1   New tools for diet analysis: nanopore sequencing of metagenomic DNA from rat

2   stomach contents to quantify diet

3   Nikki E. Freed, William S. Pearman, Adam N. H. Smith, Georgia Breckell, James Dale,

4   Olin K. Silander

5   Addresses of all authors: School of Natural and Computational Sciences, Massey

6   University, Auckland 0745, New Zealand

7   4. Corresponding authors: Olin K. Silander, School of Natural and Computational

8   Sciences, Massey University, Auckland 0745, New Zealand, olinsilander@gmail.com,

9   +64 9 213 6618; Nikki E. Freed, School of Natural and Computational Sciences,

10  Massey University, Auckland 0745, New Zealand, freednikki@gmail.com, +64 9 213

11  6639

**Abstract**

12

13    Accurate determination of animal diets is difficult. Methods such as molecular barcoding

14    or metagenomics offer a promising approach, allowing quantitative and sensitive

15    detection of different taxa. Here we show that rapid and inexpensive diet quantification

16    is possible through metagenomic sequencing with the portable Oxford Nanopore

17    Technologies (ONT) MinION. Using an amplification-free approach, we profiled the

18    stomach contents from 24 wild-caught rats. We conservatively identified diet items from

19    over 50 taxonomic orders, ranging across nine phyla, including plants, vertebrates,

20    invertebrates, and fungi. This highlights the wide range of taxa that can be identified

21    using this simple approach. We calibrated the accuracy of this method by comparing the

22    characteristics of reads matching the ground-truth host genome (rat) to those matching

23    diet items, and show that at the family-level, taxon assignments are approximately

24    97.5% accurate. Some inaccuracies may arise from database biases; we suggest a way

25    to mitigate for database biases when using metagenomic approaches. Finally, we

26    implemented a constrained ordination analysis and show that we can identify the

27    sampling location of an individual rat within tens of kilometres based on diet content

28    alone. This work establishes proof-of-principle for long-read metagenomic methods in

29    quantitative diet analysis. We show that diet content can be quantified even with limited

30    expertise, using a simple, amplification free workflow and a relatively inexpensive and

31    accessible next generation sequencing method. Continued increases in the accuracy

32    and throughput of ONT sequencing, along with improved genomic databases, suggests

33    that a metagenomic approach for quantification of animal diets will become an important

34    method in the future.

## Background

Accurate quantification of animal diets can yield critical insights into ecosystem and

food web dynamics. However, unbiased and sensitive assessment of diet content is

difficult to achieve. This is largely due to the limited accuracy of many current methods.

Such methods include visual inspection of gut contents (1,2), which presents bias

against items are most easily degraded (for example, soft-bodied species);  stable

isotope analysis (3,4), which yields only broad information on diet, such as whether diet

items  are terrestrial or marine in origin (5,6); and time-lapse video (7,8), for which

species identification is difficult for small prey items or in low-light conditions.

To circumvent these issues, DNA-based methods (9,10) have become popular. Perhaps

the most widely applied DNA-based method is metabarcoding. This approach relies on

PCR amplification and sequencing of conserved regions from nuclear, mitochondrial,

or plastid genomes (9). With adequate primer selection, this method can detect a wide

range of species, and does not require specific expertise, which is often necessary for

other methods.

However, DNA metabarcoding is not free from bias: PCR primers must be specifically

tailored to particular sets of taxa or species (11). Although "universal" PCR primer pairs

have been developed (for example targeting all bilaterians or even all eukaryotes (12),

all primer sets exhibit bias towards certain taxa. Five-fold differences in fungal

operational taxonomic units (OTU) estimates have been found when using different

sets of fungal-specific PCR primer pairs  (13). It has also been shown that published

57  universal primer pairs are capable of amplifying only between 57% and 91% of tested

58  metazoan species, with as few as 33% of species in some phyla being amplified at all

59  (e.g. cnidarians)(14). Different genomic loci from the same species can exhibit up to

60  2,000-fold differences in DNA concentration, as inferred using qPCR (15). The choice

61  of polymerase can also bias diversity metrics when using metabarcoding (16). For

62  these reasons, an approach that circumvents PCR and thus avoids these biases is

63  desirable.

64  Metagenomic sequencing aims to directly sequence all of the DNA in a sample without

65  introducing bias. Although there are still biases with this approach, for example due to

66  nucleotide content affecting the likelihood of a molecule being sequenced, these are

67  inherently less than those introduced by metabarcoding. Metagenomic approaches

68  have most frequently been used to yield insights into microbial diversity and function

69  (17–24), while metagenomic applications aimed at eukaryotic taxa identification are

70  less common.  Several metagenomic diet studies have implemented filtering steps to

71  select only mitogenomic sequence or metabarcode regions, or have used abridged

72  databases before data analysis in order to mitigate database biases (25–27). However,

73  to our knowledge, very few studies have used unfiltered metagenomic sequence

74  analysis to infer diet (30,31).

75  Here, we establish a proof-of-principle methodology to accurately classify

76  metagenomic sequences from eukaryotic taxa and determine diet content using low-

77  accuracy, long-read sequencing, Oxford Nanopore (ONT). Toward this aim, we

78  quantified rat diets from several locations in the North Island of New Zealand using

79    stomach samples. Using these samples and methodology provides three distinct

80    advantages.

81    First, rats are extremely omnivorous. As such, they serve as an excellent means to

82    quantify the breadth of taxa that can be detected using a metagenomic long read

83    approach.

84    Second, the use of stomach samples means that a significant number of reads will be

85    host reads. This allows us to assess the characteristics of true positive sequence reads

86    (rat-derived reads that match rat database sequences), as well as false positive reads

87    (rat-derived reads that match non-rat database sequences). We can then determine

88    whether reads matching diet items have similar characteristics to known true positive

89    reads. This use of host reads is exactly analogous to feeding the rats a diet of known

90    content (i.e. rat) and testing whether the contents of the known diet can be accurately

91    identified.

92    Third, quantifying rat diets has important ecological implications. It is well-established

93    that the relatively recent introduction of mammalian predators to New Zealand has had

94    significant negative effects on many of the native animal populations. This ranges from

95    insects (34), to reptiles (35), to molluscs (36), to birds (37,38), with downstream effects

96    on terrestrial and aquatic ecosystems (39). To counteract the effects of mammalian

97    predators, an ambitious plan is currently being put into place that aims for the

98    eradication of all mammalian predators from New Zealand (including possums, rats,

99    stoats, and hedgehogs), by 2050 (http://www.doc.govt.nz/predator-free-2050; (40). A

100   useful step toward this goal would be to prioritise the management of predators and

101    establish in which locations native species experience the highest levels of predation.

102    To do so requires establishing the diet content of local mammalian predators.

103    Importantly, using DNA sequencing methods to profile diet does not require high depth

104    to accurately profile the breadth of diet items consumed by an animal. Here we aim to

105    quantify all diet items present in the diet at 1% or more. At this fraction, if we assume

106    that read counts are Poisson distributed, with only 2,000 reads, 99% of the time we will

107    quantify such items within 2-fold of their true amount. Thus, diet quantification

108    presents a clear example of a situation in which sequencing depth is not a critical

109    factor. This contrasts with microbiome profiling, in which there may be very rare taxa

110    (e.g. present at 0.01% frequency) that nevertheless have considerable effects on

111    phenotype of community function.

## Results
### DNA sequencing

114    We selected eight rats from each of three locations near Auckland, New Zealand for diet

115    quantification. Each location comprised a different type of habitat: undisturbed inland

116    native forest (Waitakere Regional Parklands, WP); native bush surrounding an estuary

117    (Okura Bush Walkway, OB); and restored coastal wetland (Long Bay Regional Park,

118    LB). We isolated DNA from whole homogenised stomach contents from each rat (see

119    Methods). We sequenced these DNA samples on two dates by multiplexing the

120    samples, obtaining a total of 82,977 reads (January 2017) and 96,150 reads (March

121    2017). These numbers are not far below expectations given the flow cell and kit

122    chemistry and MinKNOW software versions available at that time (47). However, these

123    read numbers are considerably below those expected for current ONT flow cells and

124    software, which has improved per flow cell output more than 100-fold above these
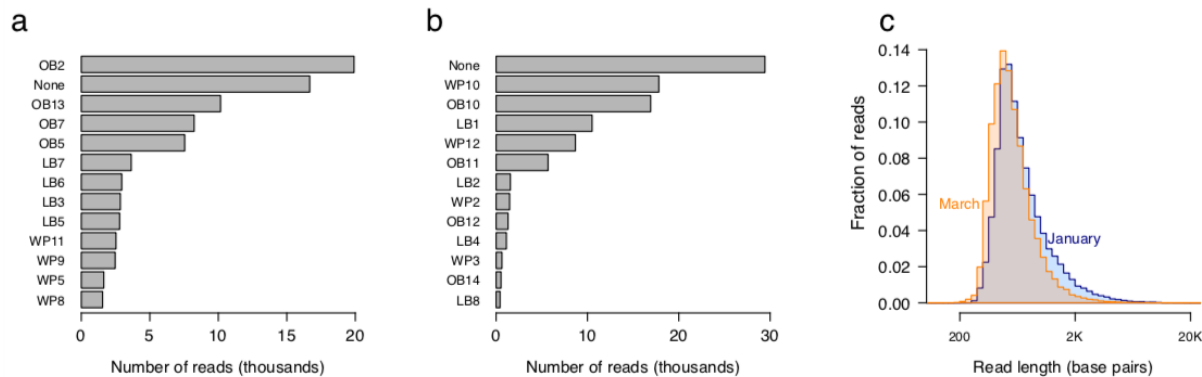
125    numbers.

