1  **Title:** Targeted sequence capture outperforms RNA-Seq and degenerate-primer PCR cloning for

2  sequencing the largest mammalian multi-gene family

3  **Authors:** Laurel R. Yohe[1,2], Kalina T. J. Davies[3], Nancy B. Simmons[4], Karen E. Sears[5], Elizabeth R.

4  Dumont[6], Stephen J. Rossiter[3], & Liliana M. Dávalos[1,7]

5  **Affiliations**:

6  [1] Department of Ecology and Evolution, Stony Brook University, Stony Brook, NY 11794, USA

7  [2] Department of Geology and Geophysics, Yale University, Stony Brook, NY 11794, USA

8  [3] School of Biological and Chemical Sciences, Queen Mary University of London, London, E1 4NS,

9  United Kingdom

10 [4] Department of Mammalogy, Division of Vertebrate Zoology, American Museum of Natural History,

11 New York, NY 10024, USA

12 [5] Department of Ecology and Evolutionary Biology, UCLA, Los Angeles, CA 90095, USA

13 [6] School of Natural Sciences, University of California Merced, Merced, CA, 95343, USA

14 [7] Consortium for Inter-Disciplinary Environmental Research, Stony Brook University, Stony Brook, NY

15 11794, USA

16 **Email addresses**:

17 LRY: laurel.yohe@yale.edu; laurel.yohe@stonybrook.edu

18 KJD: k.t.j.davies@qmul.ac.uk

19 NBS: simmons@amnh.org

20 KES: ksears@ucla.edu

21 ERD: edumont@ucmerced.edu

22 SJR: s.j.rossiter@qmul.ac.uk

23 LMD: liliana.davalos@stonybrook.edu

24 **Data availability:** RNA-seq raw reads were deposited to the NCBI GenBank Sequence Read Archive

25 (SRR8878915) and the assembled transcriptome was deposited to the NCBI GenBank Transcriptome

26 Shotgun Assembly database (PRJNA531931).

27      **Short Title:** Sequencing approaches of multigene families

28      **Keywords**: transcriptome, targeted sequence capture, olfactory receptor, multigene family, gene family

29      evolution, genome

30      **Corresponding Author**:

31      Laurel R. Yohe

32      210 Whitney Ave, New Haven, CT 06511

33      phone: (631) 632-8600; fax: (631) 632-7626

34      laurel.yohe@yale.edu

35

**Abstract**

Multigene families evolve from single-copy ancestral genes via duplication, and typically encode proteins critical to key biological processes. Molecular analyses of these gene families require high-confidence sequences, but the high sequence similarity of the members can create challenges for both sequencing and downstream analyses. Focusing on the common vampire bat, *Desmodus rotundus*, we evaluated how different sequencing approaches performed in recovering the largest mammalian protein-coding multigene family: *olfactory receptors* (*OR*). Using the common vampire bat genome as a reference, we determined the proportion of putatively protein-coding receptors recovered by: 1) amplicons from degenerate primers sequenced via Sanger technology, 2) RNA-Seq of the main olfactory epithelium, and 3) those genes "captured" with probes designed from transcriptomes of closely-related species. Our initial re-annotation of the high-quality vampire bat genome resulted in >400 intact *OR* genes, more than double the number based on original estimates. Sanger-sequenced amplicons performed the poorest among the three approaches, detecting <33% of receptors in the genome. In contrast, the transcriptome reliably recovered >50% of the annotated genomic *ORs*, and targeted sequence capture recovered nearly 75% of annotated genes. Each sequencing approach assembled high-quality sequences, even if it did not recover all putative receptors in the genome. Therefore, variation among assemblies was caused by low coverage of some receptors, rather than high rates of assembly error. Given this variability, we caution against using the counts of number of intact receptors per species to model the birth-death process of multigene families. Instead, our results support the use of orthologous sequences to explore and model the evolutionary processes shaping these genes.

56 **Introduction**

57      Multigene families, or groups of duplicated genes that have evolved from a single ancestral copy,

58 make up significant proportions of the protein-coding genome across organisms. Many of these gene

59 families underlie key roles in sensory perception and pathogen recognition (Nei and Rooney 2005; Yohe

60 *et al.* 2019). However, despite both the biological relevance and prevalence of multigene families, most

61 sequencing and assembly methods are optimized for single-copy genes. Since assembling highly similar

62 sequences is inherently problematic, many multigene families assemble poorly (Sims *et al.* 2014; Shi *et*

63 *al.* 2017). Duplicated genes are often masked from analyses and ignored, or recent gene duplicates may

64 be collapsed into single-copy genes, thus underestimating their diversity (MacRander *et al.* 2015; Holding

65 *et al.* 2018). Mapping reads back onto assembled contigs of duplicated genes is an error-prone task,

66 making it difficult to validate a well-assembled contig (Treangen and Salzberg 2012). This problem is

67 particularly evident in *de novo* assemblies, for which no reference genomes are available to validate

68 scaffolds. Even when the genome of a closely related species is available, these regions of the genome

69 may still be poorly assembled, and their highly repetitive nature results in misleading coverage estimates

70 (Yoon *et al.* 2009; Sims *et al.* 2014). While all of these issues are well known, there are few comparisons

71 of different sequencing methods and their performance in reconstructing high quality contigs from highly

72 similar sequences.

73      The mammalian *olfactory receptor* (*OR*) gene family shows one of the most extraordinary

74 patterns of gene duplication in animals (Nei and Rooney 2005), constantly expanding through

75 duplications and contracting via pseudogenization over time. Olfactory receptors account for 5% of the

76 mammalian protein-coding genome (Niimura 2012). Olfactory receptors are short ~900 basepair (bp)

77 intronless G-protein coupled receptors with divergent binding sites that reflect the diversity of potential

78 odorants to which they bind (Niimura 2012). The variation in the number of receptors among mammals is

79 enormous—humans have around 400 ORs in their genome, while rodents and elephants have thousands

80 (Niimura *et al.* 2014). Mammalian olfactory receptors can be classified into distinct subfamilies based on

81 conserved regions of the genes (Hayden *et al.* 2010). Class I receptors are shared across vertebrates and

82    can be further subdivided into four subfamilies (51, 52, 55, 56), while the much more diverse Class II

83    subfamilies are mammalian-specific and subdivided into nine subfamilies (1/3/7, 2/13, 4, 5/8/9, 6, 10, 11,

84    12, and 14) (Hayden *et al.* 2010, 2014). Because of the duplicative nature of these genes, olfactory

85    receptors pose a challenge to sequence assemblers. It is critical to obtain reliable olfactory receptor

86    sequences to infer gene duplication and loss, and even for comparing the size of repertoires across

87    species. Within a population, the sensitivity to odorant stimuli of the same receptor with segregating

88    alleles is highly variable (Logan 2014; Mainland *et al.* 2014). Thus, accurate and reliable sequences are

89    also necessary for identifying within-population evolutionary processes that shape chemosensory

90    receptors.

