# High-throughput genotyping of a full voltage-gated sodium channel gene via genomic DNA using target capture sequencing and analytical pipeline MoNaS to discover novel insecticide resistance mutations

Kentaro Itokawa[1,2,3], Koji Yatsu[1], Tsuyoshi Sekizuka[1,3], Yoshihide Maekawa[2], Osamu Komagata[2,3], Masaaki Sugiura[4], Tomonori Sasaki[5], Takashi Tomita[2], Makoto Kuroda[1], Kyoko Sawabe[2], and Shinji Kasai[2*]

1. Pathogen Genomics Center, National Institute of Infectious Diseases, Japan
2. Department of Medical Entomology, National Institute of Infectious Diseases, Japan
3. Antimicrobial Resistance Research Center, National Institute of Infectious Diseases, Japan
4. Global Research and Development Department, Fumakilla Limited, Japan
5. Research and Development Department, Fumakilla Limited, Japan

*To whom correspondence should be addressed: kasacin@nih.go.jp

17

18          Abstract

19     In insects, voltage-gated sodium channel (VGSC) is the primary target site of pyrethroid

20   insecticides. Various amino acid substitutions in the VGSC protein, which are selected

21   under insecticide pressure, are known to confer insecticide resistance. In the genome, the

22   *VGSC* gene consists of more than 30 exons sparsely distributed across a large genomic

23   region, which often exceeds 100 kbp. Due to this complex genomic structure of *VGSC* gene,

24   it is often challenging to genotype full coding nucleotide sequences (CDSs) of *VGSC* from

25   individual genomic DNA (gDNA). In this study, we designed biotinylated oligonucleotide

26   probes from annotated CDSs of *VGSC* of Asian tiger mosquito, *Aedes albopictus*. The probe

27   set effectively concentrated (>80,000-fold) all targeted regions of gene *VGSC* from pooled

28   barcoded Illumina libraries each constructed from individual *A. albopictus* gDNAs. The

29   probe set also captured all homologous *VGSC* CDSs except some tiny exons from the

30   gDNA of other Culicinae mosquitos, *A. aegypti* and *Culex pipiens* complex, with comparable

31   efficiency as a result of the high nucleotide-level conservation of *VGSC*. To enhance

32   efficiency of the downstream bioinformatic process, we developed an automated pipeline to

33   genotype *VGSC* after capture sequencing—MoNaS (Mosquito Na$^+$ channel mutation

34   Search)—which calls amino acid substitutions and compares those to known resistance

35   mutations. The proposed method and our bioinformatic tool should facilitate the discovery

36   of novel amino acid variants conferring insecticide resistance on VGSC and population

37   genetics studies on resistance alleles (with respect to the origin, selection, and migration

38   etc.) in both clinically and agriculturally important insect pests.

39

40     Key words: Targeted enrichment, Variant analysis, Voltage Gated Sodium Channel,

41   Insecticide resistance, Mosquitos

42

<div align="center">Introduction</div>

Synthetic pyrethroids are currently the most frequently used insecticides for the control of clinically important mosquitos. The mode of pyrethroid's toxicity is inhibition of the voltage-gated sodium channel (VGSC) in the nervous system (Lund & Narahashi, 1983). Developed resistance against pyrethroids, which is known as knockdown resistance (*kdr*), was first reported in housefly, *Musca domestica*, in 1950s (Busvine, 1951). The *kdr* phenotype as well as another distinct phenotype, super-kdr, was eventually linked to amino acid (aa) substitutions on the two positions, L1014F and L1014F+M918T, respectively, on the gene coding VGSC protein (Miyazaki, Ohyama, Dunlap, & Matsumura, 1996; Williamson, Martinez-Torres, Hick, & Devonshire, 1996). Currently, same aa substitutions as well as various other aa substitutions have been found to be associated with resistance in many medical and agricultural insect pests (Dong et al., 2014; Rinkevich, Du, & Dong, 2013). With this historical background, aa substations in variety of insect species are often projected to corresponding aa position in *M. domestica* VGSC for comparison. Actually, VGSC is highly conserved among insects, and many resistance-conferring aa substitutions are seen parallelly in different species (Davies, Field, Usherwood, & Williamson, 2007). Therefore, it is relatively straightforward to infer the effect of certain aa substitutions in any species if the effect of those substitutions has already been elucidated. However, sequencing analysis of an entire coding sequence (CDS) of the *VGSC* gene from genomic DNA (gDNA) is complicated because *VGSC* genes typically consist of many (>30) small exons sparsely distributed across a large genomic region, which often exceeds 100 kbp  Therefore, strategies employing direct sequencing of PCR-amplified genomic fragments usually target only restricted regions where the known resistance-conferring substitutions are frequently found; e.g., IIS5–6 (Dong et al., 2014; Rinkevich et al., 2013). Such a bias may lower the chance of discovering novel resistance mutations existing outside the region investigated.