126    After de-multiplexing the reads, we found large variation in the numbers of reads per

127    multiplex barcode: approximately 10-fold for the January samples, and up to 40-fold in

128    March (**Fig. 1A** and **1B**). We hypothesise that this is due to the highly variable quality of

129    DNA in each sample. This did not appear to have strong effects on read accuracy, as

130    the median quality scores per read ranged from 7-12 (0.80 - 0.94 accuracy) for both

131    runs. This quality is theoretically sufficient for accurate inference of taxa at the Family or

132    even Genus level (32); we investigate this accuracy in analyses below.

133    The degradation of the DNA during digestion in the stomach, as well as fragmentation

134    during DNA isolation (48) and sequencing library preparation led to relatively short

135    median read lengths of 606 bp and 527 bp for the January and March datasets,

136    respectively (**Fig. 1C**). However, there was wide variation in length, with almost 10% of

137    all reads being longer than 1200 bp. Notably, reads of this length can allow more

138    precise taxonomic identification than accurate shorter more accurate reads (e.g.

139    Illumina) (32).

## Assignment of reads to taxa
140
141    To quantify diet contents we first BLASTed all sequences against a combined database

142    of the NCBI nt database (the partially non-redundant nucleotide sequences from all

143    traditional divisions of GenBank excluding genome survey sequence, EST, high-

144    throughput genome, and whole genome shotgun

145    (ftp://ftp.ncbi.nlm.nih.gov/blast/db/README)) and the NCBI other_genomic database

146    (RefSeq chromosome records for non-human organisms

147     (ftp://ftp.ncbi.nlm.nih.gov/blast/db/README)).  We used BLAST as it is generally viewed

148     as the gold standard method in metagenomic analyses (50). Of the 133,022 barcoded

149     reads, 30,535 (23%) hit a sequence in the combined nt and other_genomic database at

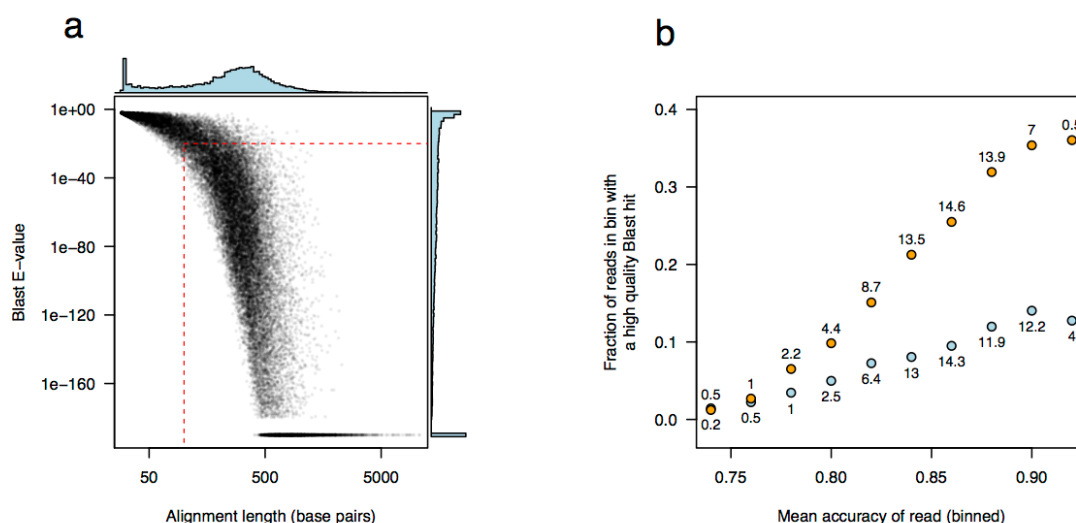150     an e-value cut-off of 1e-2.

151



152

**Fig. 1. Run statistics of nanopore metagenomic sequencing of rat stomach contents. (A) and (B) Barcode distributions for January and March runs, respectively.** We multiplexed the samples on the flow cells, using 12 barcodes per flow cell. The distribution of read numbers across barcodes varied by up to 40-fold. 20% (January) and 30% (March) of all reads could not be assigned to a barcode ("None"). The inability to assign these reads to a barcode is due primarily to their lower quality. **(C) Read length distribution for January and March nanopore runs.** 90% of the reads were between 350 bp and 1,580 bp in length, with only 0.55% being longer than 4,000 bp.

162

163     We first aimed to assess the quality of these hits. We found a bimodal distribution of

164     alignment lengths and e-values (**Fig. 2A**). We also noticed that mean read quality had

165     substantial effects on the likelihood of a read yielding a BLAST hit, with almost 40% of

166     high accuracy reads (greater than 92%) having hits in the March dataset, compared to

167 1% of low accuracy reads (less than 75% accuracy; **Fig. 2B**). We found the same

168 pattern in the January dataset, although to a lesser extent.

169 We hypothesized that many of the short alignments with high e-values were clear false

170 positives given the reported taxa. We thus first filtered the BLAST results, only retaining

171 hits with e-values less than 1e-20 and alignments greater than 100 bp. To give some

172 intuition for this length cut-off, a 100bp read with five 1 bp indels and seven mismatches

173 (88% identity) would have an identical raw score and e-value (given the default match,

174 mismatch, and gap parameters) to a single-end 70 bp Illumina read with a 100%

175 identical database match.  Similar or quality filters based on length and identity have

176 been imposed previously (30). A total of 22,154 hits passed this e-value filter.



177

178 **Fig. 2. BLAST hits of metagenomic reads. (a)** Alignment lengths and e-values were
179 bimodally distributed. The y-axis is plotted on a log scale, with zero e-values
180 suppressed by adding a small number (1e-190) to each e-value. The horizontal red
181 dotted line indicates the e-value cut-off we implemented and the vertical red dotted line
182 indicates the length cut-off (e-value < 1e-20 and alignment length of 100, respectively)
183 to decrease false positive hits. **(b)** The fraction of reads with high quality BLAST hits (e-
184 value < 1e-20) increased as a function of read accuracy. We binned the data according
185 to mean read accuracy (bin width = 0.02) and calculated the fraction of reads within

186 each bin that have a high quality BLAST hit (alignment length greater than 100bp and e-
187 value less than 1e-20) for the January and March runs separately (blue and orange
188 points, respectively). The number of reads in each bin is indicated above each point (in
189 thousands). There is a clear positive correlation between mean accuracy and the
190 likelihood of a high quality BLAST hit, reaching almost 40% for high accuracy reads
191 (>92.5%) for the March dataset.

192

193 We next used MEGAN6 (41) to assign reads to specific taxa. MEGAN6 employs an

194 LCA algorithm to assign reads to a taxon. For example, if a read has BLAST hits to five

195 different species, three of which have bit scores within 20% of the best hit, the read will

196 be assigned to the genus, family, order, or higher taxon level that is the LCA of those

197 best-hit three species (51). If a read matches one species far better than any other, by

198 definition, the LCA is that species.

199 16,820 reads (76%) were assigned to a taxon by MEGAN. Of these, 31% were

200 assigned by MEGAN as being bacterial, and 55% of these were *Lactobacillus spp*.

201 These results match previous studies on rat stomach microbiomes, which have found

202 lactobacilli to be the dominant taxa (52–55). Plant-associated *Pseudomonas*, as well as

203 *Lactococcus* taxa, were also common, at 7% and 6%, respectively.

204 MEGAN assigned reads to a wide range of eukaryotic taxa. To conservatively infer

205 taxon presence, we first reclassified MEGAN species-level assignments to the level of

206 genus. After this, several clear false positive taxon assignments remained (e.g. hippo

207 and naked mole rat). These matches were generally short and of low identity. To reduce

208 such false positive taxon inferences, we used information from reads assigned to the

209 genera *Rattus* (rat) and *Mus* (mouse), using the following strategy.

210    We inferred that the reads assigned to *Rattus* (2,696 reads in total) were true positive

211    genus-level assignments (deriving from DNA isolated from host stomach tissue), and

212    that the reads assigned to *Mus* (2,798 reads in total) were false positive genus-level

213    assignments (i.e. they were derived from Rattus host tissue and not *Mus*-derived). By

214    using host reads, we can implement a ground-truth filtering strategy similar to that

215    achieved by feeding a diet of known content (i.e. rat) and testing whether the contents

216    of the known diet can be accurately identified.

217    First, it is critical to note that the reads assigned to *Mus* are false positive taxon

218    assignments. If these were true positive *Mus* reads, then they would necessarily be due

219    to mouse predation. Although rats are known to prey on mice (56), if this had occurred,

220    we would expect that (assuming only a subset of rats had recently predated mice) if a

221    rat had recently predated a mouse, (1) the ratio of mouse to rat reads would be higher

222    than in rats that had not predated mice; (2) the percent identity of the reads assigned to

223    *Mus* would be higher than in rats that had not predated mice. However, we found that

224    the ratio of mouse to rat reads and percent identity of reads assigned to *Mus* was

225    similar for all rats. This suggested either that all rats had predated mice very recently

226    (which we view as unlikely), or that these *Mus* hits were indeed false positives. Thus,

227    we use the *Mus* hits to delineate false positive and true positive genus-level

228    assignment using the specific read identity and alignment length characteristics of

229    each read set.

