

1 **Conservation genomics reveals possible illegal trade routes and**
2 **admixture across pangolin lineages in Southeast Asia**

3 Helen C. Nash ^{a*}, Wirdateti ^b, Gabriel W. Low ^a, Siew Woh Choo ^{c, d}, Ju Lian Chong ^e,
4 Gono Semiadi ^b, Ranjeev Hari ^{c, f}, Muhammad Hafiz Sulaiman ^e, Samuel T. Turvey ^g,
5 Theodore A. Evans ^{a, h}, Frank E. Rheindt ^a

6 ^a Department of Biological Sciences, National University of Singapore, 14 Science
7 Drive 4, 117543, Singapore

8 ^b Research Center for Biology, Indonesian Institute of Sciences (LIPI), Jl Raya
9 Jakarta-Bogor Km 45, Cibinong 16911, Indonesia

10 ^c Genome Informatics Research Laboratory, High Impact Research (HIR) Building,
11 University of Malaya, 50603 Kuala Lumpur, Malaysia

12 ^d Department of Biological Sciences, Science Building B, Xi'an Jiaotong-Liverpool
13 University, 111 Ren'ai Road, Suzhou Dushu Lake Science and Education Innovation
14 District, Suzhou Industrial Park, Suzhou, P. R. China, 215123

15 ^e Department of Biological Sciences, Faculty of Science & Technology, Universiti
16 Malaysia Terengganu, 21030 Kuala Terengganu, Terengganu, Malaysia

17 ^f Centre for Bioinformatics, School of Data Sciences, Perdana University, 43400
18 Serdang, Selangor, Malaysia

19 ^g Institute of Zoology, Zoological Society of London, Regent's Park, London NW1
20 4RY, UK

21 ^h School of Biological Sciences, The University of Western Australia (M092), 35
22 Stirling Highway, Crawley, WA 6009, Australia

23

24 *Corresponding author (helencatherinenash@yahoo.co.uk, +65 84313054)

25 **Abstract**

26 The use of genome-wide genetic markers is an emerging approach for informing
27 evidence-based management decisions for highly threatened species. Pangolins are
28 the most heavily trafficked mammals across illegal wildlife trade globally, but
29 Critically Endangered Sunda pangolins (*Manis javanica*) have not been widely
30 studied in insular Southeast Asia. We used > 12,000 single nucleotide polymorphic
31 markers (SNPs) to assign pangolin seizures from illegal trade of unknown origin to
32 possible geographic sources via genetic clustering with pangolins of known origin.
33 Our SNPs reveal three previously unrecognized genetic lineages of Sunda pangolins,
34 possibly from Borneo, Java and Singapore/Sumatra. The seizure assignments
35 suggest the majority of pangolins were traded from Borneo to Java. Using
36 mitochondrial markers did not provide the same resolution of pangolin lineages, and
37 to explore if admixture might explain these differences, we applied sophisticated
38 tests of introgression using > 2,000 SNPs to investigate secondary gene flow
39 between each of the three Sunda pangolin lineages. It is possible the admixture
40 which we discovered is due to human-mediated movements of pangolins. Our
41 findings impact a range of conservation actions, including tracing patterns of trade,
42 repatriation of rescue animals, and conservation breeding. In order to conserve
43 genetic diversity, we suggest that, pending further research, each pangolin lineage
44 should as a precaution be protected and managed as an evolutionarily distinct
45 conservation unit.

46

47 **Keywords:** SNPs; mitochondrial markers; gene flow; illegal wildlife trade; population
48 assignment; conservation breeding

49

50 **1. Introduction**

51 The importance of using evidence-based conservation to inform effective
52 management decisions is increasingly recognized by conservation researchers and
53 practitioners (Sutherland et al. 2004; Segan et al. 2010; Nash et al. 2016). The use
54 of genetics to inform evidence-based management decisions for highly threatened
55 species can improve conservation outcomes (Allendorf et al. 2010; Corlett 2016;

56 Pierson et al. 2016). For example genetic tools have been used to define species
57 delimitations accurately (e.g. in passerine birds, Lohman et al. 2010; in crocodiles;
58 Shirley et al. 2014), trace illegal wildlife trade (e.g. in sharks, Clarke et al. 2006; in
59 elephant ivory, Wasser et al. 2008), assess population viability (e.g. for the Komodo
60 dragon Ciofi et al. 1999; and sturgeon Schueller and Hayes 2011), and to inform
61 conservation breeding, which includes captive breeding and genetic rescue of wild
62 populations (e.g. Florida panther, Johnson et al. 2010; Burmese roofed turtle, Cilingir
63 et al. 2017).

64 The conservation of pangolins would benefit from the application of genetic tools.
65 There are eight pangolin species which are insectivorous, scaly mammals, native to
66 Asia and Africa. Populations of all species are in decline due to habitat clearing and
67 high levels of poaching driven by demand for traditional medicines and meat
68 (Challender et al. 2014a). Consequently, pangolins are considered to be the most
69 heavily trafficked mammals across illegal wildlife trade globally (Challender et al.
70 2015). The Sunda pangolin (*Manis javanica*) is distributed across several countries
71 in Southeast Asia. 'Sunda' refers to Sundaland, which is a biogeographical region
72 including the Malay Peninsula and Indonesian archipelago. Sunda pangolins are
73 listed as Critically Endangered in the IUCN Red List, especially populations in
74 Indonesia (Challender et al. 2014b). The need for increased genetic research is
75 highlighted as a priority activity for pangolins in the global conservation action plan of
76 the IUCN Species Survival Commission's Pangolin Specialist Group (Challender et
77 al. 2014a).

78 A broad variety of genetic methods and different genetic markers can help to inform
79 conservation actions for pangolins (Alacs et al. 2010; Hassanin et al. 2015; Tan et al.
80 2016). Among these methods, genome-wide markers such as single nucleotide
81 polymorphisms (SNPs) are a powerful tool to provide detailed information about
82 population structure, sometimes with greater resolution than other markers including
83 microsatellites (Malenfant et al. 2014). Moreover, next-generation sequencing
84 methods such as double digest restriction site-associated DNA sequencing
85 (ddRADseq) have been used successfully to investigate population structure of
86 mammals (Knowles et al. 2016).

87 There are existing genetic studies for the Sunda pangolin. A whole-genome of the
88 Sunda pangolin was recently sequenced and published (Choo et al. 2016), providing
89 the genomic infrastructure for further genetic research to inform robust conservation
90 actions and management plans. The population genetic structure of Sunda pangolins
91 across Indonesia has previously been investigated using the mitochondrial (mtDNA)
92 control region (Wirdateti & Semiadi, 2017), with this study suggesting that there
93 might be more than one pangolin lineage across Indonesia. However, mtDNA is only
94 a single marker, prone to biases such as selective sweeps and introgression (Ballard
95 & Whitlock, 2004). Wider-scale techniques such as ddRADseq could instead be
96 used to look more rigorously at Sunda pangolin populations with genome-wide
97 genetic markers.

