviewed

## Ultra-rapid metagenotyping of the human gut microbiome

Zhou Jason Shi<sup>1,2</sup>, Boris Dimitrov<sup>3</sup>, Chunyu Zhao<sup>1</sup>, Stephen Nayfach<sup>4,5,\*</sup> and Katherine S. Pollard<sup>1,2,6,\*</sup>

<sup>1</sup>Chan Zuckerberg Biohub, Data Science, San Francisco, CA, <sup>2</sup>Gladstone Institutes, San Francisco, CA,

<sup>3</sup>Chan Zuckerberg Initiative, CA, <sup>4</sup>Department of Energy, Joint Genome Institute, Walnut Creek, CA,

<sup>5</sup>Lawrence Berkeley National Laboratory, Environmental Genomics and Systems Biology Division, Berkeley, CA, <sup>6</sup>University of California San Francisco

\*e-mail: [snayfach@lbl.gov](mailto:snayfach@lbl.gov); [katherine.pollard@gladstone.ucsf.edu](mailto:katherine.pollard@gladstone.ucsf.edu)

### Abstract

Sequence variation is used to quantify population structure and identify genetic determinants of phenotypes that vary within species. In the human microbiome and other environments, single nucleotide polymorphisms (SNPs) are frequently detected by aligning metagenomic sequencing reads to catalogs of genes or genomes. But this requires high-performance computing and enough read coverage to distinguish SNPs from sequencing errors. We solved these problems by developing the GenoTyper for Prokaryotes (GT-Pro), a suite of novel methods to catalog SNPs from genomes and use exact k-mer matches to perform ultra-fast reference-based SNP calling from metagenomes. Compared to read alignment, GT-Pro is more accurate and two orders of magnitude faster. We discovered 104 million SNPs in 909 human gut species, characterized their global population structure, and tracked pathogenic strains. GT-Pro democratizes strain-level microbiome analysis by making it possible to genotype hundreds of metagenomes on a personal computer.

**Software availability:** GT-Pro is available at <https://github.com/zjshi/gt-pro>.

### Introduction

Microbial species harbor extensive genetic variation, including single nucleotide polymorphisms (SNPs), structural variants (SVs), and mobile genetic elements. SNPs in particular are useful for population genetic analyses<sup>1</sup>, such as tracking transmission of strains between environments or locations, reconstructing strain phylogenetic relationships, resolving mixtures of genotypes within a host, and depicting population diversity or structure along environmental gradients. Additionally, SNPs can result in changes in protein function. For example, a single SNP in the *Dadh* gene of the human commensal *Eggerthella lenta* can predict activity of levodopa, the primary medication used to treat Parkinson's disease<sup>2</sup>. Quantifying intra-species genomic variation in the human microbiome is a prerequisite to the potential application of microbiome genomics to precision medicine.

Several approaches exist for identifying SNPs in microbiomes. The gold standard<sup>3</sup> is to sequence individual isolate genomes and identify mismatches in whole-genome alignments. In contrast, metagenomes are a rich source of strain level diversity for uncultivated taxa. In a landmark study, Schloissnig et al<sup>4</sup> discovered 10.3 million SNPs for 101 human gut species by aligning short reads from shotgun metagenomes to reference genomes. This approach is known as “metagenotyping” and has since been featured in several tools, including Constrains<sup>5</sup>, MIDAS<sup>6</sup>, metaSNV<sup>7</sup>, DESMAN<sup>8</sup> and StrainPhlAn<sup>9</sup>. While algorithms for read alignment have improved, the approach is still computationally costly. Exact matching algorithms such as Kraken<sup>10</sup>, CLARK<sup>11</sup> and bfMEM<sup>12</sup>, have been developed as a more efficient solution to the read mapping problem, achieving speedups by orders of magnitude. However, these tools have thus far been used to quantify the abundance of microbiome taxa, rather than identify intra-species genetic variation. Genotyping by exact matches between reads and short sequences covering SNPs was

51 implemented in the method LAVA<sup>13</sup> for human whole-genome sequencing data. Our goal was to extend  
52 this approach to metagenomes by addressing the challenges presented by complex mixtures of species and  
53 strains within microbiome samples, while also making software that could run on a personal computer.

54

## 55 **Results**

56

### 57 *A novel framework for in silico genotyping of microbiome species*

58

59 We introduce the GenoTyper for PROkaryotes (GT-Pro), which is a novel computational pipeline that  
60 utilizes an exact matching algorithm to perform ultra-rapid and accurate genotyping of known SNPs from  
61 metagenomes. Our proof-of-principle initial implementation of this approach focuses on the human gut  
62 microbiome. We created a reference database of 104 million common SNPs that we identified using  
63 112,904 high-quality genomes from 909 human gut microbiome species. Then we used this catalog to  
64 perform reference-based SNP calling for 25,133 publicly available metagenomes, providing insight into  
65 strain variation across individuals and geographic regions. Our results demonstrate the feasibility of  
66 performing large-scale metagenotyping without need for high-performance computing.

67

68 To overcome the low throughput, sensitivity and species coverage of current alignment-based  
69 metagenotyping methods, we developed the GT-Pro framework (Fig. 1). Our key innovations are (i)  
70 capturing the majority of common variation found in microbiome genomes with a compact database of  
71 SNP covering k-mers (sck-mers), (ii) selecting highly species-specific sck-mers, overcoming high false  
72 positives associated with k-mer exact-matching methods, and (iii) developing and optimizing algorithms  
73 and data structures for exact matching of metagenomic sequencing reads to these sck-mers, enabling  
74 SNPs to be detected rapidly and accurately in microbiome samples. Building a version of GT-Pro for a  
75 given environment involves 1) discovering common SNPs in assembled genomes for each species, 2)  
76 optionally identifying linkage disequilibrium (LD) blocks and “tag” SNPs that capture most variation  
77 within each block, and 3) designing species-specific sck-mers. We focus on common SNPs because this  
78 allows us to create a virtual genotyping “array” that is a data structure small enough to fit in computer  
79 memory while still capturing the majority of prevalent genetic variation for each of many species.

80

### 81 *A database of common SNPs for bacterial species in the gut microbiome*

82

83 As a case study, we applied GT-Pro to the human gut microbiome due to the large number of microbial  
84 genomes from this environment and its important role in human health. To construct a SNP catalog, we  
85 used 112,904 high-quality genomes ( $\geq 90\%$  completeness and  $\leq 5\%$  contamination<sup>14</sup>) from 909 species  
86 (minimum = 10 genomes, median = 35 genomes) that we downloaded from the Unified Gastrointestinal  
87 Genomes (UHGG) resource<sup>15</sup> (Fig. S1, S2 and Table S1). These include both metagenome-assembled  
88 genomes<sup>16-18</sup> (i.e. MAGs, 94.1%) as well as cultivated isolates (5.9%) and were derived from  
89 geographically and phenotypically diverse human subjects. We performed whole-genome alignments for  
90 each species, revealing 104,171,172 common, core-genome SNPs (minor allele frequency  $\geq 1\%$ , site  
91 prevalence  $\geq 90\%$ ), the vast majority of which (93.4%) were bi-allelic (Fig. 2a, S3a and S4). An  
92 extremely low fraction of SNPs ( $<0.2\%$ ) either disrupted a stop codon or introduced a premature one,  
93 which is one indicator of false positives (Fig. 2a). For context, this catalog is 10-fold larger than the one  
94 established by Schloissnig et al. and 1.22-fold larger than the catalogue of all human SNPs<sup>19</sup> (Fig. S1).  
95 Consistent with previous reports<sup>4</sup>, SNP density, nucleotide diversity, and the rate of nonsynonymous  
96 versus synonymous mutations (pN/pS) varied across species and phyla (Fig. 2b and Fig. S5-8), which  
97 may reflect differences in selective pressures, population sizes, or transmission modes.