230    We first noted that the mean percent identity values of the best BLAST hits for *Rattus*

231    and *Mus* reads differed, with reads matching *Rattus* having a median identity of 86.4%,
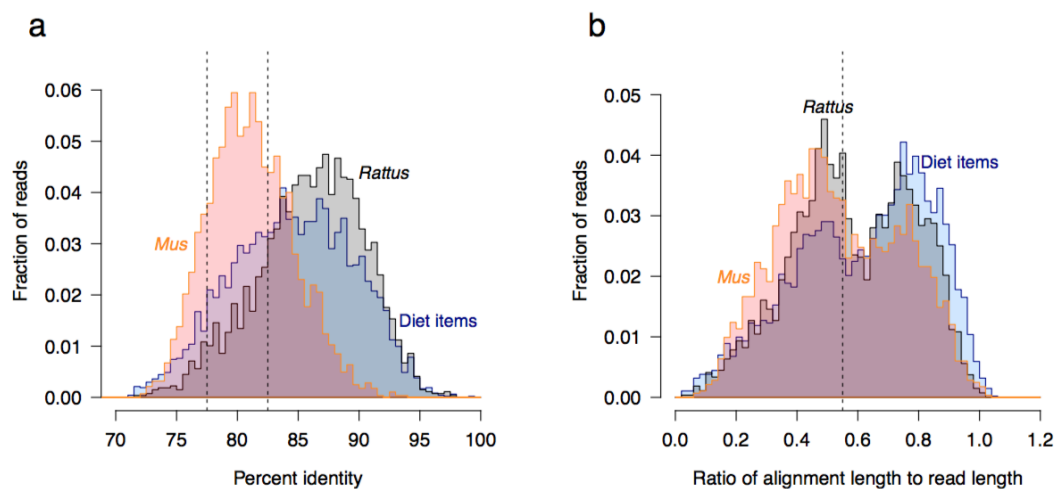
232  and reads matching *Mus* having 81.0% (**Fig. 3A**). The mean percent identity for *Rattus*

233  reads corresponds very well to that expected given the mean quality scores of the

234  reads (86.4% identity corresponds to a mean quality score of 8.7, similar to what we

235  observed; **Fig. S1A-C**).

236  A second characteristic we considered was the ratio of alignment length to read length.

237  If a read fully aligns, this ratio is one. Generally, higher quality alignments should have

238  higher ratios. Indeed, we found this ratio differed substantially between the *Rattus*- and

239  *Mus*-assigned reads: the median ratio of alignment length to read length was longer for

240  *Rattus* (0.57) than *Mus* (0.52; **Fig. 3B**).

241  We used these read characteristics to select cut-off values for assigning reads as

242  being true positive or false positive genus-level taxon assignments. For genus-level

243  assignments, we required at least 82.5% alignment identity, and an alignment length to

244  read length ratio of at least 0.55. For any alignment of lower quality, we assign reads at

245  the Family level. These cut-offs exclude 88% of the reads assigned to *Mus*, instead

246  assigning them one taxon level higher, to the Family *Muridae* (which contains the

247  Genus *Rattus*).

248  However, this rate of false positive assignments (88%) is still relatively high. We thus

249  also quantified the rate of false positive assignment at the level of family. We first

250  identified all reads classified as being from the Order *Rodentia.* Within this Order, we

251  assume that reads assigned to the Family Muridae are true positive, and that reads

252  assigned to any other Family are false positives. This assumption Is conservative when

253  calculating accuracy. We then again implemented specific cut-offs, requiring reads

254    assigned at the Family level to have database matches of at least 77.5% identity, an

255    alignment length to read length ratio of at least 0.1, and a total alignment length of at

256    least 150 bp. With these cut-offs, 97.3% of all reads assigned to the Order *Rodentia*

257    were classified in the Family *Muridae* (which contains the genus *Rattus*). The remaining

258    2.7% were assigned to the family *Cricetidae* (voles and lemmings), except for four

259    reads assigned to *Spalacidae* (mole-rats) (**Fig. S2**). All of these Family assignments are

260    clear false positives, as it is highly unlikely that these families were predated. However,

261    these results establish that by implementing specific cut-off values for



262

263    **Fig. 3. Distributions of percent identity and length for alignments of reads**
264    **matching *Rattus* (rat), *Mus* (mouse), and diet items. (a) The percent identity for**
265    **alignments of rat (*Rattus*) and diet items is much higher than for mouse (*Mus*).**
266    Histograms of the percent identity of the alignment of the top BLAST hit with the read.
267    *Mus* matches have substantially lower percent identity compared to both *Rattus* and diet
268    items. The dotted lines indicate the cut-offs that we implemented for inferring reads as
269    belonging to a specific genus (above 82.5% identity) or family (above 77.5% identity).
270    **(b)  Ratios of alignment lengths to read lengths of rat (*Rattus*) and diet items are**
271    **higher than for mouse (*Mus*).** This plot is analogous to that in (a). The dotted line
272    indicates the cut-off that we implemented for inferring reads as belonging to a specific
273    genus (above 0.55).

274

275    alignments, we can ensure a low rate of false positive assignments at the Family level.

276    However, database bias may still play a role. For example, *Mus* and *Rattus* sequences

277    are among the most common in the database. However, (as above), one expectation if

278    reads are assigned to the wrong taxon is that these false positive assignments would

279    have lower percent identities and low alignment lengths. We thus checked whether

280    read alignments of all inferred diet items had percent identities and alignment lengths

281    similar to the true positive *Rattus* alignments, or instead whether they were more

282    similar to the false positive *Mus* alignments. We found that the majority of diet items

283    had alignment percent identities that overlapped with the *Rattus* reads. Furthermore,

284    the alignment length to read length ratios often exceed those for the *Rattus* reads. This

285    suggests that the diet taxa assignments are often correct down to the level of genus

286    (as the *Rattus*-assigned reads are correct to the level of genus) and are not false

287    positive assignments. Despite this indication of Genus-level accuracy, here we

288    conservatively report diet items at the level of Family.

289    For reads that did not pass the above cut-offs, we placed taxon assignments at the

290    level of order, or used the taxon level assigned by MEGAN. Using these cut-offs, 16%

291    of all reads were classified at the Genus level (although for the analyses below we

292    consider these at the Family-level); 71% were classified at the family-level or below;

293    89% were classified at the order-level or below; and 98% were classified at the phylum-

294    level or below.

295    There were few clear false positive taxon assignments after this analysis, and most had

296    alignments lengths just above our cut-offs (**Table S4**). The exception to this were three

reads from two rats matching *Buthidae* (scorpions), which had alignment lengths of 762

bp, 664 bp, and 298 bp with identities of 83%, 88%, and 79%, respectively. It is unlikely

these are true positives, and instead we hypothesise that these rats predated

harvestmen (*Opiliones*), a closely related sister taxon within *Arachnida*, but lacking

significant amounts of genomic data. Despite the presence of these false positive taxa,

we did not further increase the stringency of our filters, as the fraction was very small.

## Diet quantification

Within each rat, we identified a wide variety of plant, animal, and fungal orders, ranging

from two to 25 Orders per rat (mean 8.7; **Fig. 4**). In total, we identified taxa from 68

different Families, 55 different Orders, 15 different Classes, and eight different Phyla

(**Fig. 5**). This is far beyond the range of diet items that could be identified using a

straightforward metabarcoding approach.

Plants were the primary diet item, with four predominant Orders: *Poales* (grasses),

*Fabales* (legumes), *Arecales* (palms), and *Araucariales* (specifically, *Podocarpaceae*, a

common native New Zealand tree Family). The dominance of plant matter (fruits and

seeds) in rat diets has been established previously (57,58). Animal taxa made up a

smaller component of each rat's diet, with *Insecta* dominating: *Hymenoptera* (bees,

wasps, and ants), *Coleoptera* (beetles), *Lepidoptera* (moths and butterflies), *Blattodea*

(cockroaches), *Diptera* (flies), and *Phasmatodea* (stick insects). In addition,

*Stylommatophora* (slugs and snails) were present in substantial numbers (**Fig. 5A** and

**5B**). Fungi were only a small component of the rats' diet, although several orders were

present: *Sclerotiniales* (commonly plant pathogens), *Saccharomycetales* (budding

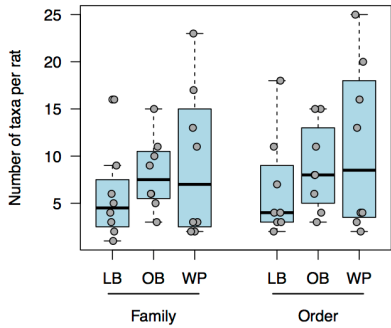yeasts), *Mucorales* (pin molds), *Russulales* (brittlegills and milk-caps), and

320 *Chytotheriales* (black yeasts)*. Finally, for many rats, a substantial proportion of the

321 stomach contents were parasitic worms (primarily *Spirurida* (nematodes) and

322 *Hymenolepididae* (tapeworms)).

323 It is important to note that due to our metagenomic approach, the fraction of each

324 element of the rats' diets may be distorted by biases in genomic databases: whole

325 genome data exists for only a few taxa, while mtDNA, rDNA, metabarcode loci, and

326 microsatellite sequence data are present in the database for many animal and plant

327 genera. However, it is possible to mitigate this bias.

328 To quantify database-driven bias, for each taxon we determined the fraction of hits that

329 mapped to mtDNA, rDNA, microsatellites, or EST libraries (we refer to this as *non-*

330 *genomic*, as these data are not from genomic sequencing projects). We also

331 determined the fraction of hits that mapped to DNA sequences arising from genome

332 sequencing projects). We expect that for animals with sequenced genomes, these two

333 fractions should be primarily determined by the relative amounts of mtDNA and nuclear

334 DNA in a diet item, rather than database bias. If a diet item consists of cells that have

335 large numbers of mitochondria (or if the animal has a small genome), we expect a large

336 fraction of reads will map to mtDNA sequences. Alternatively, if a diet item consists of

337 cells with few mtDNA, then most reads will map to genomic sequence. However, for

338 animals without sequenced genomes, there should be considerably more hits to

339 mtDNA, plastid, rDNA, and microsatellites (non-genomic sequence), and few (if any)

340 genomic hits, regardless of the relative amounts of mtDNA and nuclear DNA in the diet

341 item. By comparing these fractions for taxa that we know to have complete genome

342   sequences in the database to taxa without complete genomes we aimed to assess and

343   mitigate the effects of this bias. For this analysis, we consider Genus-level assignments.

