

Validating a laboratory pipeline for accurate reconstruction of metabarcode data

Yige Sun

MRes. Computational Methods in Ecology and Evolution

Supervised by

Prof. Alfried Vogler

Faculty of Natural Sciences, Department of Life Sciences(Silwood Park), ICL

Department of Life Sciences, Natural History Museum, London

a.vogler@imperial.ac.uk

2019-2020

1 **Keywords**

2 community barcoding, metabarcoding analysis, metagenomics, high-throughput sequencing, tropical
3 beetles, species identification.

4 **Introduction**

5 Metabarcoding is a robust technique in characterising the complex community compositions and
6 breaking through the barrier in traditional taxonomic methods. It is a highly suitable method for the
7 study of large-scale species richness and complex community composition surveys (Lebuhn et al.,
8 2013). It also is a widely used technique for the study of arthropods with High-throughput Sequencing
9 (Ji et al., 2013).

10 In insect barcoding analysis, cytochrome oxidase C subunit I (cox1) is commonly used as a barcode
11 marker (Hebert et al., 2003). Most existing studies are clustering the sequence variants into a Op-
12 erational Taxonomic Units (OTUs) that roughly correspond to the species in the Linnaean taxonomy,
13 which can accommodate PCR and sequencing errors in the individual reads, but this removes genetic
14 variation.

15 Newer approaches attempt to retain all true genetic variants, the so-called amplicon sequence vari-
16 ants (ASVs) after removing the sequencing errors. However, this leads to a further challenge from
17 amplification of nuclear mitochondrial DNA segments (NUMTs), which are pseudogenes derived from
18 the mitochondrial copies that persist in the nuclear genome to increase the apparent diversity of geno-
19 types. In addition, internalised parasites or gut contents of insects are co-amplified and represent a
20 useful source of ecological information but unhelpful for taxonomic study.

21 A previous study of UK bee survey has built up a pipeline for the taxonomic assignment which showed
22 a high congruence with morphological classification (Creedy et al., 2019). It removed the pseudo-
23 genes based on reads abundance and phylogenetic relatedness to references sets. Moreover, it re-
24 veals that there still have rooms for bioinformatically improvement in true ASVs retention and generate
25 a more accurate taxonomic assignment.

26 According to such standardised analysis have not been applied on tropical beetles survey. In this
27 project, I will analyse several metabarcode datasets of tropical beetles to improve the construction
28 of clean ASV information and detection of spurious NUMTs copies and reveal the biodiversity in the
29 regional tropical beetles community.

30 **Methods**

31 Various metabarcoding library that include up to 7,000 tropical beetles sequences generated with
32 Illumina MiSeq v.3 will be subjected to data handling and filtering by quality procedures, including
33 adaptor removal by cutadapt (Martin, 2011) and quality filtering with FASTQC (Andrews et al., 2012),
34 followed by denoising using UNOISE (Edgar et al., 2011) under various parameter settings.

35 The resulting reads will be further filtered by relatedness with reference ASVs from mitochondrial
36 whole-genome sequencing. These procedures would be carried out on individual specimens and
37 assigned to a correct haplotype. Then, specimens would be classified by relatedness in to batches
38 of 50 and set up reads abundance threshold values for true ASVs retention based on relatedness