98 The illegal trade of wildlife generates billions of USD per year, and a worryingly high
99 proportion of this trade includes pangolins (UNODC, 2016). Unsustainably large
100 seizures of pangolins have occurred across insular Southeast Asia, for example, >5
101 tonnes in Medan, Indonesia (WCS News Releases, 2015), and tonnes of scales and
102 meat destined for countries such as China (Cheng et al. 2017). A paucity of wild
103 samples of known geographic origin has hindered genetic assignments to investigate
104 the sources of illegally traded pangolins (Zhang et al. 2015), and further widespread
105 sampling of pangolins of known geographic origin with full chains of custody and
106 expert identification is required. Meanwhile, the origin of illegally traded pangolins is
107 an urgent issue which needs to be addressed.

108 In this study, we used ddRADseq to generate genome-wide genetic markers to
109 assign pangolins of unknown origin from seizures of illegal trade to genetic clusters
110 with wild samples of known origin in Sundaland. This is the first application of next-
111 generation sequencing methods to the population assignment of pangolin seizures.
112 We conducted analyses for > 12,000 SNPs and two mitochondrial genes. Our results
113 are consequential to a range of conservation actions for pangolins, such as tracing
114 illegal wildlife trade, repatriation of rescued pangolins, and conservation breeding.

115 **2. Methods**

116 2.1. Sample Collection

117 We aimed to collect wild Sunda pangolin samples of known origin from across
118 insular Southeast Asia (Figure 1), and with the Indonesian Institute of Sciences (LIPI)
119 we also acquired pangolin samples from seizures of illegal wildlife trade of unknown
120 origin from across Indonesia (Figure 1). In total, we obtained 97 Sunda pangolin
121 samples between July 2008 and January 2016. The total included 89 Indonesian
122 tissue samples from the muscle of dead pangolins, eight of which had a known origin
123 from wild populations in Java (Jember), Sumatra (Lampung) and Kalimantan
124 (Pangkalanbun, Indonesian Borneo), plus seven blood samples from wild
125 Singaporean pangolins, and one tissue sample from a dead pangolin in Sarawak
126 (Malaysian Borneo). Pangolin specialists were present at each site for the collection
127 of wild samples and were able to confirm via morphological features, such as hind
128 scale counts (Gaubert 2011), that the wild samples were all Sunda pangolins. LIPI
129 also maintain official national records and documentation of their samples.
130 Veterinarians at Wildlife Reserves Singapore (WRS) collected the blood samples
131 from anaesthetized Singaporean pangolins. Similarly, veterinarians at Kadoorie
132 Farm and Botanic Garden in Hong Kong, China, collected a blood sample from one
133 Chinese pangolin (*Manis pentadactyla*) in July 2014, which was used as an outgroup.

134 2.2. DNA extraction and ddRADseq library preparation

135 We extracted DNA at LIPI, the National University of Singapore (NUS), the Universiti
136 Malaysia Terengganu, and Kadoorie Farm and Botanic Garden, using Qiagen
137 DNeasy Blood & Tissue Kits. Chinese and Sunda pangolin DNA were exported to
138 Singapore with appropriate Convention on International Trade in Endangered
139 Species of Wild Fauna and Flora (CITES) permits in 2014 and 2015/16 respectively.
140 To measure DNA yields we used fluorometric quantitation of double-stranded DNA
141 via Qubit 2.0™. For tissue samples with low yields, we re-extracted DNA using
142 phenol chloroform.

143 To obtain genome-wide markers (SNPs) with next-generation sequencing, we
144 modified a protocol for ddRADseq (Peterson et al. 2012) following Tay et al. (2016)
145 and used the restriction enzymes EcoRI-HF and MspI because they worked well for
146 other taxa (Garg et al. 2016; Ng et al. 2017). During optimisation of the RADseq
147 protocol, we selected a Sera-Mag® bead ratio that produced DNA fragments within a
148 range of 250–650 base pairs (bp) (Appendix A). We used additional control samples

149 with molecular grade water instead of DNA throughout all procedures to confirm
150 there had been no contamination. Singapore Centre on Environmental Life Sciences
151 Engineering checked the quality of each DNA library and sequenced each library in
152 two lanes of one flowcell of an Illumina HiSeq 2500 Rapid Sequencing Run to
153 produce 2 x 150 bp paired end reads. We spiked both lanes with 5% PhiX to
154 increase the quantity of data obtained.

155 2.3. Bioinformatic Analysis

156 2.3.1. Identification of SNPs

157 To call SNPs across our reads, we first checked the quality of the 150 bp paired end
158 reads with FastQC version 0.11.5 (Andrews 2016). We used a Phred Score of 20 as
159 our quality threshold, which meant we had to truncate reads to 135 bp for further
160 analysis (Appendix B). We then demultiplexed the reads using process_radtags in
161 STACKS version 1.35 (Catchen et al. 2013), which grouped the uniquely labelled
162 reads of each sample (Appendix C). We indexed the nuclear genome sequence of a
163 Sunda pangolin (Choo et al. 2016) using Burrows-Wheeler Alignment Tool version
164 0.7.1. (Li et al. 2013), and we aligned our reads to it using bwa_memscript (Li et al.
165 2013) (Appendix D). We used Samtools version 1.4 (Li et al. 2009) to convert SAM
166 to BAM files, and to sort the BAM files (Appendix E). Our ref_map.pl pipeline in
167 STACKS included pstacks, cstacks and sstacks, with a minimum stacks depth of 5,
168 to call SNPs in each sample, and match loci across populations according to
169 alignment positions (Catchen et al. 2013) (Appendix F). In addition, we ran
170 population analysis in ref_map.pl, with a minimum of 90% of individuals in a
171 population required to process a locus for that population. For the population labels
172 we assigned Singaporean versus non-Singaporean samples. Later, we checked if a
173 default model of no population substructure changed the overall results, which it did
174 not (Appendix G).

175 For further quality control of our SNP calling, we removed SNPs with 10% or more
176 missing data, and individuals with more than 15% of loci missing, using PLINK
177 version 1.9 (Purcell et al. 2007). We also tested higher missing data cut-offs to see if
178 we could retain additional samples for analysis but it was not feasible. We pruned
179 SNPs that were correlated (Appendix H), which applied a sliding window of 25 bp
180 and removed correlated SNPs of R-squared ≥ 0.9 using a normal distribution curve

181 across each window. We also used PGDSpider version 2.1.0 (Lischer & Excoffier
182 2012) and BayeScan version 2.1 (Foll & Gaggiotti 2008) to further confirm that no
183 loci were under selection (Appendix I).

184 2.3.2. Population genomic analysis

185 2.3.2.1. Principal Component Analysis and Fst estimation

186 As a preliminary analysis to explore how many genetic clusters there were across
187 our samples, we applied Principal Component Analysis (PCA) to the remaining 83
188 Sunda pangolin samples to compare principal components of 12,150 SNPs using the
189 SNPRelate package in R version 3.2 (R Core Team, 2016) (Figure 2a, and Appendix
190 J). We also used SNPRelate to investigate pairwise Fst between clusters to get a
191 sense of the extent of genetic differentiation (Appendix J), which applied the method
192 of Weir & Cockerham (1984) to estimate Fst.