98

99 We hypothesized that the SNP database could be greatly compressed by clustering SNPs into linkage  
100 disequilibrium (LD) blocks that co-vary across reference genomes (Fig. S9) and selecting a single "tag"

101 SNP per LD block. A similar strategy is commonly used when designing genotyping arrays in human  
102 genetics. Using single-linkage clustering ( $R^2 > 0.81$ ), the 104 million SNPs were clustered into 6.8  
103 million LD blocks, representing a >15-fold reduction in database size and revealing a remarkable degree  
104 of local genomic structure. Our choice of  $R^2$  is motivated by thresholds used for high confidence SNP  
105 imputation in other species and the fact that discovery of LD blocks stabilizes in this range for gut species  
106 (Fig. S10). On average LD blocks spanned ~4.3Kbp and ~23.5 SNPs, though the number and size of LD  
107 blocks varied considerably across bacterial species (Fig. 2c, S5c and S11a and b). As expected, linkage  
108 between SNPs decayed with increasing genomic distance (Fig. 2d-f), though decay rates differed  
109 substantially across species (Fig. 2d-e). Altogether, these differences in genetic diversity and structure  
110 across species likely reflect variation in recombination rates and/or the number and relatedness of  
111 sequenced genomes.

112

### 113 *Species-specific kmers enable accurate and efficient identification of SNPs*

114

115 Having constructed a large SNP catalog of the gut microbiome, we next used GT-Pro to identify k-mers  
116 that could unique identify each SNP from shotgun metagenomes. We empirically determined that length  
117  $k=31$  ensured high specificity while limiting compute and memory requirements. Of the ~13.3 billion  
118 candidate 31-mers that overlapped a SNP (124 per SNP), we identified 5.7 billion that were unique. These  
119 kmers overlapped 51% of the 104 million SNPs for 65% of LD blocks (mean 108 sck-mers per SNP, >1  
120 sck-mer for 97% of species, Fig. S1 and S12). We refer to these as species-specific, SNP-covering kmers  
121 (sck-mers). Species with few or no SNPs that can be genotyped with this strategy include those with a  
122 very close relative and are most common within Actinobacteria (Fig. 2g and S3b). While only 50% of  
123 SNPs were tagged by a sck-mer, they capture 83% of the within species variation compared to whole-  
124 genome average nucleotide identity, and achieve a much higher level of resolution compared to individual  
125 taxonomic markers (Fig. S13). Due to the large scale of the database, GT-Pro uses a highly efficient data  
126 structure to store the sck-mers, requiring only 13 GB of RAM and permitting GT-Pro to run on most  
127 modern personal computers (Fig. S14 and S15). We also created a low memory version of the GT-Pro  
128 database (< 4 GB RAM) which just stores sck-mers for a single "tag" SNP per LD block (Methods) and  
129 still captures the majority of within species variation (Fig. S13).

130

### 131 *Optimized k-mer exact matching accelerates metagenotyping 100-fold*

132

133 To search for exact matches between billions of k-mers among metagenome reads and billions of sck-  
134 mers in the GT-Pro database, it is crucial to have a highly efficient search algorithm with low RAM and  
135 I/O requirements. To this end, we developed an exact match algorithm that leverages data structures  
136 optimized for this specific application (Fig. S16). Our approach is similar to a multi-index search with  
137 three main steps operating on bit encoded k-mers (2 bits per base) (Fig. S16a). After generating all k-mers  
138 in each metagenomic sequencing read, GT-Pro uses a 1-bit Bloom filter on the first  $l < k$  bits of each k-mer  
139 to quickly rule out the vast majority of read k-mers that have no chance to match database sck-mers  
140 because they do not share an l-mer. For the k-mers that pass through the 1-bit filter, the algorithm recruits  
141 an m-bit (last m bits of encoded k-mer) index to serve as secondary filter that locates a bucket of pre-  
142 sorted sck-mers in the database containing all possible exact matches to the full k-mer. Finally, the  
143 algorithm invokes a sequential search for exact matches between the full k-mer and these only the sck-  
144 mers in this bucket.

145

146 We next evaluated GT-Pro computational performance. First, we measured both speed and peak RAM  
147 use while tuning the values of  $l$  and  $m$ , two parameters derived from the 1-bit and m-bit filter that are  
148 expected to have a large impact on performance due to their direct relationships with query speed and  
149 peak RAM use. In general, both performance metrics increase with higher values of  $l$  and  $m$  (Fig. 3a).  
150 Within the range of the tested parameters, we found best speed and peak RAM use with  $l = 30$  and  $m = 35$

151 in the laptop environment (26.5GB RAM) and the with  $l = 32$  and  $m = 36$  on a server (56.55 GB RAM).  
152 In a boundary case ( $l = 30$  and  $m = 36$ ) on the laptop where the peak RAM use hit the hardware limit,  
153 speed drops  $>87\%$ . These results demonstrate that the values of  $l$  and  $m$  should be carefully chosen based  
154 on the hardware for optimal performance, which is handled automatically by GT-Pro.

155  
156 We then compared the computational performance of GT-Pro to traditional read alignment method as  
157 baseline (Fig. 3b). We arbitrarily selected a total of 40 stool metagenomes from a Tanzanian cohort<sup>20</sup>  
158 (Table S7) for the evaluation. For alignment, SNPs were called mapping reads to database from GT-Pro  
159 and an independent one (metaSNV<sup>7</sup>). Although GT-pro had a larger peak RAM use than alignment  
160 method ( $<10\text{GB}$ ), was 100x faster on a server and 10x faster on a laptop where peak RAM use was  
161 26.5GB (Fig. 3b).

162  
163  
164 *Accurate identification of SNPs from simulated metagenomes*

165  
166 We next evaluated the accuracy of SNP calling with GT-Pro compared to alignment using simulated  
167 metagenomes. Towards this goal, we generated Illumina sequencing reads *in silico* from 978 human gut  
168 isolates<sup>21</sup> and identified the ground truth set of SNPs based on whole-genome alignment. We first  
169 simulated reads from individual isolates with sequencing coverages ranging from 0.001x to 15x (Table  
170 S2). In this simplified scenario, genotyping errors can result from sequencing errors, insufficient  
171 coverage, or incorrect read or k-mer mapping. Across isolates and coverage levels, the false discovery  
172 rate (FDR) of genotypes was on average lower for GT-Pro (median= 0.7%, IQR=1.1%) compared to read  
173 alignment (median= 2.2%, IQR=4.7%) (Fig. 4a) while the median sensitivity of GT-Pro tended to be  
174 consistently higher (4.1-17.6%) at all coverages (Fig. 4b). While read alignment methods typically use a  
175 minimum coverage threshold to avoid false positives from sequencing error (e.g.  $>10\text{x}$ ), that would have  
176 further decreased the sensitivity in this experiment.