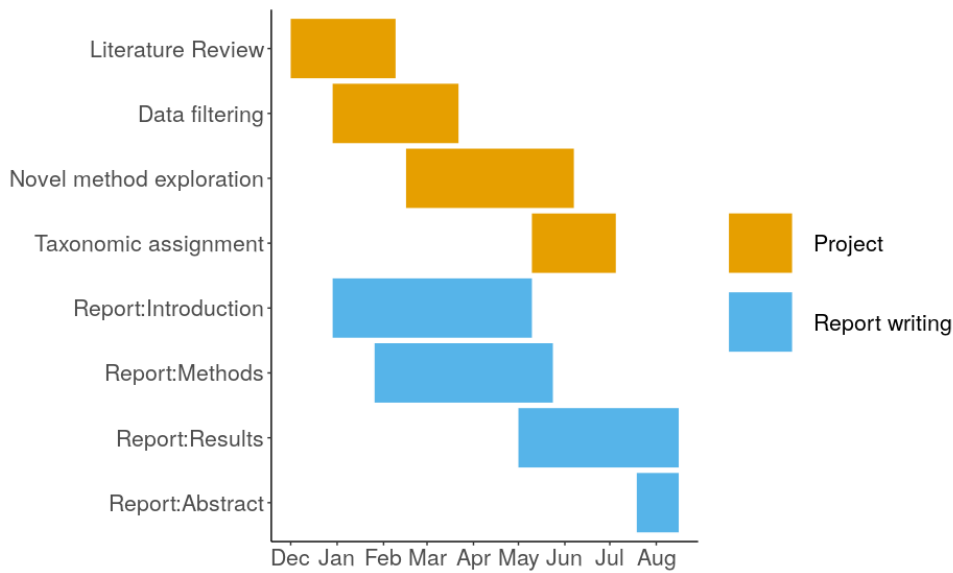
39 of representative reference ASVs. It is a most plausible method for ASV retention that founded in
40 previous studies, but there is a trade-off between insufficient removal of NUMTs and poor retention
41 of true ASVs copies, in particular if these correspond to rare species or low biomass.
42 Therefore, apart from applying the developed ASV retention method on tropical beetles dataset,
43 I am also aiming to explore a improved method to recognise the abundance pattern of true ASV
44 copies, NUMTs and other non-targeted reads(e.g. parasite or gut contents) based on relatedness to
45 reference ASVs. The reads that are phylogenetically distant from reference ASV would be considierd
46 as non-target reads then directly get omitted. The reads site on closely related branches would be
47 assign abundance threshold values for true ASV copies retention and NUMTs removals. The new
48 method would be tested with simulating a mock community or more complex community dataset.
49 Finally, the retained ASV copies will be subjected to BLASTn search against the NCBI database to
50 concludes the community composition of regional tropical beetles.

51 **Anticipated Outcomes**

52 To apply an automated and valid metabarcoding analytical pipeline for large bulk tropical beetles
53 biodiversity study. Draw an conclusion on diversity of a regional tropical beetles community. Also
54 enhance the accuracy and efficiency of ASVs retention, ultimately provides a more accurate method
55 of taxonomy assignment on complex community.

56 **Project Feasibility**

Work timeline represented in the Gantt Chart below:



57

58 **Budgets**

59

Transportation to NHM	£16.55/travel	36 weeks
Total		£595.8

60 References

- 61 Andrews, S., Krueger, F., Segonds-Pichon, A., Biggins, L., Krueger, C., Wingett, S. (2012). FastQC:
62 A quality control tool for high throughput sequence data. Babraham, UK: Babraham Institute.
- 63 Creedy T.J, Norman H, Tang C.Q, Chin K.Q, Andujar C, Arribas P, O'Connor R.S, Carvell C, Not-
64 ton D.G, Volger A.P,(2019). A validated workflow for rapid taxonomic assignment and monitoring
65 of a national fauna of bees (Apiformes) using high throughput DNA barcoding. *Molecular Ecology*
66 *Resource*,00, 1-14.
- 67 Edgar, R. C., Haas, B. J., Clemente, J. C., Quince, C., Knight, R. (2011). UCHIME improves sensi-
68 tivity and speed of chimera detection.*Bioinformatics*, 27, 2194-2200.
- 69 Hebert, P. D. N., Cywinska, A., Ball, S. L., DeWaard, J. R. (2003). Biological identifications through
70 DNA barcodes. *Proceedings of the Royal Society of London. Series B: Biological Sciences*, 270(1512),
71 313-321.
- 72 Ji, Y., Ashton, L., Pedley, S. M., Edwards, D. P., Tang, Y., Nakamura, A., Yu, D. W. (2013). Reliable,
73 verifiable and efficient monitoring of biodiversity via metabarcoding. *Ecology Letters*, 16, 1245-1257.
- 74 Lebuhn, G., Droege, S., Connor, E. F., Gemmill-Herren, B., Potts, S. G., Minckley, R. L., Parker,
75 F. (2013). Detecting insect pollinator declines on regional and global scales. *Conservation Biology*,
76 27, 113-120.
- 77 Martin, M. (2011). Cutadapt removes adapter sequences from highthroughput sequencing reads.
78 *Embnet*, 17, 10-12.

79 "I have seen and approved the proposal and the budget"

80 Prof. Alfried Vogler

81 Dec 11 2019