91         Despite some of the potential problems that emerge from sequencing olfactory receptors, this task

92    can become tractable with use of proper methodologies and access to genomic resources. Many olfactory

93    receptors have been identified from available genomes (Niimura *et al.* 2014), but when a reference

94    genome is unavailable, alternative approaches must be considered. One such approach is to use a set of

95    degenerate primers to amplify sequences using PCR, followed by cloning and Sanger sequencing

96    (Hayden *et al.* 2010, 2014). While Sanger sequencing has a very low error rate, primer bias (caused by the

97    preferential binding of degenerate primers to some genes over others), insufficient sampling of clones, or

98    insufficient sequencing depth (due to the relatively high cost per base) may limit complete recovery of the

99    profiles from amplicons (Hayden *et al.* 2010; Hohenbrink *et al.* 2014). By using degenerate primers with

100   paired-end sequencing platforms such as Illumina (Hughes *et al.* 2013), or with long read technologies

101   such as PacBio (Larsen *et al.* 2014), it may be possible to increase the number of recovered chemosensory

102   receptors, however, such high-throughput approaches can introduce higher sequencing error rates without

103   resolving the problems arising from primer bias. Transcriptomes and targeted sequence capture offer

104   alternatives to avoid primer bias or insufficient sequencing. When pooling data from multiple individuals,

105   for example, studies of the mammalian olfactory transcriptome in model organisms detected up to 95% of

106   intact olfactory receptors. However, all of these studies used well-annotated reference genomes to guide

107   their assemblies (Shiao *et al.* 2012; Kanageswaran *et al.* 2015; Olender *et al.* 2016). How *de novo*

108  olfactory sequencing assemblies perform in recovery of the hyperdiverse mammalian olfactory receptor

109  repertoire remains unknown.

110       Here we compare variation in olfactory receptors of *Desmodus rotundus*, the common vampire

111  bat, recovered from different high-throughput sequencing approaches. *Desmodus rotundus* is the only

112  vertebrate that feeds exclusively on mammalian blood and, in accordance with its dietary preference, has

113  highly modified sensory systems including thermosensation (Gracheva *et al.* 2011), reduced taste function

114  (Hong and Zhao 2014), and distinct olfactory receptors (Hayden *et al.* 2014) . Our main goal is to

115  determine whether different sequencing approaches can yield representative samples of highly similar

116  protein-coding genes, even in the absence of a reference genome, and to identify the best assembly

117  approach to achieve this goal. Our analyses use as a baseline the genome of this species to identify open

118  reading frames of olfactory receptors (Zepeda Mendoza *et al.* 2018). We then compare three sequencing

119  strategies: published olfactory receptor sequences amplified via degenerate primers and cloning, and

120  sequenced using Sanger technology (Hayden *et al.* 2014), *de novo* transcriptome sequences of the main

121  olfactory epithelium, and targeted sequence capture using probes designed from the transcriptomes of

122  twelve bat species. To characterize the completeness and sensitivity of different assembly strategies for

123  one of the most complex gene families in the mammalian genome, we mapped the receptors to the

124  genome. We discovered significant variation across methods and suggest best practices for subsequent

125  analyses based on different sequencing and assembly approaches, with implications for downstream

126  analyses of multigene family evolution.

127

128  **Materials and Methods**

129  *Approach*: The following sequencing approaches were compared to assess their ability to recover

130  maximum representation of high quality olfactory receptor contigs: (1) PCR with degenerate primers and

131  Sanger sequencing of amplicons (Hayden *et al.* 2014), (2) receptors obtained from Illumina sequencing of

132  the transcriptome of the main olfactory epithelium, and (3) receptors sequenced from targeted sequence

133  capture  with  probes  designed  from  olfactory  receptors  identified  in  bat  olfactory  epithelium

6

134    transcriptomes (Fig. 1). Because of the duplicative nature of olfactory receptors, careful consideration was

135    given to designing the pipeline for Illumina read quality control and assembly. Reads that are too short,

136    too low in quality, or do not have a matching pair, may confound the assembly. The published common

137    vampire bat genome (*Desmodus rotundus*) served as a validation of correctly assembled olfactory

138    receptors (Zepeda Mendoza *et al.* 2018). The genome was sequenced using Illumina, and after refinement

139    by the Dovetail protocol, resulted in ~2Gb genome with a mean coverage of ~233X and a final N50 =

140    26.9 Mb. Each assembly approach was compared to the genome by mapping assembled contigs to the

141    olfactory receptor locations in the genome.

142    *Tissue collection*: For RNA-Seq, we generated the tissue-specific transcriptome of the main olfactory

143    epithelium (MOE) from one male *D. rotundus* (AMNH 278722), collected in Lamanai (Belize) under the

144    Belize Forestry Department Scientific Research and Collecting Permit CD/60/3/14 (17) and protocols

145    approved by Institutional Animal Care and Use Committee at Stony Brook University (IACUC: 2012-

146    1946-NF-4.16.15-BAT). The bat was euthanized using an overdose of isofluorane and the maxilla, that

147    contains the entire nasal cavity, was immediately removed from the specimen and placed in a vial of

148    Qiagen RNAlater and left to soak overnight at 4˚C to allow complete permeation of the tissue. The

149    following morning, the tissue vial was flash-frozen in liquid nitrogen. Upon returning to the laboratory,

150    the MOE was dissected in sterile conditions on a dry ice cold counter-top under a dissecting scope and

151    under the guidance of a published video protocol (Brechbühl *et al.* 2011). RNA was immediately

152    extracted after dissection.

153    *RNA extractions*: All RNA extractions were performed using the Qiagen RNeasy Micro Kit (ID: 74004)

154    and followed the protocol for "Purification of Total RNA from Animal and Human Tissues". The

155    following modifications were made to optimize the total RNA from the delicate neural tissue of the MOE.

156    We added 20 µL of 2M dithiothreitol (DTT) per 1 mL of the lysate buffer, Buffer RLT. Prior to tissue

157    homogenization, we also added 5 µL of a 4 ng/µL working solution of carrier RNA, as total RNA yields

158    of neural tissue are generally low (Qiagen RNeasy Micro Handbook). A sterilized glass mortar and pestle

159    was used for tissue disruption and homogenization by grinding for 5 minutes in Buffer RLT and carrier

160     RNA. The tissues were homogenized by pumping the pestle and shearing the cellular components.

161     Incubation of the spin columns during the DNase treatment was reduced from 15 to eight minutes.

162     Finally, during the final extraction step, we eluted with 20 µL of RNase-free water and let the water soak

163     on the spin column membrane for 5 minutes prior to elution.

164     *cDNA library sequencing*: RNA extracts were sent to BGI in China for cDNA library preparation and

165     Illumina sequencing. RNA concentration, quality, and purity were measured using the Agilent 2100

166     Bioanalyzer. cDNA libraries were generated using standard BGI in-house protocols. Libraries were

167     sequenced using Illumina HiSeq™ 4000 to generate 6G of 100 bp paired-end reads per sample.