*Aedes albopictus*, or Asian tiger mosquito, is a medically important mosquito species ubiquitously present on most continents on the Earth. In some regions where the other effective vector, *A. aegypti*, is absent, the species often take a main role for transmitting Chikungunya and Dengue viruses (Kutsuna et al., 2015; Reiter, Fontenille, & Paupy, 2006).

<div align="center">3</div>

72    The *kdr* substitution in *A. albopictus* had not been reported until the F1534C allele was

73    discovered in Singapore 2009 (S Kasai et al., 2011). Since this discovery, F1534C and other

74    *kdr* substitutions at the same aa position, F1534S and F1534L, were reported from in *A.*

75    *albopictus* in various geographic locations worldwide (H. Chen et al., 2016; Marcombe,

76    Farajollahi, Healy, Clark, & Fonseca, 2014; Xu et al., 2016). More recently, we also

77    discovered the new *kdr* substitution V1016G in *A. albopictus* by extending the region of the

78    search for mutations (Shinji Kasai et al., 2019).

79    The next-generation sequencing (NGS) technology has revolutionary reduced the cost and

80    time of DNA sequencing by orders of magnitude. The recent *Anopheles gambiae* 1000

81    Genomes project, Ag1000G, has uncovered a number of previously unknown

82    nonsynonymous mutations in the *VGSC* gene in *A. gambiae* and *A. coluzzii* (Clarkson et al.,

83    2018), some of which have been suspected to cause resistance directly or indirectly. The

84    study also showed that even neutral variations within or flanking the *VGSC* locus represent

85    valuable information to infer the origin and evolution of resistance. Although whole-genome

86    sequencing may discover novel variants of *VGSC* unequivocally, this naive approach is still

87    too costly per sample just for analyzing *VGSC*. Alternatively, we considered an enrichment

88    approach involving hybridization of oligo DNA/RNA (Gnirke et al., 2009) which is often

89    employed to selectively sequence targeted genomic regions for studies e.g. on genotyping

90    of disease-related genes in humans. This technology is aimed at increasing the depth of

91    reads and the number of samples to be multiplexed per given sequencing capacity in return

92    for limiting the region to be analyzed. In this study, we designed biotinylated oligonucleotide

93    DNA probes from *A. albopictus VGSC* CDSs. The probe set efficiently concentrated targeted

94    regions from the gDNA of individual *A. albopictus*. Although the probe set was designed

95    from the *A. albopictus VGSC* gene, the same probe set captured most CDSs of *A. aegypti*

96    and *Culex pipiens* complex as a result of the high nucleotide conservation of *VGSC*. This

97    technology allows for full-CDS analysis of the complex *VGSC* gene in a relatively low-cost

98    and highly multiplexed manner, which is expected to promote discoveries of novel

99    resistance-conferring aa substitutions both in medical and agricultural insect pests.