193 2.3.2.2. Bayesian clustering approaches and Network Analyses

194 Informed by the PCA results, we next applied a Bayesian clustering approach, using
195 STRUCTURE and CLUMPP, which required an a priori understanding of the
196 potential number of clusters. For STRUCTURE (Pritchard et al. 2000), we tested K =
197 1 to K = 7 to investigate whether there might be 1 to 7 genetic clusters across our 83
198 samples (Appendix K). We then used CLUMPP version 1.1.2 to determine the
199 optimal alignment of clusters (Jakobsson & Rosenberg 2007) (Appendix L).

200 We were aware that STRUCTURE sometimes generates erroneous results due to
201 uneven sample sizes between subpopulations (Puechmaille 2016), so we double-
202 checked our genetic cluster results by using Network Analyses in NetView version
203 1.0, available in RStudio version 0.99.903 (RStudio Team 2015). The Network
204 Analyses were based on a genetic distance matrix which we made from PLINK
205 version 1.9 for our 83 samples using the 12,150 SNPs. We investigated three
206 network algorithms: Fast-greedy, Infomap and Walktrap (Appendix M). Network
207 Analyses provide a range of results which are all valid clustering arrangements
208 (Appendix M). We visualized the clustering arrangement with the highest genetic
209 distances as a mutual k-nearest neighbour graph in RStudio (Figure 2b).

210 During the course of our research, new methods to understand population genetic
211 structure became available, so in addition we applied fineRADstructure package v0.2
212 (Malinsky et al. 2016) to quantify the ancestry sources in each population (Figure 2d).
213 FineRADstructure utilizes a fineSTRUCTURE MCMC clustering algorithm (Lawson
214 et al. 2012) to infer a co-ancestry matrix, which is a summary of nearest neighbour
215 haplotype relationships across the dataset (Malinsky et al. 2016). We used our
216 haplotypes.tsv file generated from the above populations analysis in ref_map.pl, and
217 we converted it to a FineRADstructure input using Python scripting contributed by
218 Emiliano Trucchi (Appendix S, Malinsky et al. 2016). The conversion script was run
219 in PyCharm 2017.3 (Professional Edition), using 1 as the maximum number of SNPs
220 per locus, and 50% as the maximum percentage of missing loci to be included in the
221 PCA. The amount of missing data to allow was based on the missing data plot from
222 fineRADstructure (Appendix S) in order to exclude only samples with high missing
223 data. The clustered fineRADstructure co-ancestry matrix for 80 samples was
224 visualized in RStudio (Figure 2c, and Appendix S).

225 2.3.3. Phylogenetic Analyses

226 2.3.3.1. Phylogenomics using SNPs

227 We felt the consistent emerging trend across our genetic clustering results warranted
228 some further investigation to try and better understand the evolutionary trajectory of
229 these lineages. In order to generate nucleotide sequences containing SNPs for the
230 construction of maximum likelihood phylogenies in RAxML version 8.2.9 (Stamatakis,
231 2014), we aligned demultiplexed 135bp sequence reads in pyRAD version 3.0.64
232 (Eaton, 2014). This included the Chinese pangolin as an outgroup. SNPs were called
233 using the ddrad option with the following parameters: clustering threshold of 95%,
234 minimum cluster coverage of five, maximum of six low-quality sites per locus,
235 minimum of 79 samples present in a final locus, and maximum of three individuals
236 with a shared heterozygous site per locus (Appendix N). A total of 2,365 SNPs was
237 generated and inputted in RAxML version 8.2.9 with the following parameters:
238 GTRGAMMA option provided, 1000 rapid bootstraps inferences, and a final
239 maximum likelihood search (Figure 3a, and Appendix O). We tested other runs with
240 increasing and decreasing amounts of missing data, but the bootstrap support

241 values did not improve, so we only present the concatenation method with 2,365
242 SNPs (Figure 3a).

243 2.3.3.2. Phylogenetics using mitochondrial DNA

244 It was useful to compare the genetic clustering results of our genome-wide markers
245 with mitochondrial DNA (mtDNA) markers which are commonly used to assign
246 species in illegal wildlife trade, so we sequenced two mtDNA coding genes,
247 cytochrome b (*Cytb*) and cytochrome oxidase c subunit 1 (*CO1*). We didn't have
248 sufficient DNA remaining from every pangolin sample, but we were able to include
249 the majority, 59 pangolins. We conducted Sanger sequencing at NUS and LIPI using
250 primer sequences provided by LIPI (Appendix Q). Sequences were aligned and then
251 *Cytb* and *CO1* sequences were manually concatenated in MEGA 7.0 (Kumar et al.
252 2016). The complete concatenation sequence consisted of 1575 base pairs (*Cytb* =
253 787 bp, *CO1* = 788 bp), and the sequences are available on GenBank (respective
254 accession numbers MG825495-MG825551 and MG825552-MG825610). We added
255 the mitogenome of a Chinese pangolin (GenBank KT445978.1) and one additional
256 Sunda pangolin (NC_026781.1), then used MEGA 7.0 (Kumar et al. 2016) to
257 construct a phylogenetic tree with maximum likelihood (ML), with a General Time
258 Reversible Model, partial deletion and 1000 bootstraps (Figure 4a). In order to
259 facilitate comparison of results, the first letter of the mtDNA tree labels reflect the
260 genetic clustering result from the SNPs, e.g. J = Java, and the seizure location is
261 also given at the end of each label. We used asterisks beside the labels to indicate
262 when the seizure location and genetic cluster result from the SNPs were the same
263 geographic area (Figure 4a). We also generated a haplotype network from the
264 concatenated mtDNA sequences using the Median-Joining method in PopART
265 (Leigh & Bryant 2015) (Figure 4b). The colour of each sample label represents the
266 SNP cluster results to further aid comparison between the mtDNA and SNP results
267 (Figure 4b).

268 2.3.4. Tests for Introgression

269 The contrasts between our genetic clustering results from the SNPs (Figure 2)
270 versus the topology of the mtDNA tree/haplotype network (Figure 4) raised our
271 suspicion that potential introgression might explain these differences. It was not
272 possible to test every sample for introgression because the computational run time is

273 prohibitive. Instead we selected a few anomalous results from the mtDNA tree, and
274 we applied ABBA BABA tests to investigate secondary gene flow in those samples
275 from our wild samples of known origin (Zinenko et al. 2016). The ABBA BABA tests
276 required nucleotide sequences containing SNPs so we used our PHYLIP files from
277 pyRAD including 2,365 SNPs (Appendix R). Our first ABBA BABA test used a wild
278 sample from Borneo, MZBR 1163, as group A; an anomalous sample, MZBR 1040,
279 which falls within the Bornean cluster of SNPs, but has suspected introgression from
280 Java based on the mtDNA result, as group B; and a wild sample from Java, MZBR
281 1184, as group C; the outgroup was our Chinese pangolin. This tested MZBR 1040
282 for suspected introgression from the Javan lineage. Our second and third tests used
283 a wild sample from Java, MZBR 1184, as group A; and an anomalous sample,
284 MZBR 0270, which falls within the Javan cluster of SNPs, but has suspected
285 introgression from both Singapore/Sumatra and Borneo based on the mtDNA result,
286 as group B; group C was initially a wild sample from Singapore/Sumatra, rescue 1,
287 which we then switched to MZBR 1163 a wild sample from Borneo in another run;
288 the outgroup was always our Chinese pangolin. This initially tested MZBR 0270 for
289 suspected introgression from the Singapore/Sumatra lineage, and next tested the
290 same sample for suspected introgression from the Bornean lineage. The ABBA
291 BABA results were summarised in a simple cartoon figure (Figure 3b). These ABBA
292 BABA tests provided examples of introgression.