344



345

346   **Fig. 4. Numbers of taxa in individual rats.** Each boxplot indicates the range of
347   families (left boxes) or orders (right boxes) consumed by each rat in each location (OB:
348   Okura Bush, native bush; LB: Long Bay Park, restored wetland; WP: Waitakere Park,
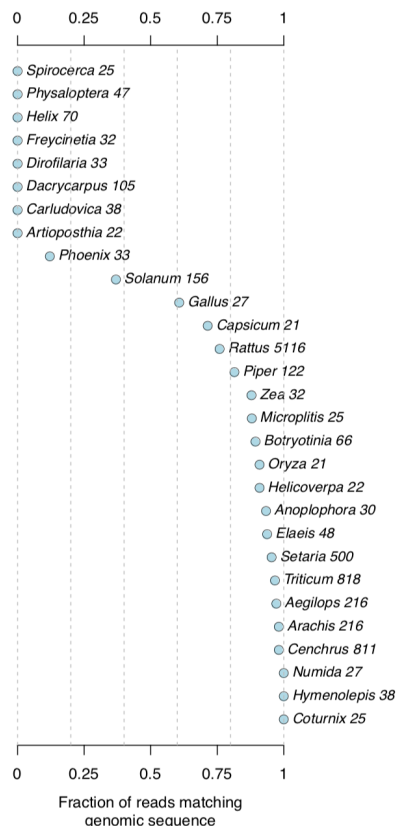349   native forest). The numbers for individual rats (eight per location) are plotted in grey.
350



351

352

**Fig. 5. Proportions of taxa in the diets of individual rats. (a) Reads assigned to taxa at the family and (b) order level.** The rows correspond to a single rat, with the proportions of reads for that rat assigned to each family or order indicated in shades of blue and yellow. Reads that were not assigned to a specific family or order are indicated at the right side of the figure. The families and orders have been sorted so that the most common diet components appear on the left. Only the 55 most common families are shown. Note that the color gradations presented on the scale are not linear.

360

For animal Genera with at least of species that had sequenced genomes in the database, we found that the fraction of reads that mapped to non-genomic sequence (mtDNA, rDNA, microsatellite, plastid, and EST library) ranged from 0% (*Coturnix* (quail)) to 39% (*Gallus* (chicken)) (**Fig. 6**). This is a considerable range, and we hypothesise that variation in the fraction of non-genomic reads is due to the type of tissue sequenced (affecting mtDNA content). However, it is also possible that the results are biased by the relative frequency of specific types of studies, such as those generating microsatellite data. For *Rattus* (27% non-genomic), the sequenced tissue was primarily stomach muscle, which has a relatively high fraction of mtDNA, perhaps explaining the large fraction of reads mapping to non-genomic DNA.

371   For plants with sequenced genomes, the fraction of reads matching non-genomic

372   sequence (mostly mtDNA, plastid, and rDNA) was generally lower: between 2%

373   (*Cenchrus* buffelgrass)) and 12% (*Zea* (corn)). Thus, on average, for animals with

374   sequenced genomes present in the database, approximately 30% of all reads mapped

375   to non-genomic sequences; for plants, approximately 5% mapped to non-genomic

376   sequences.

377   In contrast, for taxa with little or no genomic sequence in the database, the vast majority

378   of matches were non-genomic (mtDNA, plastid, rDNA, or microsatellite loci): 90% of

379   *Phoenix* (date palm) hits; all *Helix* (snail); and all *Rhaphidophorae* (endemic cave weta)

380   hits. All *Arthurdendyus* (endemic New Zealand flatworm) hits were solely to rDNA loci.

381   We note that these data indicate the accuracy of the read classifications, as several of

382   these are endemic New Zealand species.

383   These ratios are in strong contrast to animals with sequenced genomes, for which ana

384   average of only 30% of all reads should map to non-genomic sequence. This suggests

385   that for animal taxa with little or no genomic sequence data, we have underestimated

386   the actual number of sequences from that taxon by two- to three-fold. For plant taxa

387   with little or no genomic sequence data, we have underestimated read abundance by

388   approximately 20-fold. In terms of diet biomass, there is considerable uncertainty in both

389   of these estimates.

390

**Figure 6.** Fractions of reads matching genomic and non-genomic sequence for the best BLAST hit of each read. For the species with complete genomes, the fraction of reads matching genomic sequence ranges from 40% (*Solanum*) to 100%. This large range is likely due to the tissue from which the DNA was isolated. For example, muscle tissue has a higher fraction of mtDNA to nuclear DNA than egg. For species without fully sequenced genomes, this fraction ranges from 0% to 20% (*Phoenix*, which has a small amount of genomic data present in the database).

It was apparent after read classification, specific taxa were overrepresented in the diets

of rats from particular locations. For example, six out of eight rats from the native

estuarine bush habitat (OB) consumed *Arecaceae*, while only one in the restored

wetland area (LB) did. All three rats that consumed *Phaseanidae* were from the native

estuarine habitat (OB). All five rats that consumed *Solanales* were from the restored

404    wetland area. These patterns suggested that it might be possible to use diet

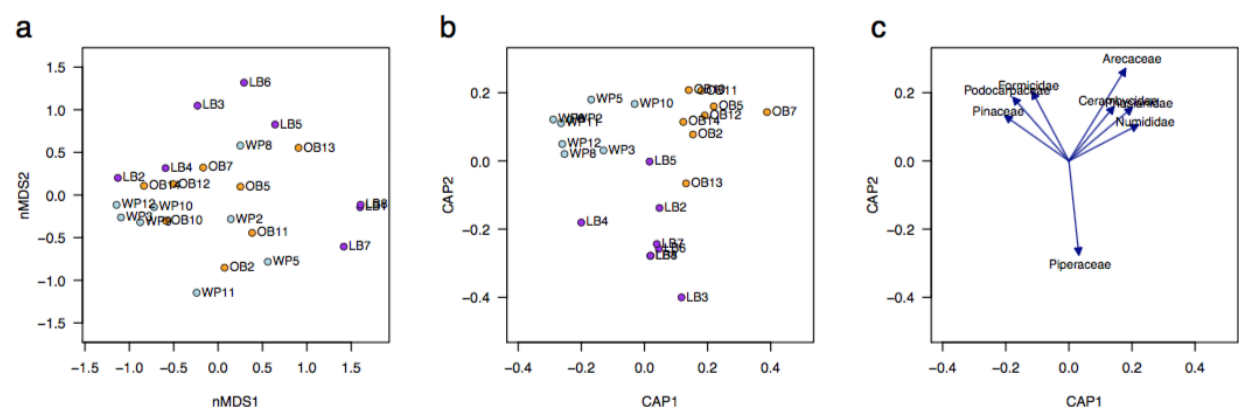405    components alone to pinpoint the habitat from which each rat was sampled.

## nMDS and CAP analysis by location

407    In order to determine if diet composition of the rats differed consistently between

408    locations, we first performed an unconstrained analysis using non-metric

409    multidimensional scaling (nMDS) on taxa assigned at the family level. The input for the

410    nMDS was the dissimilarity matrix (Bray-Curtis distance of diet at the family level).

411    NMDS uses rank-based distances to cluster samples that are most similar.

412    The family-level unconstrained ordination (nMDS) showed no obvious grouping of rats

413    with respect to the locations (**Fig. 7A**), indicating that locations did not correspond to the

414    predominant axes of variation among the diets. We next performed a constrained

415    ordination method, Canonical Analysis of Principal coordinates (CAP, see Methods).

416    CAP identifies axes of variation, if any, that distinguish the diets of rats from different

417    locations (**Fig. 7B**). We found that the CAP axes correctly classified the locations of 19

418    out of 24 (79%) rats using a leave-one-out procedure. The families having the largest

419    correlations with the first two principal coordinates, and thus most responsible for the

420    separation between groups, were primarily plants: *Arecaceae*, *Podocarpaceae*,

421    *Piperaceae*, and *Pinaceae*. In addition, insect groups (*Cerambycids* and *Formicids*) and

422    birds (*Phaseanidae* and *Numididae*) played a role (**Fig. 7C**).

423    The families driving similarity within the three locations (i.e., those that had the greatest

424    within-location SIMPER scores, see Methods) varied among locations. LB had average

425    Bray-Curtis within-location similarity of 13%; mostly attributable to *Hymenolepidae*

426     (accounting for 51% of the within-group similarity), *Solanaceae* (11%), and *Fabaceae*

427     (11%). The average similarity for OB was 21%, with the greatest contributing taxa

428     being *Arecaceae* (33%), *Poaceae* (23%), *Fabaceae* (9%), and *Phasianidae* (8%). The

429     average similarity for WP was 24%, with the greatest contributing taxon being *Poaceae*

430     (72%) (**Table S4**).



431

**Fig. 7. Unconstrained nMDS (a) and constrained CAP (b) ordinations of the diets**
**of rats from three locations. Both ordinations were based on Bray-Curtis**
**dissimilarities of square root transformed proportions of reads attributed to each**
**family.** The locations were a native estuarine bush (OB, orange); a restored marine
wetland (LB, purple); and a native forest (WP, light blue). The CAP ordination is
repeated in panel **(c)** as a biplot with the rats omitted to show the Pearson correlations
between families and the first two CAP axes. The eight families with the strongest
correlations are shown, indicating the taxa associated with each location.

## Discussion

### Accuracy and sensitivity

443     Here we have shown that using a simple metagenomic approach with error-prone long

444     reads allows rapid and accurate classification of rat diet components (approximately

445     2.7% error in taxon assignments at the family-level). We expect that this technique can

446     be used to infer diet for a wide variety of animal and sample types, including samples

447     that use less invasive collection methods, such as fecal matter. The accuracy of this

448     approach will likely improve as the accuracy and yield of ONT sequencing continues to

449     increase. The analysis here is based on fewer than 200,000 reads from two flow cells.