168     *RNA-Seq assembly*: Using the BBTools bioinformatics package (https://sourceforge.net/projects/bbmap/),

169     low quality reads were filtered using the bbduk.sh script, in which reads less than 25 bp (minlen = 25)

170     were discarded. Reads were trimmed from both ends (qtrim = rl) until the average read quality was 10 or

171     greater (trimq = 10); otherwise, the read was discarded. All other settings for this function were set to

172     defaults. To assemble the RNA-Seq data *de novo*, the Oyster River Protocol v. 2.1.0 was implemented

173     (MacManes 2018). This recently developed assembly strategy uses several assembly programs under a

174     variety of different parameters to overcome the biases incurred by different assembly algorithms (Vijay *et*

175     *al.* 2013). The Oyster River Protocol streamlines this approach and provides different benchmarking

176     measures to evaluate the quality of each transcript assembled, as well as overall assembly quality

177     assessment. Briefly, the protocol performs the following analyses: (1) additional trimming and error

178     correction; (2) assembly using Trinity v. 2.8.4 (Grabherr *et al.* 2011), Trans-Abyss v. 2.0.1, and SPAdes

179     v. 3.13.0; (3) merging of assemblies via OrthoFinder v. 2.2.6 (Emms and Kelly 2015); and (4) assembly

180     evaluation using TransRate v. 1.0.2 (Smith-Unna *et al.* 2016) and BUSCO v. 3.0.1 (Waterhouse *et al.*

181     2017). The overall TransRate score is calculated using the product of the four following measures: the

182     proportion of nucleotides with zero coverage, how the bases are ordered correctly based on information

183     from read pairs, how well the nucleotides of mapped reads match those in the assembled contig, and

184     univariate coverage depth that quantifies the probability all reads come from the same transcript. Mapping

185     reads back to the transcriptome can be particularly problematic for duplicated genes, and this TransRate

8

186   score identifies particularly questionable assembled contigs. BUSCO measures the completeness of each

187   assembled contig by searching for orthologous annotated proteins and measuring the standard deviation

188   of each transcript contig from its reciprocal hit in the ortholog database. The assembled transcriptome was

189   compared against an ortholog database for mammals that includes 4,104 BUSCO groups

190   (http://busco.ezlab.org/).

191   *Olfactory receptor identification*: A published pipeline, Olfactory Receptor family Assigner (ORA) v.

192   1.9.1, tailored to specifically identify mammalian olfactory receptors and classify each receptor into its

193   respective subfamily (Hayden *et al.* 2010) was used to characterize the olfactory receptors of each

194   sequencing approach. ORA is a set of Bioperl (v. 1.006924) scripts that implement hidden Markov

195   models trained on conserved protein sequence motifs of mammalian olfactory receptors via HMMER v.

196   3.1b2 (Eddy 2010). This method has been shown to be robust, with low false positives rates, and has been

197   used to identify olfactory receptors and their open reading frames across mammals, including bats

198   (Hayden *et al.* 2010, 2014). An E-value threshold of 1e-10 for sequences matched in the database was

199   used. For the transcripts, we discarded all olfactory receptor sequences with open reading frames <650

200   bp, as it is impossible to distinguish transcribed pseudogenes from degraded transcripts at short lengths.

201   *Comparison with Sanger-sequenced OR amplicons*: A previous study amplified the olfactory receptors of

202   *D. rotundus* using PCR with two pairs of degenerate primers (for Class I and Class II *OR* genes), isolated

203   each gene by cloning, and sequenced the receptors using Sanger sequencing (Hayden *et al.* 2014). Given

204   the low error rates of Sanger sequencing, this provided an opportunity to explore the different methods for

205   sequencing olfactory receptors, and to assess whether the higher error rates of Illumina significantly

206   affected sequencing of *OR* genes. As the degenerate primers bind to conserved regions within the reading

207   frame, recovered sequences were incomplete, only ~700-750 bp (Fig. 1).  The amplicons were obtained

208   using degenerate primers and only a few clones were selected. Since the olfactory receptor repertoire may

209   be quite large, the amplification step has the potential to introduce primer bias, which is then exacerbated

210   by reduced representation.

9

211   *Targeted sequence capture from genomic DNA:* Olfactory receptors identified from the transcriptome

212   were used to design probes for an olfactory receptor targeted sequence capture. Pooling 3,814

213   chemosensory genes from twelve species of bats (Table S1), probes were designed from RNA-Seq data to

214   make 120-bp probes with 2X tiling density. The initial raw number of probes was 45,052, and given the

215   duplicative nature of the genes, we clustered similar probes with 95% nucleotide identity of one another.

216   The final probe count was 16,468 custom targets designed for chemosensory genes. All but one species of

217   bat used in the probe design were sampled from the Noctilionoidea superfamily, a monophyletic clade

218   that shared a common ancestor within the last 40Ma. Probes were designed and synthesized by Arbor

219   Biosciences (Ann Arbor, Michigan) using myBaits technology; they also performed library preparation,

220   target enrichment and oversaw sequencing of the resultant products. To avoid unfair bias, as different

221   individuals were used for the Sanger-sequenced amplicons and genomic datasets, a different *D. rotundus*

222   individual than the one used for the transcriptome was also sequenced here. DNA was extracted from

223   liver tissue sampled from a bat obtained in La Selva, Costa Rica in 2014 (Permit: R-018-2013-OT-

224   CONAGEBIO; IACUC: 2013-2034-R1-4.15.16-BAT) using the DNeasy Blood and Tissue Kit Protocol

225   from Qiagen (69504). Target sequences captured by the probes were sequenced using Illumina

226   sequencing technology following enrichment. Reads were first trimmed for quality using the same

227   bbduk.sh script from the transcriptome assembly and exact duplicate reads were removed using ParDRe

228   v. 2.2.5. To assemble the reads into receptor contigs, a target from the probe design were used to map and

229   align reads with HybPiper v.1.2 (Johnson *et al.* 2016) reads_first.py pipeline with the "-bwa" option

230   selected.

231   *Genome mapping and recovery sensitivity analyses*: Mapping to the same location in the genome was

232   used to assess whether the same receptor was recovered in sequencing and assembly approaches. We

233   mapped all identified olfactory receptors from RefSeq sequences to the *D. rotundus* genome using GMAP

234   v. 2017-01-14 (Wu and Watanabe 2005). We first indexed the genome with gmap_build using a kmer

235   value of 12. We then identified the olfactory receptor coding sequences from the genome using the ORA

236   pipeline, and mapped the identified genomic olfactory receptors back to the genome with GMAP. The

237    mapping yielded genomic scaffold coordinates of the olfactory receptors in the genome to be compared

238    against the location of the receptors from other assembly methods. Only coordinates of genomic receptors

239    that mapped with 100% identity were used. In contrast, Sanger-sequenced amplicons, transcriptome

240    receptors, and receptors assembled from targeted sequence capture were mapped using GMAP, with

241    settings for which there was at least 50% overlap with the receptor coordinates in the genome to account

242    for partially assembled receptors to map. We allowed for mappings with 95% identity, as this was the

243    average sequence nucleotide identity of post-duplication olfactory receptors within mammalian olfactory

244    subfamilies (Hughes *et al.* 2018).  Receptors sometimes mapped with different quality values, to multiple

245    locations in the genome, or in a chimeric fashion, thus a threshold for true mappings was set. If a receptor

246    mapped to multiple locations, the location with the highest sequence identity and mapping quality was

247    used. Receptor mapping localities that intersected with those in the genome were determined using the

248    "intersectBed" in bedtools v. 2.26.0 (Quinlan 2014).

249        We performed a sensitivity analysis to quantify the recovery of all assembled olfactory receptors.

250    Some receptors recovered in each sequencing approach mapped to the genome, but to locations not yet

251    annotated. Thus, there were more olfactory receptors discovered than were previously identified in the

252    published genomic protein-coding sequences for *D. rotundus*. Any receptor from any method that mapped

253    to the genome was considered a "true positive". A receptor that was present in the genome, but not found

254    in another method was considered a "false negative".  Specificity in this case should be interpreted with

255    caution, as there is no variation between sequencing methods in the number of "true negatives", *i.e.* any

256    gene not identified as an olfactory receptor is not an olfactory receptor under this approach. Confidence

257    intervals were calculated using 2000 bootstrap replicates of sensitivity. Sensitivity values were calculated

258    using the "pROC" v. 1.1.0 package in R v. 3.3.2 Scripts for all assemblies and *post hoc* analyses are

259    available on Dryad [XXXXX].