293 2.3.5. Tracing illegal trade

294 The genetic clustering results from each genome-wide population genomic analysis,
295 including PCA, STRUCTURE, NetView and fineRADstructure, all gave compatible
296 results. Hence, any of those methods could be used to assign the seized pangolin
297 samples of unknown origin to the wild samples of known origin (Figure 2 and
298 Appendix L). The assignment results were the same for every method, except that
299 fineRADstructure only used 80/83 samples due to differences in sample filtering
300 (Figure 2c). The additional three results came from only PCA (Figure 2a),
301 STRUCTURE (Appendix L) and NetView (Figure 2b). The phylogenetic trees
302 (Figures 3 and 4) were not used for the illegal trade assignments due to the
303 presence of introgression and the poor resolution of the trees.

304 There is perhaps some further grey support for our genetic assignments according to
305 the locations of the illegal seizures. Therefore, we highlighted in bold cases where
306 the sample seizure location is similar to its genetic cluster result, and the wild
307 samples of known origin labelled as WILD are also highlighted in bold (Figure 5).

308 All necessary research permits and ethics approvals were granted prior to
309 commencement of this project. In particular, NUS Institutional Animal Care and Use
310 Committees (IACUC) approved research methods, and we obtained CITES permits
311 for all pangolin samples.

312 **3. Results**

313 3.1. ddRAD sequencing and SNP discovery

314 Our collection of samples included 97 Sunda pangolins and 1 Chinese pangolin.
315 Following library preparation, 2 Sunda pangolin tissue samples did not yield
316 sufficient DNA fragments within the 250–650 bp range, so only 96 samples
317 underwent Illumina sequencing (Appendix P). In total we obtained 49.67 GB of data
318 from 96 samples across two lanes (reads one = 24.12 GB, reads two = 25.55 GB),
319 and there were no lane differences in the results.

320 After demultiplexing our reads, 7 Sunda pangolin samples were discarded due to low
321 sequencing coverage (Appendix P). Across the remaining 89 samples 60,197 SNPs
322 were called via the STACKS pipeline. In PLINK 43,439 loci were removed due to
323 having $\geq 10\%$ missing data. Only 83 samples met our $< 15\%$ missing data
324 requirement; the Chinese pangolin sample was excluded at this stage. A further
325 4,608 loci were removed due to correlation. Consequently, for further analysis we
326 had 83 Sunda pangolin samples with 12,150 SNPs. BayeScan did not detect
327 selection.

328 3.2. Population Genomic Structure

329 3.2.1. Principal Component Analysis and F_{st} estimation

330 The default settings of SNPRelate removed a further 223 SNPs prior to PCA of
331 11,927 SNPs. The PCA results suggested that we were dealing with three distinct
332 genetic clusters of Sunda pangolins, possibly from Sumatra/Singapore, Java, and
333 Borneo according to the genetic cluster labels of wild pangolins of known origin

334 (Figure 2a). The first two eigenvectors of PCA held the largest percentage of
335 variance among the population, principal component 1 = 4.47 % and principal
336 component 2 = 4.38 %. The mean F_{st} between clusters is: Java versus
337 Sumatra/Singapore = 0.0684, Borneo versus Java = 0.0556, Borneo versus
338 Sumatra/Singapore = 0.0446, and across all three clusters = 0.0545. The mean F_{st}
339 values increased when we removed samples showing signals of introgression (see
340 section 3.3.2 for details of introgression), for example, the mean F_{st} across all three
341 clusters increased when MZBR 0270 was removed, mean F_{st} = 0.0556 (Appendix J).

342 3.2.2. Bayesian clustering and Network Analyses

343 The results from STRUCTURE, Structure Harvester and CLUMPP supported the
344 PCA results, similarly indicating that three genetic clusters across the 83 Sunda
345 pangolins is the most likely arrangement (Appendix L). Moreover, all of the clustering
346 approaches across PCA (Figure 2a), STRUCTURE (Appendix L), NetView (Figure
347 2b) and fineRADstructure (Figure 2c) produced compatible clustering results.

348 3.2.2.1. Additional insight to population substructure

349 The co-ancestry matrix from fineRADstructure used 46, 274 loci to illustrate the
350 levels of co-ancestry across 80 Sunda pangolin samples (Figure 2c and Appendix S),
351 with yellow representing the lowest levels of co-ancestry. The matrix confirmed that
352 there are three main clusters, Borneo, Java and Sumatra/Singapore, and revealed
353 additional differences in co-ancestry among clusters/samples. The black colour
354 represents the highest level of co-ancestry and suggests that MZBR 1030 and 1031
355 might be highly related pangolins from Sumatra. The purple coloured MZBR 1166
356 and 1167 could be closely related pangolins from Borneo. The deep red colour
357 represents samples MZBR 1184, 1185, 1190, 1063 and 1060 and suggests these
358 might be close relatives perhaps from a similar area in Java. We were aware of the
359 relationships between MZBR 1030 and 1031, also 1166 and 1167, from the NetView
360 analyses (Appendix M). However, we were not aware about the groupings of
361 pangolins labelled with deep red colours which show very high levels of co-ancestry.
362 FineRADstructure provided more compelling detail about population substructure
363 than our other population genomic analyses.

364 3.3. Phylogenetic trees

365 3.3.1. Phylogenomic tree using SNPs

366 Phylogenomic analysis with the concatenation method RAxML used 2,365 SNPs,
367 and although the Javan samples in the RAxML tree did not form a monophyletic
368 group (Figure 3a), the results were otherwise compatible with our population
369 clustering results from PCA, STRUCTURE, NetView and fineRADstructure (Figures
370 2a-c and Appendix L). Overlaid coloured shading illustrates the population genomic
371 clustering results on top of the RAxML tree to aid comparison (Figure 3a). The
372 RAxML trees generally had low bootstrap support. A few samples also had high
373 missing data, including MZBR 1189, 1183 and 1179 (Appendix O). We tested a wide
374 variety of parameters in both pyRAD and RAxML to try to improve the bootstrap
375 support, for example, minimum coverage for a cluster (5 to 10), maximum number of
376 sites with quality score less than 20 (4 to 6), clustering threshold (0.85 to 0.95), and
377 minimum samples in a final locus (70 to 81), however, the results did not improve.