450     Current yields for similar read length distributions are in excess of ten million reads per

451     flow cell. As ONT modal sequencing accuracy is currently just above 96%, and

452     continues to improve. This increase in read accuracy will clearly affect the accuracy of

453     taxon assignment, illustrated by the fact that the fraction of reads yielding BLAST hits

454     increases substantially for high accuracy reads, approaching 40% for high quality

455     reads in our dataset (reads with greater than 92.5% accuracy, **Fig. 2B)**. With current

456     ONT sequencing techniques, 92.5% is at the lower end of read accuracy.

457     Furthermore, as the species sampling of genomic databases increases (59), the taxon-

458     level precision of this method will improve. Given the current rate of genomic

459     sequencing, with careful sampling, the vast majority of multicellular plant and animal

460     families (and even genera) will likely have at least one type species with a sequenced

461     genome within the next decade. Continued advancement in sequence database search

462     algorithms as compared to current methods (23,24,60) should considerably decrease

463     the computational workload necessary to find matching sequences.

464     ## Methodological advantages

465     As genomic databases become more complete, metagenomic approaches will offer

466     significant advantages due to decreased bias as compared to other methods. We found

467     that rats consumed many soft-bodied species (e.g. mushrooms, flat worms, slugs, and

468     lepidopterans) that would be difficult to identify using visual inspection of stomach

469  contents.  Achieving data on such a wide variety of taxa (across multiple phyla) would

470  also be difficult to quantify using metabarcoding, as there are no universal 18S or COI

471  universal primers capable of amplifying sequences in all these taxa. While it might be

472  possible to use several different primer sets targeted at different phyla or orders,

473  quantitatively comparing diet components across these using sequences amplified with

474  different primer sets is extremely difficult due to differences in primer binding and PCR

475  efficiency.

476  The ONT-based sequencing method has several unique advantages. Perhaps the most

477  obvious is the accessibility of the platform. Compared to other high throughput

478  sequencing technologies (e.g. Illumina, IonTorrent, or PacBio), there is no initial capital

479  investment required. On a per-sample basis, data generation is inexpensive (assuming

480  12 multiplexed samples, approximately $150 USD per sample, and half this price if

481  reagents are purchased in bulk). Library preparation and sequencing can be extremely

482  rapid, going from DNA sample to sequence in less than two hours (61). Furthermore,

483  the sequencing platform itself is highly portable. Given (1) that ONT-based methods are

484  now similar in cost-per-read as the most accessible Illumina method (we estimate $650

485  for 10 million reads using ONT, versus $1300 for 20 million reads using MiSeq); and (2)

486  that even marginal increases in read length are likely to significantly improve species

487  identification, we expect that ONT-based methods should soon become useful for

488  routine ecological monitoring of species (62).

489  Some modifications to our approach might further increase the precision of our ability to

490  infer community composition. Any error-prone long read dataset (i.e. PacBio or ONT)

491     has both short (e.g. 500 bp) and long (e.g. 5,000 bp) reads, as well as high quality (e.g.

492     mean accuracy greater than 90%) and low quality (e.g. mean accuracy less than 80%)

493     reads. When inferring community composition, a null expectation is that taxa should be

494     equally represented by long, high quality reads as they are by short, low quality reads. If

495     some taxa are represented only by short, low quality reads, this suggests that these

496     taxa may be false positive inferences. Similarly, the difficulty in correctly mapping short

497     inaccurate reads could be mitigated by weighting the probability of taxon mapping by

498     the number of long, accurate reads that map to certain taxa. Thus, the fact that not all

499     reads are extremely long and accurate does not mean that they cannot all be used to

500     infer taxon presence in metagenomic analyses.

501     Finally, it is critical to note that for many diet studies, the aim is to resolve biomass,

502     nutritional content, or prey numbers. However, estimating these numbers is constrained

503     by the fact that different tissues and different taxa have different amounts of DNA (both

504     nuclear and mitochondrial) per gram of biomass. It is nearly to impossible to fully

505     account for this variation using any DNA-based method. Regardless, there is

506     considerable utility in using DNA-based approaches for diet assessment, not least

507     because it is one of the few methods that allows the full breadth of the diet to be

508     observed, as illustrated here by the number of different orders we find that rats predate.
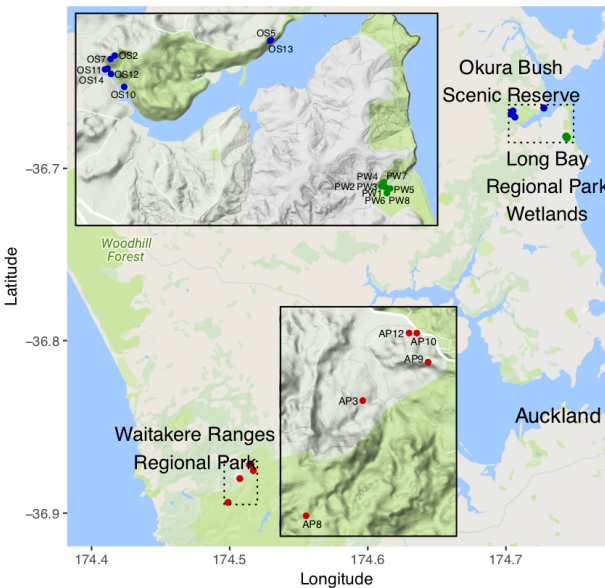
## Conclusions

509

510     Here we have shown that using a rapid error-prone long read metagenomic approach

511     we can accurately characterise diet taxa at the Family level, and distinguish between

512     the diets of rats according to the locations from which they were sourced. This

513     information may be used to guide conservation efforts toward specific areas and

514   habitats in which native species are most at risk from this highly destructive introduced

515   predator.

## Methods
## Study Areas

518   We trapped rats from three locations near Auckland, New Zealand. Each location

519   comprised a different type of habitat: undisturbed inland native forest (Waitakere

520   Regional Parklands, WP); native bush surrounding an estuary (Okura Bush Walkway,

521   OB); and restored coastal wetland (Long Bay Regional Park, LB) (**Fig. 8**).  Snap traps in

522   OB and LB were baited with peanut butter, apple, and cinnamon wax pellets; or bacon

523   fat and flax pellets.



524
**Fig. 8. Location of rat sampling sites in the greater Auckland area in the North Island of New Zealand**. Each point indicates a trap where one rat was captured, with the colour of the points indicating the three broad locations: the native estuarine bush habitat of Okura Bush (OB), the restored wetland of Long Bay (LB), and the native forest of Waitakere Park (WP). The two insets show the three locations in higher resolution with topographical details. Green indicates park areas. Precise geographical coordinates were only available for five out of eight rats in WP.

532

533     Traps in WP were baited with chicken eggs, rabbit meat, or cinnamon scented poison

534     pellets. From 16 November to 16 December 2016, traps were surveyed by established

535     conservation groups at each site every 48 hours. A total of 36 rats were collected from

536     these locations. Three of the rats from WP had poison in their stomachs. However, all

537     three of these were killed by the snap traps, and as such are unlikely to have swallowed

538     any bait. In addition, none of these rats were identified as having chicken or rabbit in

539     their diet. The majority of rats collected (34/36) were determined to be male *Rattus*

540     *rattus* by visual inspection. These 34 rats were selected for further analysis.

541     ## DNA Isolation

542     Within 48 hours of trapping, rats were stored at either -20°C or -80°C until dissection.

543     We dissected out intact stomachs from each animal and removed the contents. After

544     snap freezing in liquid nitrogen, we homogenised the stomach contents using a sterile

545     mini blender to ensure sampling was representative of the entire stomach.

546     We purified DNA from 20 mg of homogenised stomach contents using the Promega

547     Wizard Genomic DNA Purification Kit, with the following modifications to the Animal

548     Tissue protocol: after protein precipitation, we transferred the supernatant to a new tube

549     and centrifuged a second time to minimise protein carryover. The DNA pellet was

550     washed twice with ethanol. These modifications were performed to improved DNA

551     purity. We rehydrated precipitated DNA by incubating overnight in molecular biology

552     grade water at 4°C and stored the DNA at -20°C. DNA quantity, purity, and quality was

553     ascertained by nanodrop and agarose gel electrophoresis. The DNA samples were

554     ranked according quantity and purity (based on A260/A280 and secondarily, A230/A280

555    ratios). The eight highest quality DNA samples from each of the three locations were

556    selected for sequencing. We did no size selection on the purified DNA.

## DNA Sequencing

558    Sequencing was performed on two different dates (24 January 2017 and 17 March

559    2017) using a MinION Mk1B device and R9.4 chemistry. For each sequencing run, DNA

560    from each rat was barcoded using the 1D Native Barcoding Kit (Barcode expansion kit

561    EXP-NBD103 with sequencing kit SQK-LSK108) following the manufacturer's

562    instructions. This included an AMPure bead purification step to remove adaptors, which

563    also likely removed very short reads (less than 200 bp; see **Fig. 1A**). Twelve samples

564    were pooled and run on each flow cell, for a total of 24 individual rats. The flow cells had

565    1373 active pores (January 2017) and 1439 active pores (March 2017). Both runs were

566    re-basecalled after data collection using Albacore 2.2.7 with demultiplexing performed in

567    Albacore and filtering disable*d (optio*ns *--barcoding  --disable_filtering*).

## Sequence classification

569    All sequences were BLASTed (blastn v2.6.0+) against a locally compiled database

570    consisting of the combined NCBI other_genomic and nt databases (downloaded on 13th

571    June 2018 from NCBI). Default blastn parameters were used (match 2, mismatch -3,

572    gapopen -5, gapextend -2). Due to the predominance of short indels present in

573    nanopore sequence data, we used an initial set of basecalled data to test whether

574    changing these default penalties affected the results (gapopen -1, gapextend -1). We

575    found that these adjusted parameters did not qualitatively change our results.

576 We assigned sequence reads to specific taxon levels using MEGAN6 (v.6.11.7 June

577 2018) (41). We only used reads with BLAST hits having an e-value of $1\times10^{-20}$ or lower

578 (corresponding to a bit score of 115 or higher given the databases we used) and an

579 alignment length of 100 base pairs or more. To assign reads to taxon levels, we

580 considered all hits having bit scores within 20% of the bit score of the best hit (MEGAN

581 parameter Top Percent).