260    **Results**

261    *RNA-Seq and transcriptome assembly*: Extracted RNA from the MOE sample resulted in 1.09 µg at 91

262    ng/µL and an RNA integrity number (RIN) of 9.6 for *Desmodus rotundus*, enough quantity and quality

263    for library preparation. After trimming and removal of low-quality reads, the sample produced more than

264    56 million total reads, with a median insert size of 330. The average read quality for the set of pairs

265    indicated low error rate, with a mean quality score of 39.6 ±1.3 for the right and 38.8 ±1.3 for the left. As

266    expected, different assembly methods within the Oyster River Protocol resulted in different numbers of

267    genes, and ultimately 564 unique genes were identified across all assemblies. The pooled assembly

268    consisted of 255,295 sequences, in which 49% of the contigs had an open reading frame and the mean

269    contig length was 733 bp.

270         Both the transRate and BUSCO scores indicated a high-quality assembly. The optimal transRate

271    score was 0.59 and the empirical score was 0.51. Over 91% of the reads were considered "good

272    mappings" back to contigs, and only 1.3% of assembled contigs had no coverage. The lower transRate

273    score was mostly affected by the 79% of contigs considered to have low coverage, defined by a mean per-

274    base read coverage of less than 10, but this is to be expected for lowly expressed transcripts. The

275    assembly also resulted in a BUSCO score of 82.1% complete (46.5% single copy, 35.6% duplicated),

276    indicating that nearly all orthologs from the database matched to an ortholog within the assembly. Only

277    8% of the ortholog database matched to transcripts considered to be fragmented and 9.9% of the database

278    was missing.

279    *Olfactory receptor detection*: There were 424 intact ORs identified in the *D. rotundus* genome. The

280    Sanger-sequenced amplicons made available to us from previously published work consisted of 132 intact

281    olfactory receptor sequences (Hayden *et al.* 2014). From the transcriptome, 291 olfactory receptors were

282    recovered and, of these, 267 had a "good" transRate score indicating high coverage and low rates of

283    fragmentation for most of these genes. From targeted sequence capture, 424 intact olfactory receptors

284    were also recovered, though despite the exact number as those found in the genome, not all of these

285    receptors were detected in the genome and *vice versa* (see below).

286    *Olfactory receptor genome mapping*: By mapping intact olfactory receptors to the *D. rotundus* genome,

287    we assessed whether the same olfactory receptor was assembled across different sequencing and assembly

288    approaches. First, the olfactory receptor coding sequences identified from the genome were mapped back

12

289   onto the genome to obtain the location of each olfactory receptor. Of the 424 identified coding sequences,

290   only 384 sequences mapped with 100% identity to the genome, indicating a discrepancy between the

291   post-processing of the coding sequence identification (*e.g.* open reading fame editing) from the genome

292   and the actual published genome (Fig. 2). Thus, because we could only be certain of 384 olfactory

293   receptor locations, these receptor localities were used to match the receptors in the *de novo* sequencing

294   data sets. Of these 384, 5% of the receptors mapped to multiple locations. Although the genome is not

295   assembled into chromosomes, having the same scaffold index indicates receptors relatively close together.

296   The distribution of mapped reads showed most receptors were clustered by subfamily on the same

297   scaffold (Fig. 3). For the majority of subfamilies with multiple receptors, the distribution of these

298   receptors was restricted to two or three scaffolds. Class I genes in particular, which are homologous with

299   olfactory receptors across vertebrates, are mostly distributed along only two scaffolds.

300   The quality of mapping differed across sequencing approaches (Fig. 2). The Sanger-sequenced

301   amplicons had the highest proportion of failed mapped receptors compared to any other approach (Fig. 2).

302   Nearly 19% of the 132 amplicon olfactory receptors failed to map to the genome, compared to 11% of the

303   transcriptome contigs and 6% of the targeted sequence capture. Targeted sequence capture had the highest

304   proportion of uniquely mapped receptors, with nearly 84% of the receptors matching to a locality in the

305   genome.

306   To determine if different approaches recovered the same receptor, we matched the index of each

307   mapped receptor in each sequencing method to the index of the 384 genomic receptors with known

308   locations (Fig. 4). We then removed sequences that failed to map, and receptors that redundantly mapped

309   to the same position. Redundantly mapped receptors are distinct from a single receptor mapping to

310   multiple locations. Instead, receptors deemed unique in each sequencing approach data set (perhaps due to

311   a sequencing error) are considered the same receptor if they map to the same genomic location with up to

312   95% sequence identity. We report the minimum number of receptors confidently identified in the genome

313   that confidently match those in the *de novo* approaches. After filtering, 56 receptors from the genome

314   matched a Sanger-sequenced amplicon (Fig. 4; 5). In other words, a recovery rate of 42% of the Sanger-

13

315    sequenced amplicons mapped to a receptor annotated in the genome. For the transcriptome, 53% of the

316    genes were recovered and 73% of receptors were recovered for the targeted sequence capture (Fig. 4).

317    Only 20 receptors of the 384 protein-coding genomic sequences were consistently recovered by the three

318    approaches, spread across different *OR* subfamilies. The amplicon data has a clear underrepresentation of

319    certain subfamilies, particularly in the Class I receptors (Fig. 4), while the transcriptome provides a more

320    even representation survey of olfactory receptors in different subfamilies.

321        Some receptors recovered by the sequencing approaches mapped to the genome but did not map

322    to the localities of the protein-coding genes identified from the genome (*i.e.*, the receptors mapped to

323    unannotated locations in the genome). We still considered these "true" receptors since they exist in the

324    genome. Figure 5 summarizes these receptors from other sequencing approaches that mapped but were

325    not annotated in the genome. Three "true" receptors were found in the Sanger-sequenced amplicons,

326    transcriptome, and targeted bait capture and six "true" receptors were found in the Sanger-sequenced

327    amplicons and targeted bait capture but were not annotated in the genome (Fig. 5). There were five

328    receptors from Sanger-sequenced amplicons, six receptors from the transcriptome, and 28 receptors from

329    the targeted bait capture that mapped to the genome but were not recovered in any other sequence

330    approach (Fig. 5).

331        We performed sensitivity analyses to quantify the assembly of receptors within the scope of all

332    possible receptors that may be in the genome (Fig. 6). From the pool of all possible receptors determined

333    from locations in which at least one receptor mapped from one of the sequencing approaches, a total of

334    430 intact receptors were found in the genome. Sensitivity analyses represent the "true positive" results

335    for each assembly approach. The highest sensitivity was for the protein-coding genomic sequences at 0.83

336    (95% confidence intervals: 0.79, 0.87), followed by the targeted sequence capture at 0.77 (0.73, 0.81), the

337    transcriptome at 0.45 (0.40, 0.50), the Sanger-sequenced amplicon receptors were the least sensitive at

338    0.15 (0.12, 0.19) (Fig. 6).

339    **Discussion**

14

340    In this study, we compared three methods to recover high-quality sequences for multigene families in

341    non-model species lacking reference genomes. We used olfactory receptors, the largest protein-coding

342    gene family in the mammalian genome, to illustrate advantages and differences in sequencing and

343    assembly approaches. By comparing to genomic sequencing data, we showed that targeted sequence

344    capture is the most comprehensive method for recovering multigene sequences across different

345    sequencing approaches, recovering up to 72% of the receptors annotated in the genome. High-coverage

346    MOE specific transcriptomes can also recover a proportion (~48%) of olfactory receptors; however, we

347    found that no method, including high-coverage, high-quality whole-genome sequencing, resulted in a

348    complete inventory of olfactory receptors. We also found that amplicon-based approaches previously

349    used to characterize olfactory receptor repertoires produced inventories that were both the least complete

350    and the most biased in terms of olfactory subfamily representation.