378 3.3.2. Mitochondrial DNA tree and the issue of introgression

379 The results of the mtDNA tree (Figure 4a) and the mtDNA haplotype network (Figure
380 4b) were congruent. The mtDNA tree was not well resolved (Figure 4a) and the
381 haplotype network clearly shows the differences to the SNP results (Figure 4b). Our
382 ABBA BABA tests of secondary gene flow using 2,365 SNPs (Appendix R)
383 demonstrated that there are signals of introgression across our samples, for example,
384 MZBR 0270 showed a signal of introgression from both Singapore/Sumatra and
385 Borneo. The test result for MZBR 1040 was not statistically significant. The
386 statistically significant ABBA BABA test results are summarised in a simple cartoon
387 diagram (Figure 3b).

388 3.4. Tracing illegal trade

389 Based on the congruent genetic clustering results from PCA (Figure 2a),
390 STRUCTURE (Appendix L), NetView (Figure 2b) and fineRADstructure (Figure 2c),
391 which all showed three key distinct clusters which were geographically labelled by
392 the wild samples of known origin, we conclude that 20 of our samples possibly
393 originated from Java, 21 are possibly from Sumatra/Singapore, and 42 are possibly
394 from Borneo (Figure 5). These three key clusters are represented by black squares
395 (Figure 5) beneath which a full list of each sample's seizure location or wild location

396 has been provided. The samples highlighted in bold do not appear to have been
397 translocated overseas, their seizure location is similar to their genetic origin. The pie
398 charts (Figure 5) show the trade route at point of seizure of the samples, they do not
399 include the wild samples. The arrows help to clarify the pie charts by pointing in the
400 geographic direction of the trade.

401 Among the 49 pangolins seized from illegal wildlife trade in Java (Figure 1), only 18
402 may have been sourced within Java, while 23 may instead have been imported from
403 Borneo, and 8 from Sumatra. The 6 pangolins seized in Medan, North Sumatra
404 (Figure 1), may have originated in Sumatra as these samples group with the
405 Sumatra/Singapore cluster. All of the 16 pangolins seized in Borneo (Figure 1) likely
406 originated there. Only one pangolin across our dataset seems to have been traded
407 from Borneo into Sumatra.

408 **4. Discussion**

409 Our study is the first application of next-generation sequencing methods to the
410 critically endangered Sunda pangolin. The SNP data from each analysis that we
411 have employed suggests there are three key genetic clusters across our samples,
412 which are likely from Borneo, Java and Singapore/Sumatra. The majority of trade
413 across our samples seems to be from Borneo into Java. The presence of
414 introgression across our samples likely explains the poor bootstrap support of the
415 phylogenetic trees. Only the population genomic analyses, PCA, STRUCTURE,
416 NetView and fineRADstructure, provided sufficient clarity to genetically match the
417 pangolins of unknown origin from illegal trade to the wild samples of known origin.

418 The most compelling detail about population substructure was generated by the
419 fineRADstructure method, and we recommend the use of this programme for further
420 research. FineRADstructure not only confirmed our three main clusters, Borneo,
421 Java and Singapore/Sumatra, it revealed further insights into the substructure of
422 these clusters, such as highly related pangolins and other shared co-ancestry. There
423 might be a geographic basis for some of the further substructures which
424 fineRADstructure identified, such as the deep red clusters showing very high levels
425 of shared co-ancestry (Figure 2c), but we require further wild samples of known
426 origin to geographically label these clusters to more concise geographic areas.
427 Currently, the deep orange coloured group within the Singapore/Sumatra cluster is

428 the only subgroup which we can geographically label. The wild Singaporean samples
429 label that cluster as Singaporean.

430 Genetic data are essential to help trace illegal wildlife trade and create forensic
431 genetic databases (Ogden et al. 2009; Wasser et al. 2015; Ogden & Linacre 2015).
432 Although we understand that without the use of trained conservation dogs it is very
433 difficult to capture rare and nocturnal pangolins, which cannot be easily baited, we
434 urge further sampling of wild pangolins of known origin to help inform geographic
435 population assignments. In Singapore our wild sampling was more extensive than
436 elsewhere in Southeast Asia due to a 24hr /7 day per week rescue service for wildlife,
437 which is managed by ACRES, as well as rehabilitation of pangolins at Singapore Zoo.
438 We hope that other partnerships such as this between NGOs and researchers could
439 facilitate improved pangolin sampling elsewhere. It is important that during the
440 collection of reference samples full chains of custody are documented with expert
441 verification.

442 The possible directions of illegal trade which we have revealed in this study (Figure 5)
443 are similar to patterns of illegal trade of other wild species, for example, it has been
444 documented that birds are also poached in large numbers across Indonesia and
445 transported to markets in Jakarta and Java for sale or onwards distribution (Chng et
446 al. 2015). There is some further grey support for our population assignments from
447 our seizure locations, for example, from the local police seizures in Borneo all of
448 those seized pangolins clustered within our probable Bornean cluster (Figure 5).

449 The distinct genetic clustering of populations from Borneo, Java and
450 Singapore/Sumatra matches divergence patterns seen in other vertebrate taxa
451 (Wilson & Reeder, 2005; Leonard et al. 2015). It is possible that our Sunda pangolin
452 lineages might not only be genomic but also ecologically differentiated, since the
453 habitats on Java are primarily monsoonal dryland and savannah, whereas the
454 habitats on the other sampled areas are primarily rainforest (Sundevall, 1843;
455 Whitten et al. 1996). Studies of ecology and morphology across the three Sunda
456 pangolin lineages revealed in this study should be conducted to help inform their
457 taxonomic classification (Gaubert & Antunes 2005), and phenotypic inquiry may
458 uncover traits that support the genomic divergences shown by our data.

459 The genetic divisions were not well detected by our mtDNA data, partly due to
460 admixture and introgression. It is possible that genetic introgression explains the
461 differences between our tentative RAxML (Figure 3a) and mtDNA trees (Figure 4a).
462 Also mtDNA is a single marker that only reflects the maternal history (Ballard &
463 Whitlock, 2004). We suggest that for the population assignment of illegally traded
464 pangolins the use of genome-wide SNP markers can provide higher resolution than
465 mtDNA markers alone. We suggest that future population clustering results using
466 SNP data could be obtained quickly with just fineRADstructure (Malinsky et al. 2016)
467 and more georeferenced samples of known origin are needed.