## Multivariate analyses

583 Multivariate analyses were done using the software PRIMER v7 (42). The data used in

584 the multivariate analyses were in the form of a sample- (i.e. individual rat) by-family

585 matrix of read counts. All bacteria, rodent, and primate families were removed as these

586 are not diet content. The majority of the primate hits (32 in total) were assigned to

587 Hominidae (19), which likely resulted from sample contamination (**Table S3**).

588 The read counts were converted to proportions per individual rat by dividing by the

589 total count for each rat, to account for the fact that the number of reads varied

590 substantially among rats (43). The proportions were then square-root transformed so

591 that subsequent analyses were informed by the full range of taxa, rather than just the

592 most abundant families (44). We then calculated a matrix of Bray-Curtis dissimilarities,

593 which quantified the difference in the gut DNA of each pair of rats based on the

594 square-root transformed proportions of read counts across families (43).

595 We used unconstrained ordination, non-metric multidimensional scaling (nMDS)

596 applied to the dissimilarity matrix to examine the overall patterns in the diet

597 composition among rats. To assess the degree to which the diet compositions of rats

598   were distinguishable among the three locations, we applied canonical analysis of

599   principal coordinates (CAP) (45) to the dissimilarity matrix. CAP is a constrained

600   ordination which aims to find axes through multivariate data that best separates *a priori*

601   groups of samples (in this case, the groups are the locations from which the rats were

602   sampled); CAP is akin to linear discriminant analysis but it can be used with any

603   resemblance matrix. The out-of-sample classification success was evaluated using a

604   leave-one-out cross-validation procedure (45).

605   We used Similarity Percentage (SIMPER; (46)) to characterise and distinguish between

606   the locations. This allowed us to identify the families with the greatest percentage

607   contributions to (1) the Bray-Curtis similarities of diets within each location (**Table S5**)

608   and (2) the Bray-Curtis dissimilarities between each pair of locations (**Table S6**).

## Declarations

### Ethics approval

611   Sample collection was performed under (Auckland Council Permit to Undertake

612   Research WS1064).

### Availability of data and materials

614   Sequence data are available in the SRA archive (accession number PRJEB27647).

### Competing interests

616   WP received funding from Oxford Nanopore Technologies (1000$USD) to present this

617   work at a conference (Ecological Society of Australia 2018).

## Funding

## Authors' contributions

WP, JD, NF, and OS conceived the project. WP performed the stomach dissections. WP and NF optimised the genomic DNA isolation and library preparation. NF performed the nanopore sequencing. GB and OS processed and performed quality control on the sequencing data. WP and OS performed the sequence classification. WP, AS, NF, and OS analysed the data. WP, NF, AS, and OS wrote the paper, with input from all authors.

## Acknowledgements

## References

1. Daniel MJ. Seasonal Diet Of The Ship Rat (Rattus Rattus) In Lowland Forest In New Zealand. Proc. 1973;20:21–30.

2. Pierce GJ, Boyle. A review of methods for diet analysis in piscivorous marine mammals. Oceanogr Mar Biol Annu Rev. 1991;29:409–86.

3. Major HL, Jones IL, Charette MR, Diamond AW. Variations in the diet of introduced Norway rats (Rattus norvegicus) inferred using stable isotope analysis. J Zool. 2007 Apr 1;271(4):463–8.

4. Carreon-Martinez L, Heath DD. Revolution in food web analysis and trophic ecology: diet analysis by DNA and stable isotope analysis. Mol Ecol. 2010

642        Jan;19(1):25–7.

5.    Hobson KA. Use of stable-carbon isotope analysis to estimate marine and terrestrial protein content in gull diets. Can J Zool. 1987 May 1;65(5):1210–3.

6.    Basha WA, Chamberlain AT, Zaki ME, Kandeel WA, Fares NH. Diet reconstruction through stable isotope analysis of ancient mummified soft tissues from Kulubnarti (Sudanese Nubia). Journal of Archaeological Science: Reports. 2016 Feb 1;5:71–9.

7.    Dunlap M, Pawlik JR. Video-monitored predation by Caribbean reef fishes on an array of mangrove and reef sponges. Mar Biol. 1996 Mar 1;126(1):117–23.

8.    Brown KP, Moller H, Innes J, Jansen P. Identifying predators at nests of small birds in a New Zealand forest. Ibis . 2008;140(2):274–9.

9.    King RA, Read DS, Traugott M, Symondson WOC. Molecular analysis of predation: a review of best practice for DNA-based approaches. Mol Ecol. 2008 Feb;17(4):947–63.

10.  Soininen EM, Valentini A, Coissac E, Miquel C, Gielly L, Brochmann C, et al. Analysing diet of small herbivores: the efficiency of DNA barcoding coupled with high-throughput pyrosequencing for deciphering the composition of complex plant mixtures. Front Zool. 2009 Aug 20;6:16.

11.  Jarman SN, Gales NJ, Tierney M, Gill PC, Elliott NG. A DNA-based method for identification of krill species and its application to analysing the diet of marine vertebrate predators. Mol Ecol. 2002 Dec;11(12):2679–90.

12.  Jarman SN, Deagle BE, Gales NJ. Group-specific polymerase chain reaction for DNA-based analysis of species diversity and identity in dietary samples. Mol Ecol. 2004 May;13(5):1313–22.

13.  Tedersoo L, Anslan S, Bahram M, Põlme S, Riit T, Liiv I, et al. Shotgun metagenomes and multiple primer pair-barcode combinations of amplicons reveal biases in metabarcoding analyses of fungi. MycoKeys. 2015;10:1.

14.  Leray M, Yang JY, Meyer CP, Mills SC, Agudelo N, Ranwez V, et al. A new versatile primer set targeting a short fragment of the mitochondrial COI region for metabarcoding metazoan diversity: application for characterizing coral reef fish gut contents. Front Zool. 2013 Jun 14;10:34.

15.  Pawluczyk M, Weiss J, Links MG, Egaña Aranguren M, Wilkinson MD, Egea-Cortines M. Quantitative evaluation of bias in PCR amplification and next-generation sequencing derived from metabarcoding samples. Anal Bioanal Chem. 2015 Mar;407(7):1841–8.

16.  Pereira RPA, Peplies J, Brettar I, Hoefle MG. Impact of DNA polymerase choice on

677 assessment of bacterial communities by a Legionella genus-specific next-
678 generation sequencing approach [Internet]. bioRxiv. 2018 [cited 2018 Jan 31]. p.
679 247445. Available from:
680 https://www.biorxiv.org/content/early/2018/01/12/247445.abstract

681 17. Hover BM, Kim S-H, Katz M, Charlop-Powers Z, Owen JG, Ternei MA, et al.
682 Culture-independent discovery of the malacidins as calcium-dependent antibiotics
683 with activity against multidrug-resistant Gram-positive pathogens. Nat Microbiol.
684 2018 Apr;3(4):415–22.

685 18. Anantharaman K, Brown CT, Hug LA, Sharon I, Castelle CJ, Probst AJ, et al.
686 Thousands of microbial genomes shed light on interconnected biogeochemical
687 processes in an aquifer system. Nat Commun. 2016 Oct 24;7:13219.

688 19. Xu Z, Knight R. Dietary effects on human gut microbiome diversity. Br J Nutr. 2015
689 Jan;113 Suppl:S1–5.

690 20. Fierer N, Leff JW, Adams BJ, Nielsen UN, Bates ST, Lauber CL, et al. Cross-biome
691 metagenomic analyses of soil microbial communities and their functional
692 attributes. Proc Natl Acad Sci U S A. 2012 Dec 26;109(52):21390–5.

693 21. Breitwieser FP, Salzberg SL. KrakenHLL: Confident and fast metagenomics
694 classification using unique k-mer counts. bioRxiv [Internet]. 2018; Available from:
695 https://www.biorxiv.org/content/early/2018/06/06/262956.abstract

696 22. Huson DH, Mitra S, Ruscheweyh H-J, Weber N, Schuster SC. Integrative analysis
697 of environmental sequences using MEGAN4. Genome Res. 2011 Sep;21(9):1552–
698 60.

699 23. Kim D, Song L, Breitwieser FP, Salzberg SL. Centrifuge: rapid and sensitive
700 classification of metagenomic sequences. Genome Res. 2016 Dec;26(12):1721–9.

701 24. Wood DE, Salzberg SL. Kraken: ultrafast metagenomic sequence classification
702 using exact alignments. Genome Biol. 2014 Mar 3;15(3):R46.

703 25. Paula DP, Linard B, Crampton-Platt A, Srivathsan A, Timmermans MJTN, Sujii ER,
704 et al. Uncovering Trophic Interactions in Arthropod Predators through DNA
705 Shotgun-Sequencing of Gut Contents. PLoS One. 2016 Sep 13;11(9):e0161841.

706 26. Srivathsan A, Ang A, Vogler AP, Meier R. Fecal metagenomics for the
707 simultaneous assessment of diet, parasites, and population genetics of an
708 understudied primate. Front Zool. 2016 Apr 21;13:17.

709 27. Arribas P, Andújar C, Hopkins K, Shepherd M, Vogler AP. Metabarcoding and
710 mitochondrial metagenomics of endogean arthropods to unveil the mesofauna of
711 the soil. Yu D, editor. Methods Ecol Evol. 2016 Sep 7;7(9):1071–81.

28. Linard B, Crampton-Platt A, Gillett CPDT, Timmermans MJTN, Vogler AP. Metagenome Skimming of Insect Specimen Pools: Potential for Comparative Genomics. Genome Biol Evol. 2015 May 14;7(6):1474–89.

29. Meiser A, Otte J, Schmitt I, Grande FD. Sequencing genomes from mixed DNA samples - evaluating the metagenome skimming approach in lichenized fungi. Sci Rep. 2017 Nov 2;7(1):14881.