351        Comparisons of the performance of sequencing and assembly for large gene families are rare,

352    though a few studies have quantified variation in success rates outside of model organisms. A previous

353    study of orchid bees, for example, identified chemosensory genes from *de novo* antennal transcriptomes

354    and compared different assemblers in their ability to recover the maximum high quality chemosensory

355    genes (Brand *et al.* 2015). This study found that Trinity (Grabherr *et al.* 2011) outperformed other

356    assembly approaches, but intensive permutations of different Trinity parameters were required to recover

357    the maximum number of unique receptors. Another study compared assembly and sequencing approaches

358    of the major histocompatibility complex (MHC) class I-like (Ib) genes in voles (Migalska *et al.* 2016).

359    This study compared *de novo* assemblies of all reads, *de novo* assemblies guided by the mouse reference

360    genome, and assemblies of reads that only mapped to MHC-Ib loci in the mouse genome. In this analysis,

361    genome-guided assemblies outperformed all other approaches, but there was extensive variation between

362    individual samples. Some individuals yielded 38 MHC-Ib gene copies out of ~130 copies, while no

363    contigs were detected in other samples. The authors also discovered high rates of chimeric sequences, and

364    incorrect bases at loci even when coverage suggested otherwise, though this may be the result of the

365    mouse reference genome diverging from the vole RNA-Seq reads. The authors found *de novo*

15

366     transcriptome data was not ideal for sequencing copies in a highly polymorphic gene family and found

367     more success in designing primers from the transcript reads and sequencing amplicons using Sanger

368     technology. Similarly, we found the probes designed from the transcriptomes recovered many more high-

369     quality olfactory receptors than the sample obtained from the transcriptome (almost three quarters vs. half

370     of the known intact receptors in the genome, Fig. 4).

371         Our study demonstrates that challenges for *de novo* sequencing and assembly of multigene

372     families are not rooted in mis-assembled reads, but rather in the recovery of the complete inventory of

373     genes within the gene family. Despite the incomplete and variable presence of receptors across methods,

374     the majority of intact receptors assembled had high coverage and high transcriptome quality scores, with

375     low rates of chimeric and failed mappings (Fig. 2). With sufficient read depth, then, transcriptome data of

376     the main olfactory epithelium can reliably assemble highly similar olfactory receptors *de novo*,

377     accounting for at least half of the intact receptors present in the genome. Targeted sequence capture

378     provides even more comprehensive recovery of the true number of receptors (Fig. 2). It is clear, however,

379     that in all approaches a significant proportion of the intact genomic receptors were missing, and in some

380     cases more than half of the receptors were absent (Fig. 4). For the transcriptome, for example, only 53%

381     of known olfactory receptors were expressed, though it is important to note that the entire receptor

382     repertoire is not expected to be expressed at all times. For example, a previous study of human olfactory

383     epithelium transcriptome discovered 88.6% of intact olfactory receptors were expressed, though these

384     data were pooled across multiple individuals (Olender *et al.* 2016). A study in mice showed 94% of the

385     olfactory receptors were expressed in mice, but these too were pooled across multiple individuals (Ibarra-

386     Soria *et al.* 2014). The study also noted that aside from a handful of receptors, most receptors were

387     expressed at very low abundance, and a receptor was considered "expressed" even if only a single

388     fragment of the known gene was present in the transcriptome. Hence, the stringent criteria we used for

389     considering an olfactory receptor expressed likely underestimates the number of receptors in the

390     transcriptomes. At the same time, the great proportion of receptors with mapped locations in the genome

391    provides greater confidence in future *de novo* transcriptome applications for species lacking a sequenced

392    genome.

393         Another objective of this study was to assess the performance of transcriptomes assembly

394    methods in characterizing the olfactory receptor repertoire. One advantage of our approach was the

395    application of the Oyster River Protocol, in which multiple assembly approaches were implemented,

396    pooled, and then filtered for quality across approaches (MacManes 2018). This consideration is

397    particularly important for large gene families with highly repetitive sequences. For example, a previous

398    analysis of the transcriptome of orchid bee olfactory receptors demonstrated that different assemblers and

399    different parameters within each assembler recovered different receptors, and ultimately the study

400    combined receptors from up to nine different assemblies (Brand *et al.* 2015). We found hundreds of

401    olfactory receptor sequences in each assembly, though only ~15% of annotated olfactory receptors had a

402    sufficiently long reading frame to be considered an intact olfactory receptor. Many olfactory receptor

403    sequences discarded from this analysis may have had coverage too low to provide a sufficiently long

404    sequence, or the transcript itself may have been degraded, especially given the tropical field conditions

405    under which the tissue was obtained (though the RIN value suggests otherwise). It is also possible that

406    many of these discarded and truncated olfactory receptors are expressed pseudogenes, as the number of

407    pseudogenized olfactory receptors is often just as diverse as the number of functional olfactory receptors

408    (Niimura 2012). Olfactory receptor pseudogenes do get transcribed (Flegel *et al.* 2013; Verbeurgt *et al.*

409    2014; Olender *et al.* 2016), and it has been recently shown that these expressed pseudogenes may actually

410    be functional (Prieto-Godino *et al.* 2016). Though outside the scope of this study, it will be worthwhile to

411    take a closer look at the patterns of pseudogene expression in these data sets.

412         Some receptors were present in some assemblies, but not in others (Fig. 4; 5). Even though we

413    described 424 receptors from the protein-coding sequences of the genome, only 384 perfectly mapped

414    back to the genome. This may be due to subsequent annotation methods of the raw genome assembly

415    during detection of protein-coding sequences. The common vampire bat genome was sequenced from two

416    individuals (Zepeda Mendoza *et al.* 2018), and thus some of the variation may have collapsed in post-

17

417   processing. Some degree of variation in copy number of certain receptors between individuals is

418   expected. Olfactory receptors are highly polymorphic in both sequence (Mainland *et al.* 2014), and the

419   number of receptors present in an individual genome (Hasin *et al.* 2008; Young *et al.* 2008). In humans,

420   an average of eleven copy number variants occur across individuals (Nozawa *et al.* 2007), and these

421   values tend to be higher in olfactory receptor pseudogenes (Nozawa *et al.* 2007; Hasin *et al.* 2008; Young

422   *et al.* 2008). Some loci may be functional in some individuals, but pseudogenized in others (Gilad and

423   Lancet 2003; Menashe *et al.* 2003; MacArthur *et al.* 2012). In our study, each sequencing approach was

424   derived from a different individual sampled from quite different localities, which may contribute to the

425   variation observed across methods. Besides this biological variation, low coverage of some receptors

426   probably caused differences among assemblies. From visual inspection, reads from transcripts found in

427   multiple assemblies were often uniquely mapped, and the corresponding transcripts had an order of

428   magnitude higher coverage of perfectly matched reads than receptors that were either chimeric or mapped

429   to multiple loci. The low coverage of the latter receptors may have led to the incorporation of wrong reads

430   into the assembly and resulted in chimeras, or the reads may have been too few to sufficiently recover the

431   contig under a different assembly condition

432         Olfactory receptors recovered from Sanger sequencing of amplicons from degenerate primers

433   performed poorly relative to other methods (Fig. 4, 5). The amplicon data exhibited the highest failure

434   rate of receptors mapped to the genome and highest rate of receptors mapping to multiple loci (Fig. 2).