177  
178 Next, we simulated metagenomes containing pairs of conspecific isolates to evaluate performance on  
179 samples with strain mixtures, exploring a range of coverage ratios from 0.001x to 15x, where one strain is  
180 always at 15x coverage and the other varying (Table S3 and S4). In terms of detecting heterozygous sites  
181 (strains with different alleles), the false discovery rate (FDR) of GT-Pro (median= 0.9%, IQR=0.7%) was  
182 slightly higher compared to alignment (median= 0.3%, IQR=0.6%) (Fig. 4c), however, median sensitivity  
183 was higher (50.5-81.6%) for GT-Pro at all coverage ratios (Fig. 4d). A higher FDR for GT-Pro is likely  
184 caused by sequencing errors that match the alternative allele by chance, which also could cause a slightly  
185 lower sensitivity for GT-Pro at homozygous sites (Fig. S17).

186  
187 To evaluate genotype calls imputed from tag SNPs, we found low FDR $<5\%$  comparing to true genotypes  
188 for the vast majority of isolates ( $>95\%$ ) (Fig. 4e). SNPs belonging to an LD block were 5 times more  
189 likely to be detected (non-zero read count) when their tag SNPs were also detected than when they were  
190 not (Fig. S18). To show that GT-Pro is highly quantitative, we compared average coverage at SNPs in the  
191 GT-Pro output to the known genome coverage using metagenomes we simulated from individual isolates  
192 and pairs of conspecific isolates. Even at low sequencing coverage ( $<1\text{X}$ ), GT-Pro was able to accurately  
193 estimate the true coverage of each species (Fig. 4f) and the ratio between two strains (Fig. 4g). These  
194 results suggest that GT-Pro allele calls and counts could be used to impute genotypes and estimate  
195 relative abundances of species and strains accurately.

196  
197 *Accurate metagenotyping and gene imputation from gut metagenomes*

198  
199 To compare GT-Pro to existing approaches, we metagenotyped gut metagenomes<sup>16,20,22,23</sup> (Table S5-10)  
200 with alignment and compared the number of genotyped SNPs plus estimates of allele frequencies and



201 genetic distances. We found that GT-Pro genotyped more species and SNPs per metagenome (Fig. S19a-  
202 c), despite being limited to species with  $\geq 10$  genomes. This is likely due to GT-Pro having better  
203 sensitivity for low coverage species and using a human gut focused database (comparing to metaSNV).  
204 For species genotyped by both methods, within-sample heterozygosity (Fig. S18) and across-sample  
205 allele presence and frequency (Fig. 5a-d) were highly correlated. For high coverage species, alignment  
206 method detected some SNPs absent from the GT-Pro database, whereas GT-Pro detected more sites as  
207 polymorphic in medium and low coverage species (Fig. 5a-d). Despite these differences in genotyped  
208 sites, GT-Pro and alignment produce highly similar estimates of pairwise genetic distances (Jaccard  
209 index) between samples, likely because rare variants missed by GT-Pro but with sufficient coverage to be  
210 genotyped with alignment-based methods represent a small fraction of overall genetic diversity. For  
211 comparison, we repeated this analysis using only SNPs in the 16S gene and observed much lower genetic  
212 differences between samples (Fig. S21), emphasizing that GT-Pro provides strain resolution close to that  
213 of alignment and greatly exceeding that of marker gene approaches. Altogether these results are  
214 consistent with our simulations and underscore the high sensitivity of GT-Pro.

215  
216 Next, we sought to determine if GT-Pro SNPs could be used to infer the presence of nearby genes or  
217 operons, thereby serving as biomarkers for structural variants. As a case study, we used GT-Pro SNPs in  
218 flanking genes to predict presence/absence of toxicity controlling genes in *Clostridium difficile* (*C.*  
219 *difficile*). We used the GT-Pro SNPs from two 5' (CD2601 and CD2602), one 3' (trpS) gene and  
220 intergenic region to train a Random Forest classifier to predict the presence of the genomic region  
221 (CdtLoc) of three toxin genes (CD196\_cdtA, CD630\_cdtAB and cdtR) in a set of *C. difficile* isolate  
222 genomes downloaded independently from NCBI (Fig. S22a and S23a). In another example, we  
223 demonstrated that SNPs (cdd2, cdu1, and intergenic) flanking a pathogenicity locus (PaLoc) region could  
224 predict its presence (Fig. S22b and S23b). Next we applied these models to GT-Pro metagenomes from  
225 7,459 samples (Fig. 5e and f). Our predictions of CdtLoc and PaLoc region presence were highly  
226 correlated with estimated presence based on read alignment to the *C. difficile* genome, especially in  
227 metagenomes where this species was more abundant, and weaker predictions were made when not all of  
228 the genes in CdtLoc or PaLoc were present (Fig. S24a and b). These results show that GT-Pro can detect  
229 structural and strain variants when they are in high LD with flanking common SNPs.

### 230 231 *Depicting novel and global intra-species genetic structure with GT-Pro*

232  
233 To evaluate the commonality of SNPs in GT-Pro database and how GT-Pro perform in metagenotyping  
234 unknown metagenomes. We next used GT-Pro's common SNPs to perform dimension reduction on the  
235 genomes in the database as well as metagenomic samples from a North American IBD cohort<sup>24</sup> (n=220;  
236 Table S11) that did not contribute genomes to the GT-Pro database. Looking for evidence of subspecies  
237 genetic structure, we observed that for most species the metagenomes clustered with the genomes (Fig. 6a  
238 and b), suggesting that GT-Pro's database represents the common diversity across diverse metagenomes.  
239 For a few species, however, we observed clusters comprised only of metagenomes (Fig. 6c and d),  
240 demonstrating that novel subspecies genetic structure can be discovered using GT-Pro common SNPs.

241  
242 Having shown GT-Pro is faster and at least as accurate as alignment-based methods for genotyping  
243 common SNPs from metagenomes, we leveraged GT-Pro metagenotypes to conduct the most  
244 geographically diverse intra-species genetic variation meta-analysis to date, encompassing 51.8 million  
245 SNPs for 881 species found in 7,459 gut samples from 31 locations across six continents (Table S13).  
246 Consistent with prior studies<sup>4,6,9</sup>, we observed much less allele sharing between hosts (median=0.03,  
247 IQR=0.05) than within a host over time (median=0.38, IQR=0.4), and that intra-host allele sharing varies  
248 greatly between species and hosts (Fig. S25). Inter-host allele sharing differed across countries and  
249 continents (Fig. S26a and b), generally decreasing with geographic distance (Fig. S26a and b and S27a  
250 and b) and varying across species (Fig. S28). Our results also show clear associations with degree of

251 industrialization as well as relatedness of hosts (e.g., hosts within villages in Fiji share more alleles than  
252 unrelated hosts in North American cities) (Fig. 6e). To identify gut species with high levels of inter-  
253 continental population differentiation, we calculated  $F_{ST}$  for 78 prevalent and well-detected (see  
254 methods) species and observed large differences in the degree of differentiation across species (Fig. 6f).  
255 Species with high  $F_{ST}$  show distinct clusters of hosts, some but not all of which correlate with geography  
256 (Fig 6g), consistent with lifestyle and environment playing a role in which strains colonize a host. In  
257 contrast, hosts do not cluster as clearly based on species relative abundance (Fig 6h), emphasizing that  
258 metagenotypes may reveal microbiome-host associations missed in abundance analyses.