30. Srivathsan A, Sha JCM, Vogler AP, Meier R. Comparing the effectiveness of metagenomics and metabarcoding for diet analysis of a leaf-feeding monkey (Pygathrix nemaeus). Mol Ecol Resour. 2015 Mar;15(2):250–61.

31. Søe MJ, Nejsum P, Seersholm FV, Fredensborg BL, Habraken R, Haase K, et al. Ancient DNA from latrines in Northern Europe and the Middle East (500 BC–1700 AD) reveals past parasites and diet. PLoS One. 2018 Apr 25;13(4):e0195481.

32. Pearman WS, Freed NE, Silander OK. Testing the advantages and disadvantages of short-and long-read eukaryotic metagenomics using simulated reads. BMC Bioinformatics 21, 1-15

33. Hillmann B, Al-Ghalith GA, Shields-Cutler RR, Zhu Q, Gohl DM, Beckman KB, et al. Evaluating the Information Content of Shallow Shotgun Metagenomics. mSystems [Internet]. 2018 Nov;3(6). Available from: http://dx.doi.org/10.1128/mSystems.00069-18

34. Gibbs GW. Why are some weta (Orthoptera: Stenopelmatidae) vulnerable yet others are common? J Insect Conserv. 1998 Dec 1;2(3-4):161–6.

35. Towns DR, Daugherty CH, Cree A. Raising the prospects for a forgotten fauna: a review of 10 years of conservation effort for New Zealand reptiles. Biol Conserv. 2001 May 1;99(1):3–16.

36. Stringer IAN, Bassett SM, McLean MJ, McCartney J, Parrish GR. Biology and conservation of the rare New Zealand land snail Paryphanta busbyi watti (Mollusca, Pulmonata). Invertebr Biol. 2003 Sep 1;122(3):241–51.

37. Diamond JM, Veitch CR. Extinctions and introductions in the new zealand avifauna: cause and effect? Science. 1981 Jan 30;211(4481):499–501.

38. Dowding JE, Murphy EC. The impact of predation by introduced mammals on endemic shorebirds in New Zealand: a conservation perspective. Biol Conserv. 2001 May 1;99(1):47–64.

39. Graham NAJ, Wilson SK, Carr P, Hoey AS, Jennings S, MacNeil MA. Seabirds enhance coral reef productivity and functioning in the absence of invasive rats. Nature. 2018 Jul;559(7713):250–3.

747  40. Russell JC, Innes JG, Brown PH, Byrom AE. Predator-Free New Zealand:
748       Conservation Country. Bioscience. 2015 May 1;65(5):520–5.

749  41. Huson DH, Beier S, Flade I, Górska A, El-Hadidi M, Mitra S, et al. MEGAN
750       Community Edition - Interactive Exploration and Analysis of Large-Scale
751       Microbiome Sequencing Data. PLoS Comput Biol. 2016 Jun;12(6):e1004957.

752  42. Clarke KR, Gorley RN. PRIMER v7: User Manual/Tutorial. PRIMER-E, Plymouth;
753       2015 p. 296.

754  43. Clarke KR, Robert Clarke K, Somerfield PJ, Gee Chapman M. On resemblance
755       measures for ecological studies, including taxonomic dissimilarities and a zero-
756       adjusted Bray–Curtis coefficient for denuded assemblages. J Exp Mar Bio Ecol.
757       2006;330(1):55–80.

758  44. Clarke KR, Green RH. Statistical Design and Analysis for a "biological Effects"
759       Study. Mar Ecol Prog Ser. 1988;46:213–26.

760  45. Anderson MJ, Willis TJ. Canonical Analysis Of Principal Coordinates: A Useful
761       Method Of Constrained Ordination For Ecology. Ecology. 2003 Feb 1;84(2):511–
762       25.

763  46. Clarke KR. Non-parametric multivariate analyses of changes in community
764       structure. Austral Ecol. 1993;18(1):117–43.

765  47. Jain M, Koren S, Miga KH, Quick J, Rand AC, Sasani TA, et al. Nanopore
766       sequencing and assembly of a human genome with ultra-long reads. Nat
767       Biotechnol. 2018 Apr;36(4):338–45.

768  48. Deagle BE, Eveson JP, Jarman SN. Quantification of damage in DNA recovered
769       from highly degraded samples--a case study on DNA in faeces. Front Zool. 2006
770       Aug 16;3:11.

771  49. Jain M, Olsen HE, Paten B, Akeson M. The Oxford Nanopore MinION: delivery of
772       nanopore sequencing to the genomics community. Genome Biol. 2016 Nov
773       25;17(1):239.

774  50. McIntyre ABR, Ounit R, Afshinnekoo E, Prill RJ, Hénaff E, Alexander N, et al.
775       Comprehensive benchmarking and ensemble approaches for metagenomic
776       classifiers. Genome Biol. 2017 Dec 21;18(1):182.

777  51. Huson DH, Auch AF, Qi J, Schuster SC. MEGAN analysis of metagenomic data.
778       Genome Res. 2007 Mar;17(3):377–86.

779  52. Maurice CF, Knowles SCL, Ladau J, Pollard KS, Fenton A, Pedersen AB, et al.
780       Marked seasonal variation in the wild mouse gut microbiota. ISME J. 2015
781       Nov;9(11):2423–34.

782   53.  Brownlee A, Moss W. The influence of diet on lactobacilli in the stomach of the rat.
783         J Pathol. 1961 Oct 1;82(2):513–6.

784   54.  Li D, Chen H, Mao B, Yang Q, Zhao J, Gu Z, et al. Microbial Biogeography and
785         Core Microbiota of the Rat Digestive Tract. Sci Rep. 2017 Apr 4;8:45840.

786   55.  Horáková Z, Zierdt CH, Beaven MA. Identification of lactobacillus as the source of
787         bacterial histidine decarboxylase in rat stomach. Eur J Pharmacol. 1971
788         Sep;16(1):67–77.

789   56.  Bridgman LJ, Innes J, Gillies C, Fitzgerald N, King CM. Do ship rats display
790         predatory behaviour towards house mice? Anim Behav. 2013 Aug 1;86(2)):257–68.

791   57.  Sweetapple PJ, Nugent G. Ship rat demography and diet following possum control
792         in a mixed podocarp—hardwood forest. N Z J Ecol. 2007;31(2):186–201.

793   58.  Riofrío-Lazo M, Páez-Rosas D. Feeding Habits of Introduced Black Rats, Rattus
794         rattus, in Nesting Colonies of Galapagos Petrel on San Cristóbal Island,
795         Galapagos. PLoS One. 2015 May 18;10(5):e0127901.

796   59.  Lewin HA, Robinson GE, Kress WJ, Baker WJ, Coddington J, Crandall KA, et al.
797         Earth BioGenome Project: Sequencing life for the future of life. Proc Natl Acad Sci
798         U S A. 2018 Apr 24;115(17):4325–33.

799   60.  Nasko DJ, Koren S, Phillippy AM, Treangen TJ. RefSeq database growth
800         influences the accuracy of k-mer-based species identification. 2018;1–21.

801   61.  Zaaijer S, Gordon A, Speyer D, Piccone R, Groen SC, Erlich Y. Rapid re-
802         identification of human samples using portable DNA sequencing. Elife [Internet].
803         2017 Nov 28;6. Available from: http://dx.doi.org/10.7554/eLife.27798

804   62.  Kamenova S, Bartley TJ, Bohan DA, Boutain JR, Colautti RI, Domaizon I, et al.
805         Chapter Three - Invasions Toolkit: Current Methods for Tracking the Spread and
806         Impact of Invasive Species. In: Bohan DA, Dumbrell AJ, Massol F, editors.
807         Advances in Ecological Research. Academic Press; 2017. p. 85–182.

808

809 Supplemental Tables

810 **Table S1**. Read numbers and total base pairs for each barcode in the January 2017
811 sequencing run.

| Sample | Total reads | Total Mbp | Mean length |
|--------|-------------|-----------|-------------|
| **OB2** | 19907 | 14.62 | 734 |
| **WP11** | 10164 | 9.63 | 947 |
| **WP5** | 8237 | 6.78 | 823 |
| **LB7** | 7548 | 7.04 | 933 |
| **OB13** | 3644 | 3.63 | 995 |
| **WP9** | 2954 | 2.4 | 814 |
| **OB5** | 2850 | 2.06 | 721 |
| **WP8** | 2801 | 2.32 | 827 |
| **LB6** | 2531 | 1.6 | 632 |
| **OB7** | 2473 | 1.87 | 756 |
| **LB5** | 1641 | 1.16 | 705 |
| **LB3** | 1554 | 0.99 | 636 |
| **None** | 16673 | 13.01 | 781 |
| **Total** | 82977 | 67.1 | 809 |

812

813

814 **Table S2**. Read numbers and total base pairs for each barcode in the March 2017
815 sequencing run.

| Sample | Total reads | Total Mbp | Mean length |
|--------|-------------|-----------|-------------|
| LB1 | 17820 | 9.21 | 517 |
| LB8 | 16923 | 13.13 | 776 |
| WP2 | 10511 | 7.00 | 666 |
| LB4 | 8684 | 4.92 | 567 |
| OB11 | 5689 | 3.40 | 598 |
| WP10 | 1563 | 0.99 | 633 |
| OB12 | 1479 | 0.89 | 604 |
| WP12 | 1309 | 0.78 | 596 |
| LB2 | 1127 | 0.76 | 676 |
| WP3 | 637 | 0.73 | 1141 |
| OB14 | 541 | 0.37 | 683 |
| OB10 | 435 | 0.24 | 555 |
| None | 29432 | 21.33 | 725 |
| Total | 96150 | 63.75 | 663 |

816

817

818 **Table S3. Characteristics of alignments for reads assigned to the Primate family.**
819 Many reads are both long and have high identity, suggesting that they are not false
820 positive assignments, but contamination.