435   While poor genome assembly in these repetitive regions may in part cause mapping failures, there are

436   several potential explanations for the low rates of mapping in the Sanger-sequenced amplicons, despite

437   the low error rates of Sanger sequencing. The amplicon data obtained from a previously published

438   analysis was obtained by cloning olfactory receptors amplified using two sets of degenerate primers, one

439   set for Class I genes and another for Class II genes (Hayden *et al.* 2014). The study implemented a

440   statistical "mark-recapture" analysis to determine the probability that all olfactory receptors were

441   amplified, and set the threshold for the ratio of observed olfactory receptors to expected numbers to 25%

442   (Hayden *et al.* 2010, 2014). Thus, many of the published repertoires were underrepresented. One issue

18

443    with the amplicon data is the low representation, particularly in the Class I subfamily. The low diversity

444    may be due to degenerate primer bias or clone selection bias, and this is portrayed in the clustered nature

445    of the amplicon profile in Figure 4. Targeted sequencing through primer design of multigene families has

446    been relatively successful (Hohenbrink *et al.* 2013, 2014; Larsen *et al.* 2014; Yoder *et al.* 2014; Migalska

447    *et al.* 2016), but these studies often used dozens of primer pairs. It may be that two primer sets for

448    mammalian olfactory receptors that can span over 1,000 genes is insufficient for complete representation.

449    Amplicon-based olfactory receptor analyses can be a good introductory point to documenting the

450    diversity of mammalian olfactory receptors, however, it appears caution should be used when interpreting

451    these results in the context of comparative analyses of repertoire sizes across mammalian olfactory

452    receptors.

453        Our study reveals strengths and weaknesses of different sequencing approaches for multigene

454    families in terms of completeness of the representation of each gene in the family. However, depending

455    on circumstances such as tissue availability, computing resources, and time, other factors are relevant for

456    consideration. For example, while Sanger-sequencing amplicons had the most incomplete representation,

457    Sanger sequencing has very low base calling error rates relative to high-throughput methods, does not

458    require unfeasible computing time, and uses genomic DNA that does not have to be extracted from

459    pristinely-preserved tissue as input. At the same time, while the genome is a more complete inventory, the

460    costs and resources required by Dovetail genome sequencing are beyond the capacity of many labs, and it

461    requires freshly frozen tissue, which may be unfeasible for most species. Transcriptomes are useful for

462    characterizing the expressed receptors, but also require freshly dissected epithelial tissue for RNA, which

463    may not be scalable across many species. While targeted sequence capture does require high-throughput

464    sequence data for probe design, these data can come from a subset of species or individuals. Once the

465    probes are designed, experiments only require genomic tissue and can be feasibly scalable across many

466    species or individuals. Aside from the genome, targeted sequence capture recovered a substantial

467    proportion of intact receptors and offers a promising avenue for large-scale multigene family analyses.

468    Looking ahead, as long-read sequencing becomes more tractable, this technology may also have a strong

469    influence on sequencing multigene families that are often tandem-duplicated in the same genomic region

470    (Nam *et al.* 2019).

471    Our results have several implications for studies of gene family evolution and understanding

472    olfactory receptor diversity. First, gene family evolution is frequently analyzed through birth-death

473    process, in which phylogeny-based models are applied to species and/or gene trees to understand when in

474    the evolutionary history of a group losses and duplications occurred (Hahn *et al.* 2007; Niimura and Nei

475    2007; Han *et al.* 2009; Zhao *et al.* 2015). These models rely on the assumption that all copies of the gene

476    families are known in extant species. However, although there are more than 300 intact olfactory

477    receptors in the vampire bat genome, we have shown that both the transcriptomes and the amplicon data

478    represent a severe underestimation of the total number of olfactory receptor genes with open reading

479    frames in the genome. Understanding the variance between sequencing methodologies is indispensable to

480    avoid false conclusions when studying gene family evolution. If the transcriptomic data or the amplicon

481    sequences were used in analyses with genomic olfactory receptor data from other mammals, gene losses

482    in the common vampire bat may be inferred, when the apparent loss is actually due to the failure to

483    sequence the entire intact olfactory repertoire. Therefore, we recommend using genome-based sequence

484    data or sequence capture data instead of transcriptome or amplicon data for studies of birth-death

485    evolution that require estimating the presence and absence of a receptor, as well as for any large gene

486    family.

487    While the *de novo* transcriptome sequencing of multigene families may be incomplete and

488    inappropriate for birth-death modeling, the sequence data are reliably assembled and can be used in other

489    informative ways. For example, orthologous sequences from other mammals can be identified from these

490    sequences and the strength of selection on particular receptors across species can subsequently be

491    quantified. Receptors recovered from the transcriptome can also serve as excellent starting material for

492    probe and primer design, as with our sequence capture data set. Thus, understanding the caveats and

493    strengths of different sequencing and assembly approaches, analyses molecular sequence data of

494    multigene family can be properly performed. Multigene families often compose significant proportions of

495    the genome of organisms, and often underlie mechanisms involved in immunity, metabolism, and sensory

496    perception. Thus, it is crucial to understand whether variation in multigene families is derived from

497    methodological shortcomings or whether it is biologically relevant.

498    **Acknowledgements**

508    **References**

509    Brand P., S. R. Ramírez, F. Leese, J. J. G. Quezada-Euan, R. Tollrian, *et al.*, 2015 Rapid evolution of

510        chemosensory receptor genes in a pair of sibling species of orchid bees (Apidae: Euglossini). BMC

511        Evol. Biol. 15: 176. https://doi.org/10.1186/s12862-015-0451-9

512    Brechbühl J., G. Luyet, F. Moine, I. Rodriguez, and M.-C. Broillet, 2011 Imaging pheromone sensing in a

513        mouse vomeronasal acute tissue slice preparation. J. Vis. Exp. e3311.

514        https://doi.org/doi:10.3791/3311

515    Eddy S., 2010 HMMER3: a new generation of sequence homology search software

516    Emms D. M., and S. Kelly, 2015 OrthoFinder: solving fundamental biases in whole genome comparisons

517        dramatically improves orthogroup inference accuracy. Genome Biol. 16: 1–14.

518        https://doi.org/10.1186/s13059-015-0721-2

519    Flegel C., S. Manteniotis, S. Osthold, H. Hatt, and G. Gisselmann, 2013 Expression profile of ectopic

520        olfactory receptors determined by deep sequencing. PLoS One 8.

521        https://doi.org/10.1371/journal.pone.0055368

522    Gilad Y., and D. Lancet, 2003 Population differences in the human functional olfactory repertoire. Mol.

523        Biol. Evol. 20: 307–314. https://doi.org/10.1093/molbev/msg013

524    Grabherr M. G., B. J. Haas, M. Yassour, J. Z. Levin, D. A. Thompson, *et al.*, 2011 Full-length

525        transcriptome assembly from RNA-Seq data without a reference genome. Nat. Biotechnol. 29: 644–

526        52. https://doi.org/10.1038/nbt.1883

527    Gracheva E. O., J. F. Cordero-Morales, J. A. González-Carcacía, N. T. Ingolia, C. Manno, *et al.*, 2011

528        Ganglion-specific splicing of *TRPV1* underlies infrared sensation in vampire bats. Nature 476: 88–

529        91. https://doi.org/10.1038/nature10245

530    Hahn M. W., J. P. Demuth, and S. Han, 2007 Accelerated rate of gene gain and loss in primates. Genetics

531        177: 1941–1949. https://doi.org/10.1534/genetics.107.080077

532    Han M. V, J. P. Demuth, C. L. McGrath, C. Casola, and M. W. Hahn, 2009 Adaptive evolution of young

533        duplicated genes in mammals. Genome Res. 19: 859–867. https://doi.org/10.1101/gr.085951.108

534    Hasin Y., T. Olender, M. Khen, C. Gonzaga-Jauregui, P. M. Kim, *et al.*, 2008 High-resolution copy-

535      number variation map reflects human olfactory receptor diversity and evolution. PLoS Genet. 4.