259  
260

261

## 262 **Discussion**

263

264 Here, we greatly extended the gut microbiome genomic variation landscape by identifying more than 100  
265 million common core-genome SNPs from 909 bacterial species. As our solution to the bioinformatics  
266 challenge of metagenotyping, GT-Pro avoids computationally costly alignment and overcomes computing  
267 barriers. It performs strain-level analyses of microbiomes with improved accuracy, especially for low  
268 coverage species. Studies of microbiome genetic variation on a laptop or at the scale of human genome-  
269 wide association studies will be computationally feasible with GT-Pro.

270

271 It should be noted that our method comes with several limitations. First, the GT-Pro database does not  
272 capture all human gut microbial diversity. While we used 909 species, we could not use the majority of  
273 the UHGG species due to limited availability of high-quality genomes. Second, GT-Pro is analogous to a  
274 genotyping array and hence does not identify novel SNPs, which require other methods, such as  
275 alignment-based SNP calling or single-cell genome sequencing. For some species, the common SNP pool  
276 is expected to expand through additional genome sequencing. Third, a small number of species lacked  
277 species-specific *sck*-mers due to the presence of highly related species in the genome collection. Separate  
278 strategies such as using longer *k*-mers or less common SNPs could enable GT-Pro metagenotyping for  
279 these species. Fourth, although we were very selective in the choice of genomes and SNPs used for  
280 building GT-Pro, it is impossible to exclude all imperfections (e.g. incompleteness, contaminations and  
281 species misclassification) in the genome assemblies that could contribute to false SNP calls. Finally, GT-  
282 Pro does not directly genotype structural variants, which contribute significantly to intra-species genetic  
283 diversity<sup>25</sup>. However, we did show that GT-Pro can be used to impute insertions and deletions in high LD  
284 with common SNPs. Despite these caveats, we showed that the GT-Pro framework is general, accurate  
285 and sensitive for identifying genetic variation in metagenomes.

286

287 We envision several directions for future work. First, this study applied the GT-Pro approach to human  
288 gut prokaryotic species, and the framework could easily be expanded to other kingdoms and  
289 environments. Another extension is to develop alignment-free metagenotyping for short indels and  
290 structural variants. This study barely scratches the surface in terms of interpreting microbiome genetic  
291 variation. Towards leveraging microbiomes in precision medicine, it will be critical to comprehensively  
292 identify SNPs that are associated with disease and other traits (e.g. pathogenicity, antimicrobial  
293 resistance, drug degradation). We anticipate that GT-Pro will also be useful for detecting contamination,  
294 recombination, and horizontal gene transfer events, as well as tracking variants or strains over time, host  
295 lifestyle and geography.

296

297

298 **Reference**

- 299
- 300 1. Garud, N. R. & Pollard, K. S. Population Genetics in the Human Microbiome. *Trends Genet.* **36**,  
301 53–67 (2020).
- 302 2. Maini Rekdal, V., Bess, E. N., Bisanz, J. E., Turnbaugh, P. J. & Balskus, E. P. Discovery and  
303 inhibition of an interspecies gut bacterial pathway for Levodopa metabolism. *Science* **364**, eaau6323  
304 (2019).
- 305 3. Treangen, T. J., Ondov, B. D., Koren, S. & Phillippy, A. M. The Harvest suite for rapid core-  
306 genome alignment and visualization of thousands of intraspecific microbial genomes. *Genome Biol.* **15**,  
307 524 (2014).
- 308 4. Schloissnig, S. *et al.* Genomic variation landscape of the human gut microbiome. *Nature* **493**, 45–  
309 50 (2013).
- 310 5. Luo, C. *et al.* ConStrains identifies microbial strains in metagenomic datasets. *Nat. Biotechnol.*  
311 **33**, 1045–1052 (2015).
- 312 6. Nayfach, S., Rodriguez-Mueller, B., Garud, N. & Pollard, K. S. An integrated metagenomics  
313 pipeline for strain profiling reveals novel patterns of bacterial transmission and biogeography. *Genome*  
314 *Res.* **26**, 1612–1625 (2016).
- 315 7. Costea, P. I. *et al.* metaSNV: A tool for metagenomic strain level analysis. *PLOS ONE* **12**,  
316 e0182392 (2017).
- 317 8. Quince, C. *et al.* DESMAN: a new tool for de novo extraction of strains from metagenomes.  
318 *Genome Biol.* **18**, 181 (2017).
- 319 9. Truong, D. T., Tett, A., Pasolli, E., Huttenhower, C. & Segata, N. Microbial strain-level  
320 population structure and genetic diversity from metagenomes. *Genome Res.* **27**, 626–638 (2017).
- 321 10. Wood, D. E. & Salzberg, S. L. Kraken: ultrafast metagenomic sequence classification using exact  
322 alignments. *Genome Biol.* **15**, R46 (2014).
- 323 11. Ounit, R., Wanamaker, S., Close, T. J. & Lonardi, S. CLARK: fast and accurate classification of  
324 metagenomic and genomic sequences using discriminative k-mers. *BMC Genomics* **16**, 236 (2015).
- 325 12. Liu, Y., Zhang, L. Y. & Li, J. Fast detection of maximal exact matches via fixed sampling of  
326 query K-mers and Bloom filtering of index K-mers. *Bioinformatics* **35**, 4560–4567 (2019).
- 327 13. Shajii, A., Yorukoglu, D., William Yu, Y. & Berger, B. Fast genotyping of known SNPs through  
328 approximate k-mer matching. *Bioinforma. Oxf. Engl.* **32**, i538–i544 (2016).
- 329 14. Parks, D. H., Imelfort, M., Skennerton, C. T., Hugenholtz, P. & Tyson, G. W. CheckM: assessing  
330 the quality of microbial genomes recovered from isolates, single cells, and metagenomes. *Genome Res.*  
331 **25**, 1043–1055 (2015).
- 332 15. Almeida, A. *et al.* A unified sequence catalogue of over 280,000 genomes obtained from the  
333 human gut microbiome. *bioRxiv* 762682 (2019) doi:10.1101/762682.
- 334 16. Pasolli, E. *et al.* Extensive Unexplored Human Microbiome Diversity Revealed by Over 150,000  
335 Genomes from Metagenomes Spanning Age, Geography, and Lifestyle. *Cell* **176**, 649–662.e20 (2019).
- 336 17. Nayfach, S., Shi, Z. J., Seshadri, R., Pollard, K. S. & Kyrpides, N. C. New insights from  
337 uncultivated genomes of the global human gut microbiome. *Nature* **568**, 505–510 (2019).
- 338 18. Almeida, A. *et al.* A new genomic blueprint of the human gut microbiota. *Nature* **568**, 499–504  
339 (2019).
- 340 19. Auton, A. *et al.* A global reference for human genetic variation. *Nature* **526**, 68–74 (2015).
- 341 20. Smits, S. A. *et al.* Seasonal cycling in the gut microbiome of the Hadza hunter-gatherers of  
342 Tanzania. *Science* **357**, 802 (2017).
- 343 21. Zou, Y. *et al.* 1,520 reference genomes from cultivated human gut bacteria enable functional  
344 microbiome analyses. *Nat. Biotechnol.* **37**, 179–185 (2019).
- 345 22. Turnbaugh, P. J. *et al.* The Human Microbiome Project. *Nature* **449**, 804–810 (2007).
- 346 23. Bäckhed, F. *et al.* Dynamics and Stabilization of the Human Gut Microbiome during the First  
347 Year of Life. *Cell Host Microbe* **17**, 852 (2015).