468 Sunda pangolins are also distributed across Thailand, Myanmar, Cambodia, Laos
469 and Vietnam, and it is possible that greater substructure exists across their full
470 geographic range (Zhang et al. 2015). Protecting genetic diversity is important for the
471 resilience and survival of species (Pierson et al. 2016), and this genetic diversity
472 must therefore be considered in conservation management plans, not least to inform
473 conservation breeding (Hua et al. 2015), including captive breeding and genetic
474 rescue, and repatriation of rescued pangolins. Selection of pangolin release
475 locations should be mindful of the appropriate genetic source population (Challender
476 et al. 2014b). Currently, there is no consideration of genetic substructure when
477 pangolin seizures of unknown origin are released into the wild, and the impacts of
478 introducing and mixing pangolins from different lineages into the same area are not
479 fully understood. The genetic introgression uncovered in our study might possibly be
480 due to human influences, such as translocation of pangolins outside of their natural
481 range (Pantel & Chin, 2009). Conservation breeding records also need to document
482 and consider the highest resolution of population genetic structure to ensure that
483 genetic diversity is well-managed (Allendorf et al. 2010). We wish to emphasize the
484 importance of genetic screening of all individuals involved in releases and
485 conservation breeding.

486 As stated in the IUCN Pangolin Specialist Group action plan (Challender et al. 2014),
487 all pangolin species require further genetic investigation. In our study we chose the
488 double-digest form of RADseq instead of other forms, because this method might be
489 more reproducible for other pangolin researchers to follow, as well as less expensive
490 for researchers with smaller budgets (Andrews et al. 2016).

491 The findings of our research impact a range of conservation actions, including
492 tracing illegal wildlife trade, delimitation of pangolin conservation units, repatriation of
493 rescued pangolins to appropriate locations, and conservation breeding. To maximise
494 the resolution of genetic tracing of pangolins, we recommend that genome-wide
495 markers are used in combination with mtDNA genes. Our findings provide a new
496 baseline to help begin to understand Sunda pangolin populations, and we hope
497 these findings will inspire future research and management actions that can support
498 effective conservation of pangolins in Asia.

499 **Acknowledgments**

500 All authors contributed equally to this work. All authors discussed the results and
501 implications and commented on the manuscript at all stages. We thank the
502 Indonesian Institute of Sciences (LIPI), Lee Kong Chian Natural History Museum
503 (LKCNHM), Agri-Food & Veterinary Authority of Singapore (AVA), Wildlife Reserves
504 Singapore (WRS), Department of Wildlife and National Parks Peninsular Malaysia
505 (DWNP), Universiti Malaysia Terengganu (UMT) and Kadoorie Farm and Botanic
506 Garden (KFBG) for assistance with sample collection and the arrangement of
507 relevant permits and permissions. Special thanks are given to Yulianto (LIPI) and H.
508 Zhang (KFBG) who helped to extract DNA, and S. Oh (WRS), A. Ali (WRS), S. Luz
509 (WRS), P. Lee (WRS), C. F. Maosheng (LKCNHM), M. Chua (LKCNHM), R. Meier
510 (LKCNHM), C.Y. Gwee (NUS), G. Ades (KFBG) and A. Gioni (KFBG). We also
511 thank R. Asher (University of Cambridge) for his encouragement and advice, and the
512 IUCN-SSC Pangolin Specialist Group and Singapore Pangolin Working Group who
513 provided logistical support throughout this research. The Rheindt Lab at NUS shared
514 a ddRADseq protocol. L. Wijedasa gave us base maps for the figures. S. Thompson
515 helped with coding for fineRADstructure. Funding was provided by the Dennis Gould
516 Foundation; H.C.N is supported by a SINGA PhD Research Scholarship at NUS;
517 S.T.T is supported by a Royal Society University Research Fellowship (UF130573).
518 We also acknowledge internal departmental funding at the Department of Biological
519 Sciences at NUS; and the Research Center for Biology-LIPI Competitive Project
520 3400.001.002.021 SEAMEO BIOTROP DIPA 060.12/PSRP/SPK-PNLT/2014.

521 **Literature Cited**

- 522 Alacs EA, et al. 2010. DNA detective: A review of molecular approaches to wildlife
523 forensics. *Forensic Science, Medicine, and Pathology*. **6**(3): 180–194.
- 524 Allendorf FW, Hohenlohe PA, Luikart G. 2010. Genomics and the future of
525 conservation genetics. *Nature Reviews Genetics*. **11**(10): 697–709.
- 526 Andrews KR, et al. 2016. Harnessing the power of RADseq for ecological and
527 evolutionary genomics. *Nature Reviews Genetics*. **17**: 81–92.
- 528 Andrews S. 2016. FastQC: a quality control tool for high throughput sequence data.
529 Version 0.11.5 Available at:
530 <http://www.bioinformatics.babraham.ac.uk/projects/fastqc>
- 531 Ballard JWO, Whitlock MC. 2004. The incomplete natural history of mitochondria.
532 *Molecular Ecology*. **13**(4): 729–744.
- 533 Catchen J, et al. 2013. Stacks: an analysis tool set for population genomics.
534 *Molecular Ecology*. **22**(11): 3124–40.
- 535 Challender DWS, Waterman C, Baillie JEM. 2014a. Scaling up pangolin
536 conservation. IUCN SSC Pangolin Specialist Group Conservation Action Plan.
537 Zoological Society of London, London, UK.
- 538 Challender DWS, et al. 2014b. *Manis javanica*. The IUCN Red List of Threatened
539 Species 2014: e.T12763A45222303. Available at:
540 <http://dx.doi.org/10.2305/IUCN.UK.2014-2.RLTS.T12763A45222303.en>.
- 541 Challender DWS, Harrop SR, MacMillan DC. 2015. Understanding markets to
542 conserve trade-threatened species in CITES. *Biological Conservation*. **187**:
543 249–259.
- 544 Cheng W, Xing S, Bonebrake TC. 2017. Recent Pangolin Seizures in China Reveal
545 Priority Areas for Intervention. *Conservation Letters*. doi:10.1111/conl.12339
- 546 Chng SCL, Eaton JA, Krishnasamy K, Shepherd CR, Nijman V. 2015. In the Market
547 for Extinction: An inventory of Jakarta’s bird markets. TRAFFIC. Petaling Jaya,
548 Selangor, Malaysia.

549 Choo SW et al. 2016. Pangolin genomes and the evolution of mammalian scales and
550 immunity. *Genome Research*. **26**(10): 1312-1322.

551 Çilingir FG, Rheindt FE, Garg KM, Platt K, Platt SG, Bickford DP. 2017.
552 Conservation genomics of the endangered Burmese roofed turtle. *Conservation*
553 *Biology*. **28**: in press.

554 Ciofi C, et al. 1999. Genetic divergence and units for conservation in the Komodo
555 dragon *Varanus komodoensis*. *Proceedings of the Royal Society B*. **266**: 2269-
556 2274.

557 Clarke S, et al. 2006. Identification of shark species composition and proportion in
558 the Hong Kong shark fin market based on molecular genetics and trade records.
559 *Conservation Biology*. **20**: 201-211.

560 Corlett RT. 2016. A Bigger Toolbox: Biotechnology in Biodiversity Conservation.
561 *Trends in Biotechnology*. **35**(1): 55-65.

562 Drummond AJ, et al. 2012. Bayesian phylogenetics with BEAUti and the BEAST 1.7.
563 *Molecular Biology and Evolution*. **29**(8): 1969–1973.