536      https://doi.org/10.1371/journal.pgen.1000249

537    Hayden S., M. Bekaert, T. A. Crider, S. Mariani, W. J. Murphy, *et al.*, 2010 Ecological adaptation

538      determines functional mammalian olfactory subgenomes. Genome Res. 20: 1–9.

539      https://doi.org/10.1101/gr.099416.109

540    Hayden S., M. Bekaert, A. Goodbla, W. J. Murphy, L. M. Dávalos, *et al.*, 2014 A cluster of olfactory

541      receptor genes linked to frugivory in bats. Mol. Biol. Evol. 31: 917–27.

542      https://doi.org/10.1093/molbev/msu043

543    Hohenbrink P., N. I. Mundy, E. Zimmermann, and U. Radespiel, 2013 First evidence for functional

544      vomeronasal 2 receptor genes in primates. Biol. Lett. 9.

545    Hohenbrink P., S. Dempewolf, E. Zimmermann, N. I. Mundy, and U. Radespiel, 2014 Functional

546      promiscuity in a mammalian chemosensory system: extensive expression of vomeronasal receptors

547      in the main olfactory epithelium of mouse lemurs. Front. Neuroanat. 8: 1–10.

548      https://doi.org/10.3389/fnana.2014.00102

549    Holding M. L., M. J. Margres, A. J. Mason, C. L. Parkinson, and D. R. Rokyta, 2018 Evaluating the

550      performance of *de novo* assembly methods for venom-gland transcriptomics. Toxins (Basel). 10: 1–

551      21. https://doi.org/10.3390/toxins10060249

552    Hong W., and H. Zhao, 2014 Vampire bats exhibit evolutionary reduction of bitter taste receptor genes

553      common to other bats. Proc. Biol. Sci. 281: 20141079. https://doi.org/10.1098/rspb.2014.1079

554    Hughes G. M., L. Gang, W. J. Murphy, D. G. Higgins, and E. C. Teeling, 2013 Using Illumina next

555      generation sequencing technologies to sequence multigene families in de *novo* species. Mol. Ecol.

556      Resour. 13: 510–21. https://doi.org/10.1111/1755-0998.12087

557    Hughes G. M., E. S. M. Boston, J. A. Finarelli, W. J. Murphy, D. G. Higgins, *et al.*, 2018 The birth and

558      death of olfactory receptor gene gamilies in mammalian niche adaptation. Mol. Biol. Evol. 35:

559      1390–1406. https://doi.org/10.1093/molbev/msy028

560    Ibarra-Soria X., M. O. Levitin, L. R. Saraiva, and D. W. Logan, 2014 The olfactory transcriptomes of

561        mice. PLoS Genet. 10: e1004593. https://doi.org/10.1371/journal.pgen.1004593

562    Johnson M. G., E. M. Gardner, Y. Liu, R. Medina, B. Goffinet, *et al.*, 2016 HybPiper: Extracting coding

563        sequence and introns for phylogenetics from high-throughput sequencing reads using target

564        enrichment. Appl. Plant Sci. 4: 1600016. https://doi.org/10.3732/apps.1600016

565    Kanageswaran N., M. Demond, M. Nagel, S. Benjamin, P. Schreiner, *et al.*, 2015 Deep sequencing of the

566        murine olfactory receptor neuron transcriptome. PLoS One 10: e0113170.

567        https://doi.org/10.1371/journal.pone.0113170

568    Larsen P. a, A. M. Heilman, and A. D. Yoder, 2014 The utility of PacBio circular consensus sequencing

569        for characterizing complex gene families in non-model organisms. BMC Genomics 15: 720.

570        https://doi.org/10.1186/1471-2164-15-720

571    Logan D. W., 2014 Do you smell what I smell? Genetic variation in olfactory perception. Biochem. Soc.

572        Trans. 42: 861–5. https://doi.org/10.1042/BST20140052

573    MacArthur D. G., S. Balasubramanian, A. Frankish, N. Huang, J. Morris, *et al.*, 2012 A systematic survey

574        of loss-of-function variants in human protein-coding genes. Science 335: 823–8.

575        https://doi.org/10.1126/science.1215040

576    MacManes M. D., 2018 The Oyster River Protocol: a multi-assembler and kmer approach for de novo

577        transcriptome assembly. PeerJ 6: e5428. https://doi.org/10.1109/ECTC.2009.5074081

578    MacRander J., M. Broe, and M. Daly, 2015 Multi-copy venom genes hidden in *de novo* transcriptome

579        assemblies, a cautionary tale with the snakelocks sea anemone *Anemonia sulcata* (Pennant, 1977).

580        Toxicon 108: 184–188. https://doi.org/10.1016/j.toxicon.2015.09.038

581    Mainland J. D., A. Keller, Y. R. Li, T. Zhou, C. Trimmer, *et al.*, 2014 The missense of smell: functional

582        variability in the human odorant receptor repertoire. Nat. Neurosci. 17: 114–20.

583        https://doi.org/10.1038/nn.3598

584    Menashe I., O. Man, D. Lancet, and Y. Gilad, 2003 Different noses for different people. Nat. Genet. 34:

585        143–4. https://doi.org/10.1038/ng1160

24

586    Migalska M., A. Sebastian, M. Konczal, P. Kotlik, and J. Radwan, 2016 De novo transcriptome assembly

587        facilitates characterisation of fast-evolving gene families, MHC class I in the bank vole (*Myodes*

588        *glareolus*). Heredity (Edinb). 118: 348–357. https://doi.org/10.1038/hdy.2016.105

589    Nam S., K. Hoff, O. K. Tørresen, A. T. Klunderud, and S. Jentoft, 2019 Long☐read sequence capture of

590        the haemoglobin gene clusters across codfish species. Mol. Ecol. Resour. 19: 245–259.

591        https://doi.org/10.1111/1755-0998.12955

592    Nei M., and A. P. Rooney, 2005 Concerted and birth-and-death evolution of multigene families. Annu.

593        Rev. Genet. 39: 121–152. https://doi.org/10.1146/annurev.genet.39.073003.112240

594    Niimura Y., and M. Nei, 2007 Extensive gains and losses of olfactory receptor genes in mammalian

595        evolution. PLoS One 2. https://doi.org/10.1371/journal.pone.0000708

596    Niimura Y., 2012 Olfactory receptor multigene family in vertebrates: from the viewpoint of evolutionary

597        genomics. Curr. Genomics 13: 103–14. https://doi.org/10.2174/138920212799860706

598    Niimura Y., A. Matsui, and K. Touhara, 2014 Extreme expansion of the olfactory receptor gene repertoire

599        in African elephants and evolutionary dynamics of orthologous gene groups in 13 placental

600        mammals. Genome Res. 24: 1485–1496. https://doi.org/10.1101/gr.169532.113

601    Nozawa M., Y. Kawahara, and M. Nei, 2007 Genomic drift and copy number variation of sensory

602        receptor genes in humans. Proc. Natl. Acad. Sci. U. S. A. 104: 20421–20426.