- 348 24. Franzosa, E. A. *et al.* Gut microbiome structure and metabolic activity in inflammatory bowel  
349 disease. *Nat. Microbiol.* **4**, 293–305 (2019).
- 350 25. Tierney, B. T. *et al.* The Landscape of Genetic Content in the Gut and Oral Human Microbiome.  
351 *Cell Host Microbe* **26**, 283-295.e8 (2019).  
352

353

## 354 **Figure legends**

355

### 356 **Figure 1. In silico metagenotyping framework**

357 Our method starts with a whole genome sequence collection and identifies species with sufficient high-  
358 quality genomes to call SNPs. For each species, a representative genome is chosen based on pairwise  
359 Average Nucleotide Identity (ANI) plus assembly quality metrics. SNPs are called per species based upon  
360 whole genome alignment of conspecific genomes to the representative genome. Common (MAF > 1%)  
361 bi-allelic SNPs are selected for genotyping. Up to 4 X k candidate k-mers are extracted per SNP site,  
362 covering both the reference and alternative allele on forward and reverse complementary strands (sck-  
363 mers; k=31 in this study). These candidate sck-mers are iteratively filtered through species-specificity  
364 filters of all unique k-mers present in the genomes of every other species, including species with  
365 insufficient high-quality genomes for genotyping. Only SNPs with sck-mers for both the reference and  
366 alternative allele are retained. SNPs are clustered based on pairwise linkage disequilibrium (LD). LD  
367 blocks are detected with a threshold of mean  $r^2 > 0.81$ , and we select a tag SNP with species-specific sck-  
368 mers and the highest LD to other SNPs in the block. Optimized algorithms and compressed  
369 representations of sck-mer data enable rapid metagenotyping. Further details in Methods and Figure 3.

370

### 371 **Figure 2. Genetic landscape of 909 human gut species.**

372 (a) Summary of common SNP characteristics across all species (from left to right): at most SNPs only  
373 two alleles are observed, bi-allelic SNPs are mostly within protein-coding genes, these are largely  
374 synonymous, and the non-synonymous ones rarely disrupt or introduce a stop codon. (b and c) Phyla  
375 differ in their median SNP density and average LD block size with significant variation in density across  
376 species within each phylum. (d) Rate of LD distance decay across gut bacterial species. (b-e) are colored  
377 by bacterial phylum and share the same color scheme. (e) Examples of LD distance decay for individual  
378 species. From top to bottom are three species (species id: 102446, 101694 and 102831) with increasingly  
379 fast LD distance decay, suggesting higher recombination rates. Curves represent the fitted exponential  
380 decay model. (f) Visualization of two distinct haplotype landscapes from (upper) species *Alistipes*  
381 *putredinis* (species id: 101302) and (lower) *Bacteroides xylanisolvens* (species id: 101345). Base axis  
382 represents and is ordered by genomic coordinate. Color indicates magnitude of LD between pairs of  
383 SNPs. The examples have the same genomic span (10,000 bp). (g) Distribution across species of the  
384 percentage of SNPs that can be genotyped by GT-Pro either directly (“without LD blocks”) or by  
385 imputation using genotyped tag SNPs (“with LD blocks”). For a typical species, ~75% of SNPs can be  
386 genotyped directly and ~95% can be imputed.

387

### 388 **Figure 3. Computational performance evaluation of GT-Pro.**

389 (a) Computational performance of GT-Pro in laptop (left) and server (right) environments across values  
390 for l (Bloom filter size parameter) and m (m-index size parameter). Color gradient: processing speed,  
391 circle size: peak RAM use, black box: optimal l and m for each computing environment. (b) Comparison  
392 of speed (upper) and peak RAM usage (lower) between GT-Pro and alignment-based metagenotyping  
393 (metaSNV and MIDAS; see methods). We ran GT-Pro on both server (green) and laptop (yellow)  
394 environments, while alignment-based methods were run only in the server (grey) environment due to not  
395 being optimized for personal computers. Peak RAM usage exceeds RAM needed to store the database due  
396 to intermediate calculations, such as applying filters.

397

### 398 **Figure 4. Metagenotyping accuracy evaluation of GT-Pro using simulations.**

399 Accuracy comparisons of GT-Pro and alignment-based metagenotyping across species based on reads  
400 simulated from isolate genomes with sequencing error. (a) False discovery rate at a combination of  
401 sequencing coverage ranged from 0.001x to 15x. Each observation is the result from a metagenome  
402 containing reads from one isolate. False discoveries are genotype calls that do not match the genome from

403 which reads were simulated. (b) Sensitivity across coverage levels from the simulations in (a). Sensitivity  
404 is the probability of detecting SNPs present in the isolate genome. (c) False discovery rate at  
405 heterozygous sites in metagenomes containing reads from two isolates of each species. A combination of  
406 sequencing coverage ratio between two isolates was simulated by fixing a more abundant isolate at 15x  
407 coverage in all simulations, and varying the other isolate's coverage from 0.001x to 15x (coverage ratio =  
408 0.001:15 to 15:15). (d) Sensitivity at heterozygous sites in metagenomes from (c). Sensitivity is the  
409 probability of correctly calling the heterozygous genotype of sites that differ between the genomes from  
410 which reads were simulated. (e) False discovery rate of genotypes imputed from tag SNPs based on allele  
411 matching in simulations in (a). Imputation is simply done by selecting the genotype associated with the  
412 observed tag SNP. (f) Sequencing coverage estimated using read counts at GT-Pro genotyped SNPs  
413 correlates with the simulated coverage, even when coverage is <1x. Each observation is the estimate from  
414 metagenomic reads simulated with sequencing error from a single isolate genome. (g) Sequencing  
415 coverage ratio estimates based on read counts for each allele at GT-Pro genotyped heterozygous sites  
416 correlate with the simulated ratios of two isolate genomes, even when one is much less abundant than the  
417 other ( $\leq 1:15$ ). The more abundant isolate is at 15x coverage in all simulations.  
418