564 Earl DA, vonHoldt BM. 2012. STRUCTURE HARVESTER: A website and program
565 for visualizing STRUCTURE output and implementing the Evanno method.
566 *Conservation Genetics Resources*. **4**(2): 359–361.

567 Eaton DAR. 2014. PyRAD: assembly of de novo RADseq loci for phylogenetic
568 analyses. *Bioinformatics*. **30**(13): 1844-1849.

569 Evanno G, Regnaut S, Goudet J. 2005. Detecting the number of clusters of
570 individuals using the software STRUCTURE: A simulation study. *Molecular*
571 *Ecology*. **14**(8): 2611–2620.

572 Foll M, Gaggiotti O. 2008. A genome-scan method to identify selected loci
573 appropriate for both dominant and codominant markers: A Bayesian perspective.
574 *Genetics*. **180**(2): 977–993.

575 Garg KM, et al. 2016. Genome-wide data help identify an avian species-level lineage
576 that is morphologically and vocally cryptic. *Molecular Phylogenetics and*
577 *Evolution*. **102**: 97–103.

- 578 Gaubert P. 2011. Family Manidae. In: Wilson, D.E., Mittermeier, R.A. (Eds.),
579 Handbook of the Mammals of the World, Vol. 2: Hoofed Mammals. Lynx
580 Edicions, Barcelona, Spain.
- 581 Gaubert P, Antunes A. 2005. Assessing the taxonomic status of the Palawan
582 pangolin *Manis Culionensis* (Pholidota) using discrete morphological characters.
583 *Journal of Mammalogy*. **86**(6): 1068–1074.
- 584 Hassanin A, Hugot JP, van Vuuren BJ. 2015. Comparison of mitochondrial genome
585 sequences of pangolins (Mammalia, Pholidota). *Comptes Rendus Biologies*.
586 **338**(4): 260–265.
- 587 Hua L, et al. 2015. Captive breeding of pangolins: Current status, problems and
588 future prospects. *ZooKeys*. **507**: 99–114.
- 589 Jakobsson M, Rosenberg NA. 2007. CLUMPP: A cluster matching and permutation
590 program for dealing with label switching and multimodality in analysis of
591 population structure. *Bioinformatics*. **23**(14): 1801–1806.
- 592 Johnson WE, et al. 2010. Genetic restoration of the Florida panther. *Science* **329**:
593 1641-1645.
- 594 Knowles LL, et al. 2016. Quantifying the similarity between genes and geography
595 across Alaska's alpine small mammals. *Journal of Biogeography*. **43**(7): 1464–
596 1476.
- 597 Kumar S, Stecher G, Tamura K. 2016. MEGA7: Molecular Evolutionary Genetics
598 Analysis Version 7.0 for Bigger Datasets. *Molecular Biology and Evolution*. **33**(7):
599 1870-4.
- 600 Lawson DJ, Hellenthal G, Myers S, Falush D. 2012. Inference of Population
601 Structure using Dense Haplotype Data. *PLoS Genetics*. **8**(1): e1002453.
- 602 Leigh JW, Bryant D. 2015. POPART: full-feature software for haplotype network
603 construction. *Methods in Ecology and Evolution*. **6**: 1110–1116.
- 604 Leonard JA, et al. 2015. Phylogeography of vertebrates on the Sunda Shelf: a multi-
605 species comparison. *Journal of Biogeography*. **42**: 871–879.

606 Li H, et al. 2009. The Sequence Alignment/Map format and SAMtools. *Bioinformatics*,
607 **25**(16): 2078–2079.

608 Li H, et al. 2013. Burrows-Wheeler Alignment Tool v0.7.1. Available at:
609 <http://bio-bwa.sourceforge.net/bwa.shtml>

610 Lischer HEL, Excoffier L. 2012. PGDSpider: An automated data conversion tool for
611 connecting population genetics and genomics programs. *Bioinformatics*. **28**(2):
612 298–299.

613 Lohman DJ, et al. 2010. Cryptic genetic diversity in “widespread” Southeast Asian
614 bird species suggests that Philippine avian endemism is gravely underestimated.
615 *Biological Conservation*. **143**: 1885-1890.

616 Malenfant RM, Coltman DW, Davis CS. 2014. Design of a 9K illumina BeadChip for
617 polar bears (*Ursus maritimus*) from RAD and transcriptome sequencing.
618 *Molecular ecology resources*. **15**(3): 587-600.

619 Malinsky M, Trucchi E, Lawson D, Falush D. 2016. RADpainter and
620 fineRADstructure: population inference from RADseq data. *BioRxiv* 057711.
621 (pre-print). doi: <https://doi.org/10.1101/057711>.

622 Nash HC, Wong MHG, Turvey ST. 2016. Using local ecological knowledge to
623 determine status and threats of the critically endangered Chinese pangolin
624 (*Manis pentadactyla*) in Hainan, China. *Biological Conservation*. **196**: 189-195.

625 Ng NS, Wilton PR, Prawiradilaga DM, Tay YC, Indrawan M, Garg KM, Rheindt FE.
626 2017. The effects of Pleistocene climate change on biotic differentiation in a
627 montane songbird clade from Wallacea. *Molecular Phylogenetics and Evolution*.
628 **114**: 353-366.

629 Ogden R, Dawnay N, McEwing R. 2009. Wildlife DNA forensics - Bridging the gap
630 between conservation genetics and law enforcement. *Endangered Species*
631 *Research*. **9**(3): 179–195.

632 Ogden R, Linacre A. 2015. Wildlife forensic science: A review of genetic geographic
633 origin assignment. *Forensic Science International: Genetics*. **18**: 152–159.

634 Pantel S, Chin SY. 2009. Proceedings of the Workshop on Trade and Conservation
635 of Pangolins native to South and Southeast Asia. TRAFFIC Southeast Asia.
636 Petaling Jaya, Selangor, Malaysia.

637 Peterson BK, et al. 2012. Double digest RADseq: An inexpensive method for de
638 novo SNP discovery and genotyping in model and non-model species. PLoS
639 ONE. **7**(5).

640 Pierson JC, et al. 2016. Genetic factors in threatened species recovery plans on
641 three continents. *Frontiers in Ecology and the Environment*. **14**(8): 433–440.

642 Pritchard JK, Stephens M, Donnelly P. 2000. Inference of population structure using
643 multilocus genotype data. *Genetics*. **155**(2): 945–959.

644 Puechmaille SJ. 2016. The program structure does not reliably recover the correct
645 population structure when sampling is uneven: Subsampling and new estimators
646 alleviate the problem. *Molecular Ecology Resources*. **16**(3): 608–627.

647 Purcell S, et al. 2007. PLINK: A tool set for whole-genome association and
648 population-based linkage analyses. *American Journal of Human Genetics*. **81**(3):
649 559–575.