4

100 <div align="center">Materials and Methods</div>

101 **Design of custom probes**

102    The full-length *VGSC* gene (AALF000723-RA in gene set: AaloF1.2) was found in scaffold

103 JXUM01S000562 in the genome assembly of an *A. albopictus* Foshan strain, AaloF1 (X.-G.

104 Chen et al., 2015) hosted on vectorbase.org (Giraldo-Calderón et al., 2015). Because the

105 annotation missed some CDSs (entire exon 19c for instance), we refined the annotation by

106 aligning shotgun-sequenced NGS reads of *VGSC* cDNA (Shinji Kasai et al., 2019) using

107 Hisat2 (Kim, Langmead, & Salzberg, 2015) and the *M. domestica* VGSC protein sequence

108 (GenBank accession No.: AAB47604) via BLASTX (Altschul, Gish, Miller, Myers, & Lipman,

109 1990). Compared to AaloF1.2, the refined annotation included three added CDSs, and four

110 extended CDSs (see detail in Table S1). Among the 35 coding CDSs in total, sizes of 34 (32

111 + 2 mutually excluding exons) CDSs matched to the *A. aegypti VGSC* CDSs annotated by

112 Davies et al. (2007). Therefore, the numbering of exons in this paper was set to be

113 concordant with the *A. aegypti VGSC* exons described by Davies et al. (2007). The

114 additional optionally used 45 bp small exon, referred to as exon 16.5 here, was found in

115 cDNA data between exons 16 and 17 in the genome. All CDS sequences, some of which

116 contained a flanking intronic region (for tiny exons less than 120 bp in size) were submitted

117 to the IDT website (https://sg.idtdna.com) to design 120 bp xGen Lockdown biotinylated

118 oligonucleotide DNA probes with2× Tiling density option. We also included some exons of

119 other genes flanking *VGSC* (AALF020128, AALF020129, AALF020130, AALF000725,

120 AALF000726, AALF000727, AALF000728, and AALF000730) or a gene nested in the

121 intronic region of *VGSC* (AALF020132) during the probe design to take advantage of the

122 population genetic analysis in future studies. From 15 kbp genomic regions in total, 229

123 probes were designed (Table S2), of which 145 target *VGSC*. Nonetheless, in the more

124 recent contiguous assembly of the C6/36 cell line (see below), AALF020132, AALF000725,

125 AALF000726, AALF000727, AALF000728, and AALF000730 are not located in the same

126 assembly with the *VGSC* locus. For this reason, in this paper, we evaluate the performance

127 of the probe set only in terms of *VGSC* CDS enrichment.

<div align="center">5</div>

**Samples**

Fifty-six mosquitos either belonging to species *A. albopictus*, *A. aegypti* or *Culex pipiens* complex—either kept in the laboratory or caught in the wild (Table 1)—served as a source of gDNA. Of those, strains Aalb-SP, Aaeg-SP, and Cpip-JPP were already known to possess haplotypes with 1534C, 989P-1016G, and 1014F aa variants, respectively (Hardstone et al., 2007; S Kasai et al., 2014; Shinji Kasai et al., 2019).

**gDNA extraction**

gDNA was individually extracted from the whole body of an adult or pupa using the MagExtractor Genome Kit (TOYOBO). The protocol was modified to conduct the extraction in 8-strip PCR tubes or a 96-well PCR plate as follows. The whole body of a single insect was homogenized in a PCR well containing 50 µl of the Lysis & Binding Solution and zirconia beads (ø 2 mm; Nikkato) in TissueLyser II (QIAGEN) at 25 Hz for 30 s. After that, the samples were centrifuged at 2000 × *g* for 1 min to precipitate large debris, and each supernatant was transferred to a new well containing 50 µl of the Lysis & Binding Solution and 5 µl of DNA-binding Magnetic Beads. The solution was shaken in MicroMixer E-36 (TAITEC) at the maximum speed (2500 rpm) for 10 min, and then, on a magnetic plate, the supernatant was discarded. The beads bound to DNA were washed twice with 100 µl of the Washing Solution and twice with 75% ethanol each. Finally, DNA was eluted with 50 µl of low-TE buffer (0.1 mM EDTA, 10 mM Tris-HCl pH 8.0) by shaking in MicroMixer E-36 at the maximum speed for 10 min. The obtained DNA was quantified with the Qubit Highly Sensitive DNA Assay Kit (Invitrogen). The obtained DNA concentration ranged from 2.3 to 8.6 ng/µl for *A. albopictus*, 5.3 to 10 ng/µl for *A. aegypti*, and 7.8 to 11 ng/µl for *C. pipiens* complex mosquitos.