### 419 **Figure 5. Metagenotyping and gene imputation from gut metagenomes**

420 Comparison of metagenotypes from GT-Pro and alignment with gut microbiome samples from a North  
421 America cohort<sup>22</sup> (HMP project; n=358; Table S8). As an example, we show the species *Bacteroides*  
422 *stercoris* (species id: 103681). Each point represents a SNP, with color indicating if the genotypes from  
423 the two methods agree (green), both methods return a genotype but the alleles disagree (purple), or only  
424 GT-Pro returns a genotype (black). Disagreements largely occur near 0.5 allele frequency, where small  
425 differences in read counts per allele can “flip” the major and minor alleles. (a) The proportion of samples  
426 in which each SNP is genotyped (prevalence) is similar with both methods. (b) Average allele frequency  
427 across samples varies across SNPs but is highly correlated between the two methods. (c and d)  
428 Comparison similar as (a and b) showing the species GCA\_000431835.1 (genus: *Succinivibrio*, species  
429 id: 100412) from a different Madagascar cohort<sup>16</sup> (n=112; Table S9). Prediction of presence/absence of *C.*  
430 *difficile* pathogenic gene sets in human gut metagenomes from a mix of cohorts (n=7459) (Table S14) a  
431 random forest classifier built using GT-Pro SNPs from flanking regions in 117 *C. difficile* isolates  
432 (Figures S23-S24) with 10-fold cross validation. Heatmaps show the predicted (first column) and  
433 observed (based on alignment, second column) presence (black) or absence (white) in each sample  
434 (rows). Barplots show *C. difficile* relative abundance (left), whole genome sequence coverage (middle),  
435 and number of detected genes from the pathogenicity locus (right), all estimated by mapping reads from  
436 each sample to a *C. difficile* representative genome. Random Forest predictions correlate with abundance,  
437 coverage, and number of detected pathogenic genes (Figure S25). (l) CdtLoc genes. (m) PaLoc genes.  
438

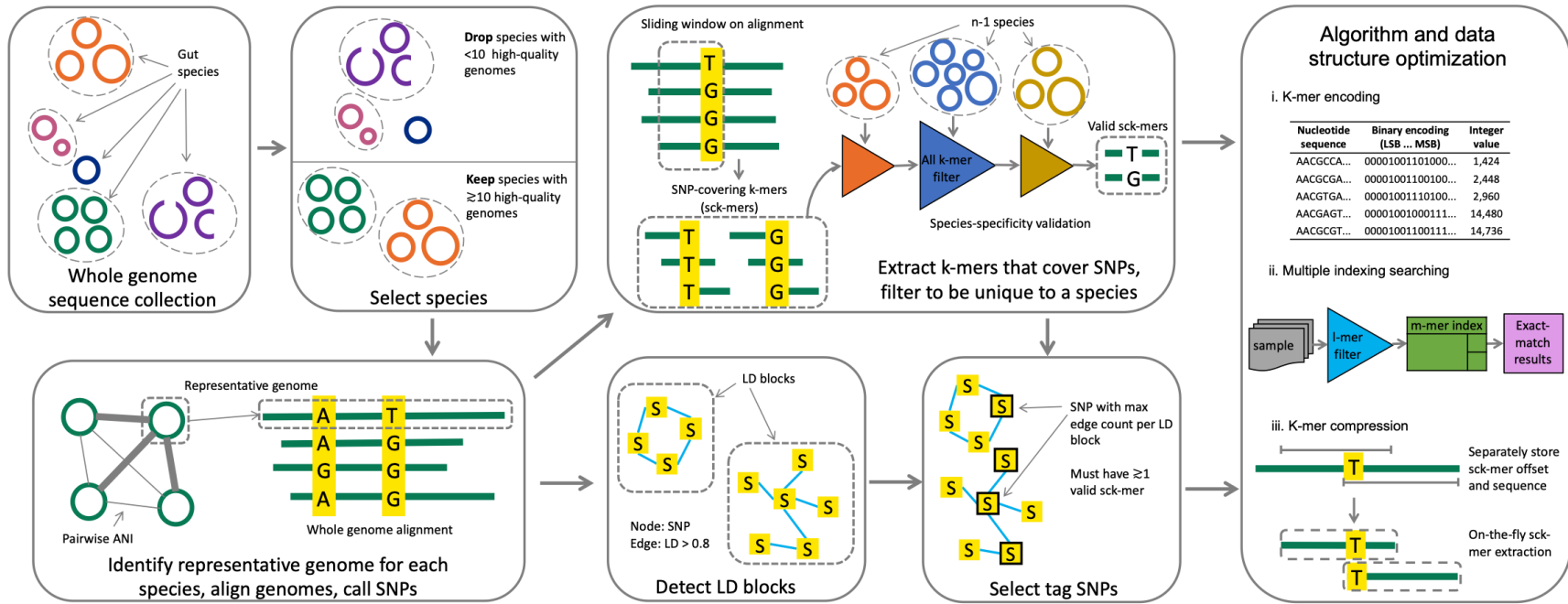
### 439 **Figure 6. Global genetic structure in 6,452 human gut metagenomes.**

440 (a-d) Gut species differ in the amount of common SNP genetic diversity already present in sequenced  
441 genomes. Metagenomic samples from a North American IBD cohort<sup>24</sup> (n=220; Table S11) (purple) are  
442 visualized in two dimensions alongside the UHGG genomes (green). Each plot is the result of applying  
443 UMAP to a matrix of genotypes at GT-Pro SNPs for one species. Each dot represents a strain of that  
444 species (major allele for heterozygous metagenomes); those closer together in UMAP space have more  
445 similar genotypes. (a) *Anaerostipes hadrus* (species id: 102528) and (b) *Ruminococcus\_B faecis* (species  
446 id: 100249) are species where metagenomes lie within the diversity previously captured by genomes. (c)  
447 *Blautia\_A obeum* (species id: 100212) and (d) *Dialister invisus* (species id: 104158) are species where  
448 metagenomes harbor combinations of common SNPs outside the range present in genomes, which may  
449 represent novel subspecies. (e) Heatmap of mean allele sharing scores over all species between  
450 metagenomes from different pairs of countries. Crossed cells indicates missing scores due to insufficient  
451 (< 5000) pairs of samples. (f) Analysis of inter-continental population differentiation (FST) for 78  
452 prevalent species. Each boxplot represents a distribution of inter-continental FST for one species, ordered

453 by medians. (g) An example of geographic patterns captured by within-species genetic variation in the  
454 GT-Pro metagenotypes of specific species. Each dot is a metagenomic sample, colored by continents.  
455 Dimension reduction and visualization performed with UMAP. The example is from *Agathobacter*  
456 *rectalis* (species id: 102492). Nearby samples in UMAP space have similar abundance profiles; the  
457 absence of distinct groups indicates that relative abundance does not show strong geographic clustering.  
458 (h) UMAP analysis based on the relative abundances of the 881 GT-Pro species in the same samples as  
459 (i).  
460  
461  
462  
463  
464  
465