650 Rambaut A, et al. 2014. Tracer v1.6. Available at: <http://beast.bio.ed.ac.uk/Tracer>

651 Rambaut A, Drummond AJ. 2016. TreeAnnotator v2.4.2. Institute of Evolutionary
652 Biology, University of Edinburgh.

653 R Core Team 2016. R: A Language and Environment for Statistical Computing. R
654 Foundation for Statistical Computing, Vienna, Austria. Available at: [http://www.R-](http://www.R-project.org/)
655 [project.org/](http://www.R-project.org/)

656 RStudio Team 2015. RStudio: Integrated Development for R. RStudio, Inc., Boston,
657 MA. Available at: <http://www.rstudio.com/>

658 Schueller AM, Hayes DM. 2011. Minimum viable population size for lake sturgeon
659 (*Acipenser fulvescens*) using an individual-based model of demographics and
660 genetics *Canadian Journal of Fisheries and Aquatic Sciences*. **68**: 62-73.

661

662 Segan DB, et al. 2010. Using conservation evidence to guide management.
663 Conservation Biology. **25**(1): 200-202.
664

665 Shirley MH, et al. 2014. Rigorous approaches to species delimitation have significant
666 implications for African crocodylian systematics and conservation. Proceedings
667 of the Royal Society B. **281**: 20132483.

668 Stamatakis A. 2014. RAxML version 8: a tool for phylogenetic analysis and post-
669 analysis of large phylogenies. Bioinformatics. **30**.

670 Sundevall CJ. 1843. Wer sicut der gattung Manis. Kongl. Svenska vetenskaps-
671 akademiens handlingar. 1842: 245-283.

672 Sutherland WJ, et al. 2004. The need for evidence-based conservation. Trends in
673 Ecology and Evolution. **19**: 305-308.

674 Tan TK, et al. 2016. PGD: a pangolin genome hub for the research community.
675 Database (Oxford). Available at:
676 <http://database.oxfordjournals.org/content/2016/baw063.full>.

677 Tay YC, Chng MWP, Sew WWG, Rheindt FE, Tun KPP, Meier R. 2016. Beyond the
678 Coral Triangle: high genetic diversity and near panmixia in Singapore's
679 populations of the broadcast spawning sea star *Protoreaster nodosus*. Royal
680 Society Open Science. **3**: 160253.

681 United Nations Office on Drugs and Crime (UNODC). 2016. World Wildlife Crime
682 Report: Trafficking in protected species. Available at:
683 [https://www.unodc.org/documents/data-and-](https://www.unodc.org/documents/data-and-analysis/wildlife/World_Wildlife_Crime_Report_2016_final.pdf)
684 [analysis/wildlife/World_Wildlife_Crime_Report_2016_final.pdf](https://www.unodc.org/documents/data-and-analysis/wildlife/World_Wildlife_Crime_Report_2016_final.pdf)

685 Wasser SK, et al. 2008. Combating the illegal trade in African elephant ivory with
686 DNA forensics. Conservation Biology. **22**: 1065-1071.

687 Wasser SK, et al. 2015. Genetic assignment of large seizures of elephant ivory
688 reveals Africa's major poaching hotspots. Science. **349**(6243): 84–88.

689 WCS News Releases. April 2015. [https://newsroom.wcs.org/News-](https://newsroom.wcs.org/News-Releases/articleType/ArticleView/articleId/6715/April-27-Indonesian-National-)
690 [Releases/articleType/ArticleView/articleId/6715/April-27-Indonesian-National-](https://newsroom.wcs.org/News-Releases/articleType/ArticleView/articleId/6715/April-27-Indonesian-National-)

691 Police-Seize-Major-Shipment-of-Pangolins-Arrest-Smuggler.aspx

692 Whitten T, Soeriaatmadja RE, Afiff SA. 1996. The Ecology of Java and Bali. The
693 Ecology of Indonesia Series, II. Periplus Editions (HK) Limited.

694 Wilson DE, Reeder DM. (editors). 2005. Mammal Species of the World. A
695 Taxonomic and Geographic Reference (3rd ed). Johns Hopkins University Press.

696 Wiradateti, Semiadi G. 2017. Genetic Variation of Confiscated Pangolins of Sumatra,
697 Java, and Kalimantan based on Control Region Mitochondrial DNA. Jurnal
698 Veteriner. **18**(2): 181-191.

699 Zhang H, et al. 2015. Molecular tracing of confiscated pangolin scales for
700 conservation and illegal trade monitoring in Southeast Asia. Global Ecology and
701 Conservation. **4**: 414–422.

702 Zinenko O, et al. 2016. Hybrid origin of European Vipers (*Vipera magnifica* and
703 *Vipera orlovi*) from the Caucasus determined using genomic scale DNA markers.
704 BMC Evolutionary Biology. **16**:76.

705

706 **Figure Captions**

707 Fig. 1 Locations of seizures of illegally traded pangolins and wild Sunda pangolins of
708 known origin across Indonesia (Java, Kalimantan and Sumatra), Malaysia and
709 Singapore that were used in this study. Numbers in brackets represent the number
710 of pangolins sampled per location

711 Fig. 2 Population genetic clusters. The pangolin samples within each geographically
712 labelled cluster are similar across each figure. (a) PCA of 11,927 SNPs across 83
713 Sunda pangolins. Principal component 1 = 4.47 % and principal component 2 =
714 4.38 %. (b) Network Analysis with the highest genetic distance of clusters using
715 12,150 SNPs across 83 Sunda pangolins from the Walktrap model. (c) Clustered
716 fineRADstructure coancestry matrix using 46, 274 loci across 80 Sunda pangolins

717 Fig. 3 (a) Maximum likelihood phylogeny in RAxML using 2,365 SNPs. The coloured
718 overlaid shading illustrates the clustering results from other methods, PCA, Structure,
719 NetView and fineRADstructure to facilitate comparison with those results: blue =
720 Borneo, green = Singapore/Sumatra, no colour = Java. (b) The directions of detected
721 introgression between pangolin lineages indicated by dashed arrows, from ABBA
722 BABA tests using 2,365 SNPs

723 Fig. 4 (a) mtDNA phylogeny, with concatenated *Cytb* and *CO1* data. Sample labels
724 begin with a key of the SNP cluster result to facilitate comparison: B = Borneo, J =
725 Java, S = Singapore/Sumatra. Seizure locations are given following the sample
726 name. Asterisks indicate samples where seizure location and SNP cluster result
727 were the same geographic area. (b) mtDNA haplotype network. The colour of
728 sample label represents the SNP cluster result to facilitate comparison: blue =
729 Borneo, green = Singapore/Sumatra, purple = Java, black = mtDNA only. Wild
730 samples of known origin are labelled as WILD.

731 Fig. 5 Inferred directions of the illegal trade of pangolins. Varied population genomic
732 analyses using > 12,000 SNPs all provided compatible results (PCA, Structure,
733 NetView and fineRADstructure). A full list of each sample's seizure location or wild
734 location is provided below each genetic cluster. The pangolins highlighted in bold
735 were not translocated overseas, their seizure location is similar to their genetic origin.
736 Pie charts show the trade route at point of seizure of the pangolin (they do not

737 include the wild samples). Arrows show the direction of trade, with the largest arrow
738 reflecting the highest volume of trade, and the dashed arrows reflecting less trade