603        https://doi.org/10.1073/pnas.0709956104

604    Olender T., I. Keydar, J. M. Pinto, P. Tatarskyy, A. Alkelai, *et al.*, 2016 The human olfactory

605        transcriptome. BMC Genomics 1–18. https://doi.org/10.1186/s12864-016-2960-3

606    Prieto-Godino L. L., R. Rytz, B. Bargeton, L. Abuin, J. R. Arguello, *et al.*, 2016 Olfactory receptor

607        pseudo-pseudogenes. Nature 539: 93–97. https://doi.org/10.1038/nature19824

608    Quinlan A. R., 2014 BEDTools: the Swiss☐army tool for genome feature analysis. Curr. Protoc.

609        Bioinforma. 11–12.

610    Shi W., J. Peifeng, and F. Zhao, 2017 The combination of direct and paired link graphs can boost

611        repetitive genome assembly. Nucleic Acids Res. 45: 1–15. https://doi.org/10.1093/nar/gkw1002

25

612    Shiao M. S., A. Y. F. Chang, B. Y. Liao, Y. H. Ching, M. Y. J. Lu, *et al.*, 2012 Transcriptomes of mouse

613        olfactory epithelium reveal sexual differences in odorant detection. Genome Biol. Evol. 4: 703–712.

614        https://doi.org/10.1093/gbe/evs039

615    Sims D., I. Sudbery, N. E. Ilott, A. Heger, and C. P. Ponting, 2014 Sequencing depth and coverage: key

616        considerations in genomic analyses. Nat. Rev. Genet. 15: 121–32. https://doi.org/10.1038/nrg3642

617    Smith-Unna R., C. Boursnell, R. Patro, J. M. Hibberd, and S. Kelly, 2016 TransRate: Reference-free

618        quality assessment of *de novo* transcriptome assemblies. Genome Res. 26: 1134–1144.

619        https://doi.org/10.1101/gr.196469.115

620    Treangen T. J., and S. L. Salzberg, 2012 Repetitive DNA and next-generation sequencing: computational

621        challenges and solutions. Nat. Rev. Genet. 46: 36–46. https://doi.org/10.1038/nrg3164

622    Verbeurgt C., F. Wilkin, M. Tarabichi, F. Gregoire, J. E. Dumont, *et al.*, 2014 Profiling of olfactory

623        receptor gene expression in whole human olfactory mucosa. PLoS One 9: 21–26.

624        https://doi.org/10.1371/journal.pone.0096333

625    Vijay N., J. W. Poelstra, A. Künstner, and J. B. W. Wolf, 2013 Challenges and strategies in transcriptome

626        assembly and differential gene expression quantification. A comprehensive *in silico* assessment of

627        RNA-seq experiments. Mol. Ecol. 22: 620–34. https://doi.org/10.1111/mec.12014

628    Waterhouse R. M., M. Seppey, F. A. Simão, M. Manni, P. Ioannidis, *et al.*, 2017 BUSCO applications

629        from quality assessments to gene prediction and phylogenomics. Mol. Biol. Evol. 35: 543–548.

630    Wu T. D., and C. K. Watanabe, 2005 GMAP: A genomic mapping and alignment program for mRNA and

631        EST sequences. Bioinformatics 21: 1859–1875. https://doi.org/10.1093/bioinformatics/bti310

632    Yoder A. D., L. M. Chan, M. dos Reis, P. A. Larsen, C. R. Campbell, *et al.*, 2014 Molecular evolutionary

633        characterization of a V1R subfamily unique to strepsirrhine primates. Genome Biol. Evol. 6: 213–

634        227. https://doi.org/10.1093/gbe/evu006

635    Yohe L. R., L. Liu, L. M. Dávalos, and D. A. Liberles, 2019 Protocols for the Molecular Evolutionary

636        Analysis of Membrane Protein Gene Duplicates, pp. 49–62 in *Computational Methods in Protein*

637        *Evolution*, edited by Sikosek T. Springer New York, New York, NY.

638    Yoon S., Z. Xuan, V. Makarov, K. Ye, and J. Sebat, 2009 Sensitive and accurate detection of copy

639        number variants using read depth of coverage. Genome Res. 19: 1586–1592.

640        https://doi.org/10.1101/gr.092981.109

641    Young J. M., R. M. Endicott, S. S. Parghi, M. Walker, J. M. Kidd, *et al.*, 2008 Extensive copy-number

642        variation of the human olfactory receptor gene family. Am. J. Hum. Genet. 83: 228–242.

643        https://doi.org/10.1016/j.ajhg.2008.07.005.

644    Zepeda Mendoza M., Z. Xiong, M. Escalera-Zamudio, A. Kathrine Runge, J. Thézé, *et al.*, 2018

645        Hologenomic adaptations underlying the evolution of sanguivory in the common vampire bat. Nat.

646        Ecol. Evol. 2: 659–668. https://doi.org/10.1038/s41559-018-0476-8

647    Zhao J., A. I. Teufel, D. A. Liberles, and L. Liu, 2015 A generalized birth and death process for modeling

648        the fates of gene duplication. BMC Evol. Biol. 15: 275. https://doi.org/10.1186/s12862-015-0539-2

649

**Figure Legends**

**Figure 1**. Sequencing and assembly approaches compared in this study. Sequencing approaches in red were mapped to the olfactory receptors found in the genome (black), and the proportion of receptors recovered from each approach were compared. The Sanger-sequenced amplicons were derived from published vampire bat olfactory receptors sequenced and amplified from degenerate primers (Hayden *et al.* 2014). Targeted sequence capture genes were sequenced using probes from *de novo* transcriptome assemblies. We show the expected length of olfactory receptor sequence recovered from each method and outline some pros (+) and cons (-) of each approach.

**Figure 2**. Number of receptors mapped using GMAP v. 2017-01-14 (Wu and Watanabe 2005) to the vampire bat genome (Zepeda Mendoza *et al.* 2018) for each sequencing and assembly approach, showing receptors mapped to unique positions in the genome, receptors mapped to more than one position, and receptors that failed to map (less than 95% sequence identity).

**Figure 3**. Number of olfactory receptors found by scaffold of the vampire bat genome, color-coded by olfactory receptor subfamily. Only scaffold indices that contained one or more olfactory receptors are shown.

**Figure 4**. Tile plot of olfactory receptors recovered from each sequencing approach relative to receptors present in the genome, grouped by olfactory receptor subfamily. Each row indicates a single olfactory gene identified in the genome. Empty boxes denote no olfactory receptor recovered in that sequencing or assembly approach mapped to the same location as the olfactory receptor from the genome.

**Figure 5**. Venn diagram of the number of intact receptors recovered from each method that were also recovered in an alternative sequencing approach.

**Figure 6**. Quantification of the sensitivity for each sequencing approach of the recovery of the number of potential intact olfactory receptors. Any receptor from any method that mapped to the genome was considered a "true positive". A receptor that was present in the genome, but not found in another method was considered a "false negative".

**SANGER**
with degenerate primers

**~750bp**

+lowest base calling error rates
-primer bias
-cloning bias

**GENOME**

**~900bp**

+ high coverage with Dovetail
+ annotated open reading frames
- composite of two individuals

**SEQUENCE CAPTURE**

**~900bp**

+discernible, yet flexible specificity
-dependent on probe design

**TRANSCRIPTOME**

**~900bp**

+ subset of genome
+ independent of probe design
- expression bias