466  
467

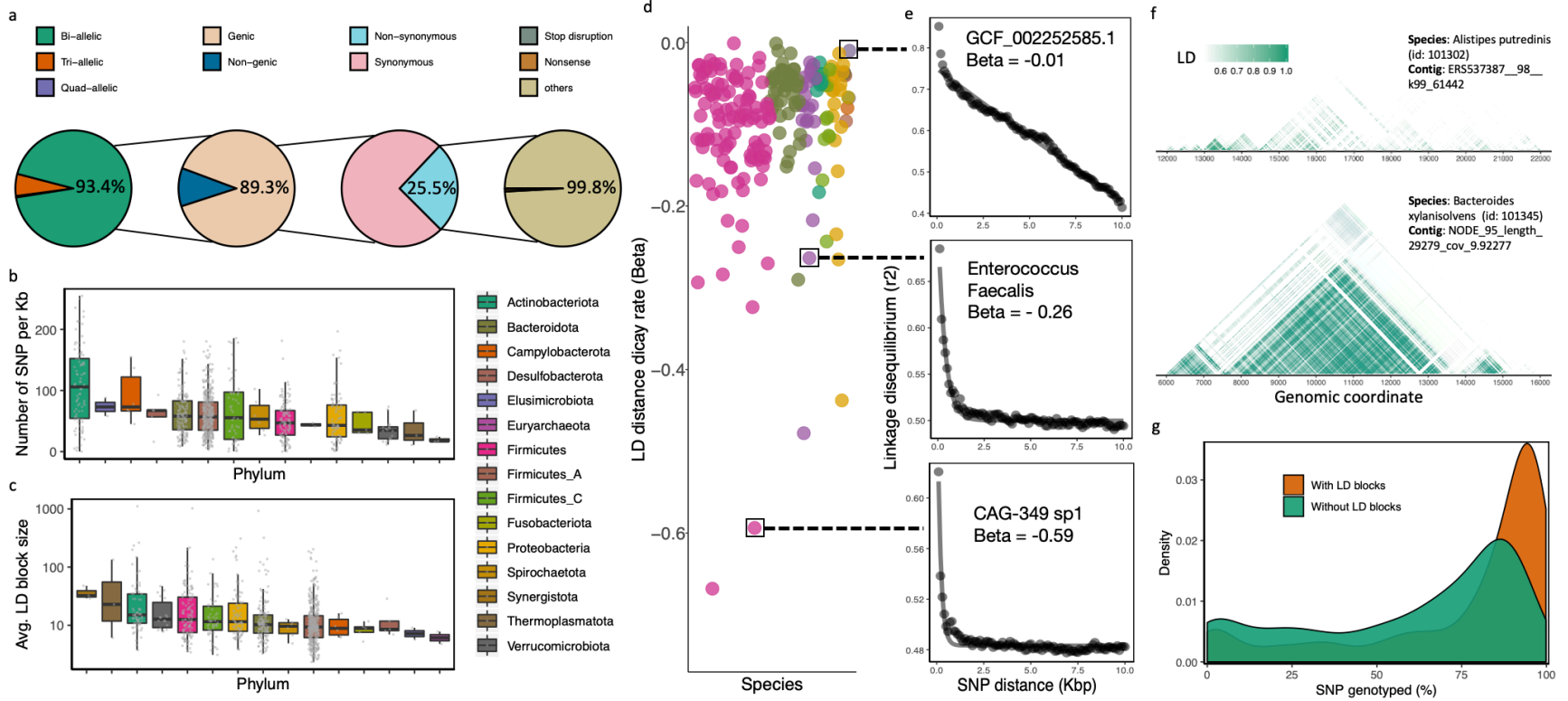
**Figure 1**



468

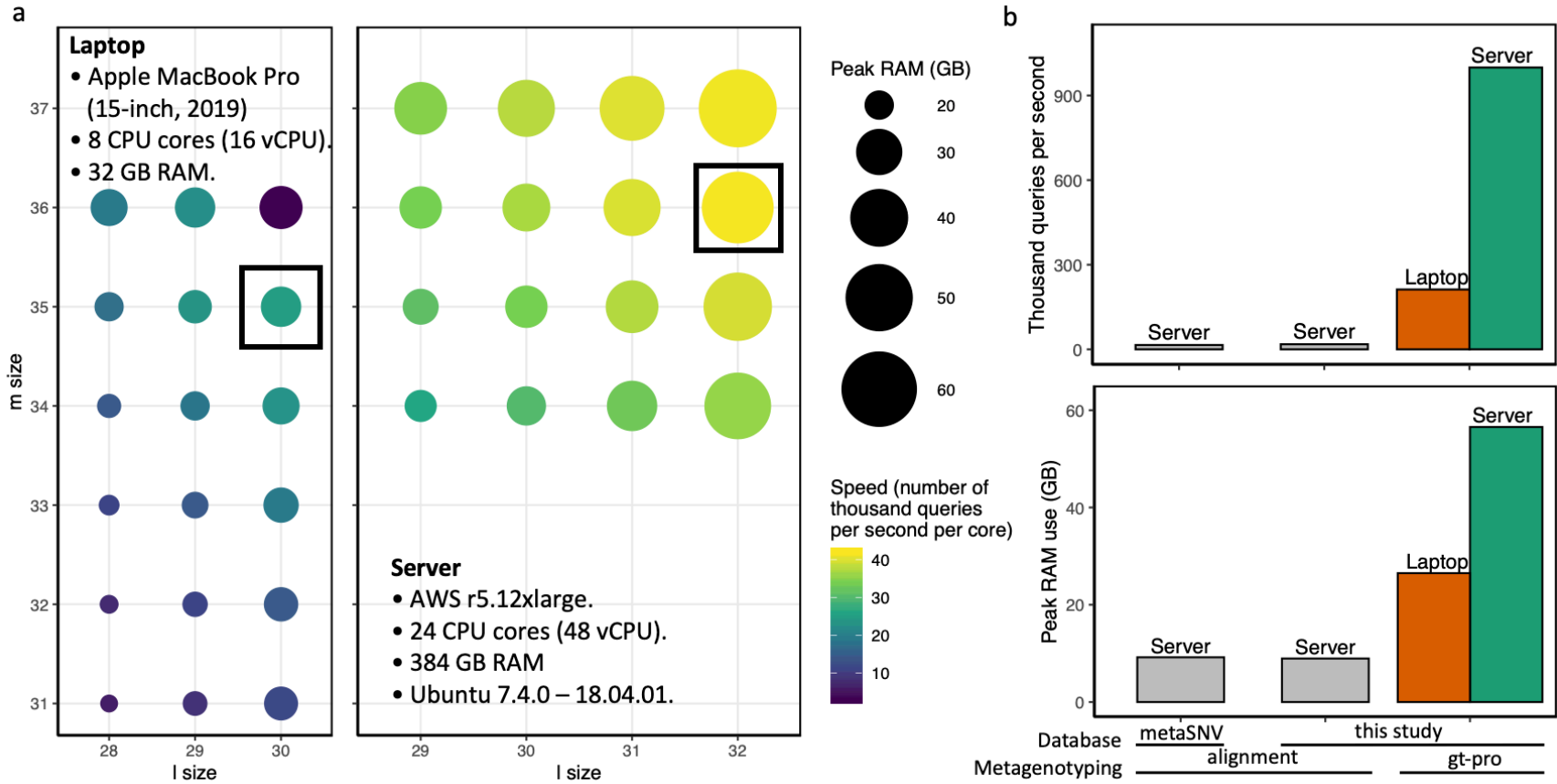


469 **Figure 2**  
 470



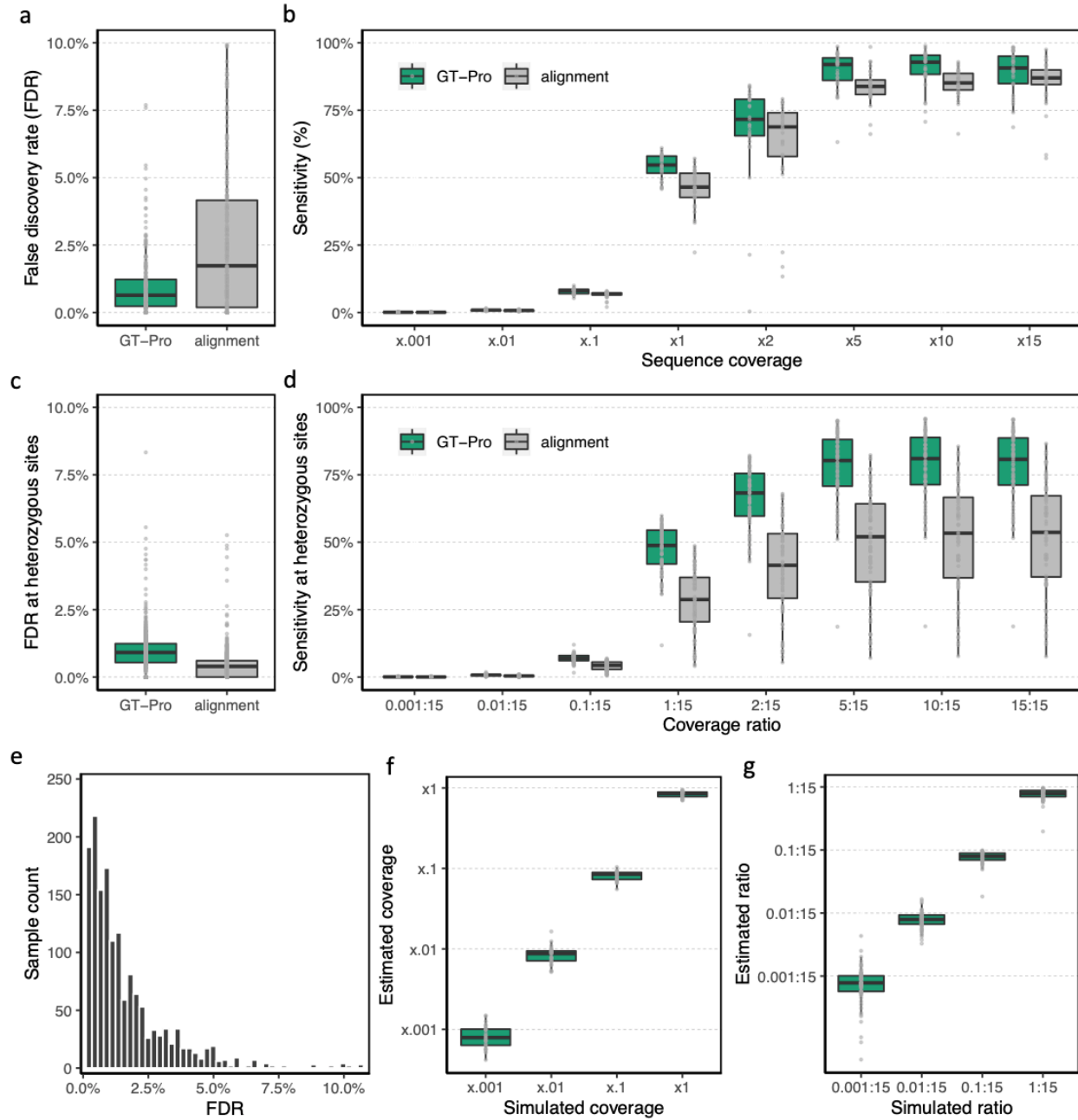