| Rat | Read length | Mean read accuracy | % ID | Alignment length | Genus |
|---|---|---|---|---|---|
| WP10 | 428 | 0.91 | 94.0 | 314 | Homo |
| WP10 | 782 | 0.92 | 93.8 | 657 | Homo |
| OB5 | 365 | 0.90 | 93.2 | 249 | Homo |
| LB3 | 515 | 0.95 | 92.6 | 462 | Homo |
| WP10 | 510 | 0.90 | 92.0 | 460 | Homo |
| WP10 | 704 | 0.90 | 91.1 | 587 | Homo |
| WP10 | 467 | 0.89 | 90.3 | 402 | Homo |
| WP10 | 494 | 0.89 | 89.7 | 388 | Homo |
| WP10 | 339 | 0.88 | 89.6 | 269 | Homo |
| WP10 | 446 | 0.88 | 89.6 | 326 | Homo |
| WP10 | 327 | 0.91 | 89.1 | 210 | Homo |
| OB5 | 415 | 0.89 | 88.8 | 277 | Homo |
| OB5 | 561 | 0.86 | 88.7 | 257 | Homo |
| WP11 | 434 | 0.86 | 88.7 | 301 | Homo |
| WP10 | 486 | 0.88 | 88.2 | 365 | Homo |
| WP10 | 613 | 0.90 | 88.2 | 526 | Homo |
| WP10 | 563 | 0.87 | 88.2 | 457 | Homo |
| WP10 | 1373 | 0.91 | 87.7 | 1337 | Homo |
| WP11 | 526 | 0.91 | 87.5 | 473 | Homo |
| OB14 | 478 | 0.89 | 86.9 | 373 | Homo |
| OB5 | 715 | 0.89 | 86.8 | 673 | Homo |
| WP10 | 475 | 0.87 | 86.7 | 362 | Homo |
| WP10 | 398 | 0.86 | 86.5 | 259 | Homo |
| WP10 | 377 | 0.88 | 85.3 | 251 | Homo |
| WP9 | 558 | 0.87 | 85.2 | 508 | Homo |
| WP10 | 429 | 0.86 | 84.8 | 276 | Homo |
| WP10 | 322 | 0.81 | 84.5 | 174 | Homo |
| WP8 | 723 | 0.84 | 83.1 | 438 | Homo |
| LB8 | 965 | 0.86 | 80.4 | 245 | Rhinopithecus |

| LB5 | 3042 | 0.94 | 79.4 | 1018 | Cebus |
|-----|------|------|------|------|-------|
| WP2 | 464  | 0.93 | 77.3 | 216  | Homo  |
| LB8 | 671  | 0.90 | 73.9 | 406  | Aotus |

821

822

823 **Table S4. Characteristics of alignments for reads that are likely false positive**
824 **assignments.** Most are short long or have low identity, suggesting that they are false
825 positive assignments. The exception are the reads matching Buthidae, which we
826 hypothesize are due to the rats predation of the sister taxa, harvestmen. Octopodidae,
827 Salmonidae, and Poeciliidae (guppies and similar aquaria fish) are possible but
828 improbable prey items.

| Rat | Read length | Mean read accuracy | % ID | Alignment length | Genus | Megan family |
|-----|-------------|--------------------|------|------------------|-------|--------------|
| WP5 | 1285 | 0.90 | 82.9 | 298 | Centruroides | Buthidae |
| OB11 | 1874 | 0.90 | 88.0 | 664 | Centruroides | Buthidae |
| WP5 | 1711 | 0.93 | 79.3 | 762 | Centruroides | Buthidae |
| WP11 | 859 | 0.93 | 86.1 | 151 | Octopus | Octopodidae |
| WP2 | 516 | 0.86 | 81.4 | 172 | Oncorhynchus | Salmonidae |
| OB12 | 424 | 0.90 | 85.7 | 140 | Xiphophorus | Poeciliidae |
| OB12 | 643 | 0.84 | 89.3 | 177 | Xiphophorus | Poeciliidae |

829

830 **Table S5.** SIMPER analysis of family contributions to group similarities.

| Family | Average | Average | Similarity/SD | Percentage | Group |
|--------|---------|---------|---------------|------------|-------|
| **Hymenolepididae** | 3.37 | 6.87 | 0.34 | 51.2 | LB |
| **Solanaceae** | 1.57 | 1.48 | 0.34 | 11.1 | LB |
| **Fabaceae** | 1.74 | 1.41 | 0.44 | 10.5 | LB |
| **Arecaceae** | 2.86 | 7.11 | 1 | 33.4 | OB |
| **Poaceae** | 2.87 | 4.82 | 0.55 | 22.7 | OB |
| **Fabaceae** | 1.17 | 1.98 | 0.51 | 9.3 | OB |
| **Phasianidae** | 1.79 | 1.67 | 0.34 | 7.9 | OB |
| **Poaceae** | 5.08 | 17.61 | 0.62 | 72.1 | WP |

831

832

833   **Table S6.** SIMPER analysis of family contributions to group dissimilarities.

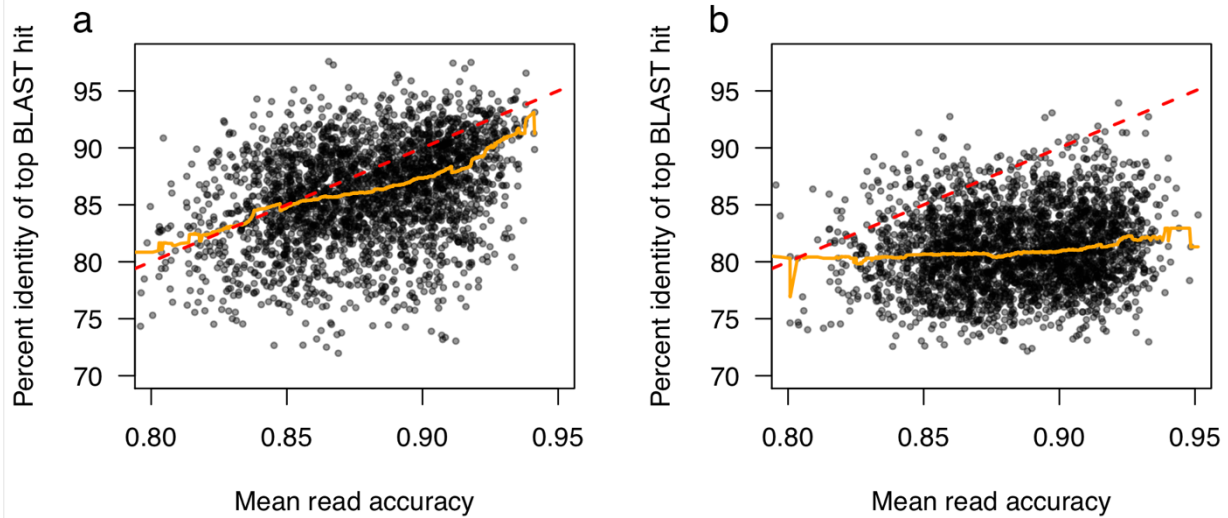| Species | Av Abund Group1 | Avg Abund Group2 | Avg. Diss | Diss/SD | % contrib | Group 1 | Group2 |
|---|---|---|---|---|---|---|---|
| **Poaceae** | 1.95 | 5.08 | 15.15 | 1.04 | 16.74 | LB | WP |
| **Poaceae** | 2.87 | 5.08 | 11.29 | 1.26 | 13.78 | OB | WP |
| **Hymenolepididae** | 3.37 | 0.48 | 10.8 | 0.73 | 11.93 | LB | WP |
| **Hymenolepididae** | 3.37 | 0.29 | 9.37 | 0.79 | 10.32 | LB | OB |
| **Poaceae** | 1.95 | 2.87 | 8.37 | 1.1 | 9.22 | LB | OB |
| **Arecaceae** | 0.05 | 2.86 | 6.99 | 1.41 | 7.7 | LB | OB |
| **Arecaceae** | 2.86 | 1.31 | 5.92 | 1.29 | 7.23 | OB | WP |
| **Fabaceae** | 1.74 | 1.05 | 6.14 | 0.67 | 6.78 | LB | WP |
| **Podocarpaceae** | 0 | 2.38 | 5.34 | 0.83 | 5.9 | LB | WP |
| **Podocarpaceae** | 0.71 | 2.38 | 4.82 | 0.99 | 5.88 | OB | WP |
| **Fabaceae** | 1.74 | 1.17 | 4.87 | 0.81 | 5.37 | LB | OB |
| **Fabaceae** | 1.17 | 1.05 | 4.31 | 0.84 | 5.26 | OB | WP |

834

835   **Datafile S1**. Table of read BLAST hits and assigned MEGAN taxa with reads
836   reclassified at the family or order level by filtering on read length to alignment length
837   ratio and percent identity.

838   **Datafile S2**. Table of read BLAST hits and assigned MEGAN taxa with no filters
839   applied.

840

Supplemental Figures

**Fig S1. Correlation of read accuracy with alignment characteristics**. Only rat reads exhibit a clear positive relationship between accuracy and percent identity. (A) indicates the relationship for reads assigned to *Rattus* and (B) for *Mus*. The orange line indicates a running median; the red dotted line is the y=x line, which is expected if accuracy corresponds exactly to percent identity.

**Fig S2. Alignment characteristics of true positive and false positive taxon assignments at the family level**. False positive taxon assignments (*Cricetideaa*, orange, and *Spalacidae*, green) have lower percent identity and shorter alignment lengths than true positive taxon assignments (*Muridae*, black). Only a single false positive taxon assignment had a read length to alignment length ratio greater than 0.5 and a percent identity greater than 85%. This suggests that with further filters or methodologies (e.g. decision tree analysis using different read and alignment characteristics) could, if necessary, decrease false positive rates even further.