471

472 **Figure 3**  
 473



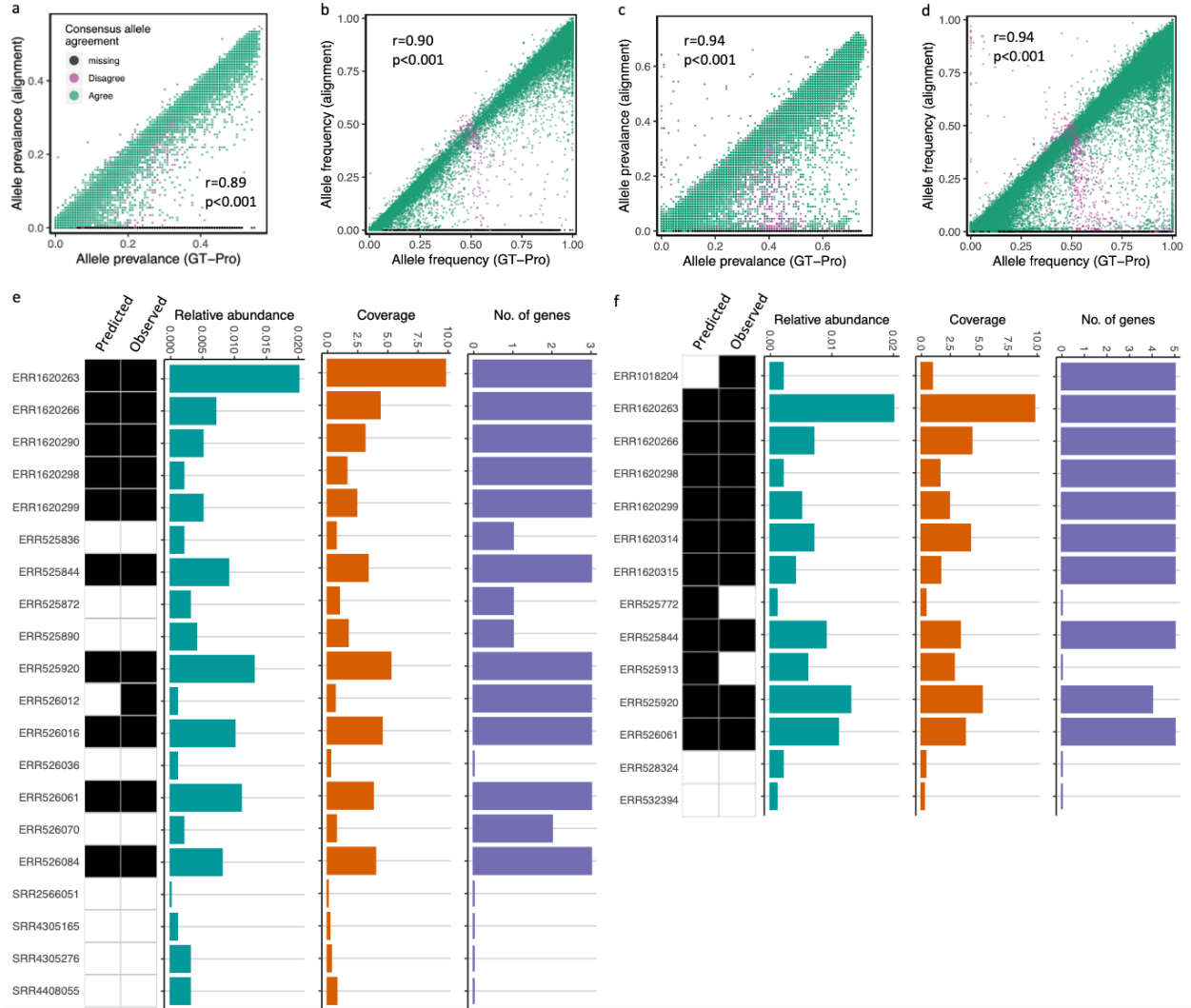
474

475 **Figure 4**  
476



477

478 **Figure 5**  
479



480  
481

Figure 6