**Library construction and hybridization capture**

Illumina libraries with TruSeq barcode adapters were prepared using NEBNext Ultra II FS DNA Library Prep Kit for Illumina (NEB) on the 1/4 scale of the manufacturer-suggested protocol. Briefly, 4 µl of the gDNA extracted above (without adjusting the concentration) was mixed with 0.4 µl of the Enzyme Mix, 1.4 µl of Reaction Buffer, and 1.2 µl of $H_2O$ on ice. The

6

156  mixture was then incubated at 37 °C for 10 min followed by incubation at 65 °C for 30 min.

157  Those end-prepped DNAs were directly ligated with Illumina adapters by the addition of 0.5

158  µl of TruSeq 96 dual-index adapters (Illumina) instead of adapters supplied with the kit, 6 µl

159  of the Ligation Master Mix, and 0.2 µl of Ligation Enhancer, with incubation at 20 °C for 15

160  min. Next, the libraries were incubated at 65 °C for 30 min to inactivate the ligase; then, all

161  the 56 libraries were pooled together in a single 1.5 ml LoBind tube (Eppendrof). The pooled

162  library was purified with 1.2× SPRIselect (Beckman Coulter) and eluted with 20 µl of low-

163  TE buffer. A 7 µl aliquot of the pooled library was aliquoted and mixed with 0.8 µl of a 10

164  µg/µl UltraPure Salmon Sperm DNA Solution (Invitrogen) in a PCR-tube. Then. the mixsture

165  was concentrated by incubation at 80 °C for 10 min while the lid of the tube and thermal

166  cycler were opened. The concentrated library mix was hybridized, captured, and washed

167  with the designed oligo DNA probe set and the xGen Hybridization and Wash Kit (IDT). After

168  that, the streptavidin magnetic beads were subjected to PCR amplification with HiFi Kapa

169  (Kapa Biosystems) for 12 cycles. The amplified library was purified with 1.2× SPRIselect

170  beads and quantified by real-time PCR using P5 and P7 adapter primers and qPCR double

171  quencher probe (6-FAM)-5′-ACACTCTTT-(ZEN)-CCCTACACGACGCTCTTC-3′-(Iowa

172  Black FQ) (IDT) in the PrimeTime Gene Expression Master Mix (IDT). Serial dilutions of the

173  phiX library (Illumina) were used for construction of the standard curve. The quantified

174  library was sequenced on Illumina MiniSeq with the Mid Output Kit (Illumina) for 151 cycles

175  from both ends along with the libraries from other studies.

176  **Reference genomes and annotation for *VGSC***

177  Although the probe sets were designed from assembly AaloF1, we chose a C6/36 cell line

178  genome assembly, canu_80X_arrow2.2 (Miller et al., 2018), as a reference genome of *A.*

179  *albopictus* for further bioinformatic analysis because this assembly has better contiguity and

180  fewer scaffolds than AaloF1 does. In the canu_80X_arrow2.2 assembly, the whole *VGSC*

181  gene was found in scaffolds MNAF02001058.1 and MNAF02001442.1 annotated as Gene

182  IDs LOC109421922 and LOC109432678, respectively, in the NCBI *Aedes albopictus*

183  Annotation Release 101. The two *VGSC* genes were assumed to be redundant haplotigs.

184  To avoid dual mapping of the NGS reads, we purged MNAF02001442 by hard-masking this

7

185    entire scaffold (replacing all bases with the "N" character) rather than MNAF02001058.1

186    because LOC109432678 in MNAF02001442.1 has a single frame-shifting nucleotide

187    deletion in the thymine homopolymer track within exon 4 (TTTTTT → TTTTT), which was

188    suspected due to an uncorrected base-calling error. LOC109421922 was defined by the

189    number of transcriptional variants in the NCBI's annotation because *VGSC* is known to have

190    complex alternative splicing patterns (Davies et al., 2007). Nevertheless, we simplified the

191    *VGSC* gene model into two possible transcriptional variants to build a GFF3 annotation file

192    for annotating aa changes. These two transcripts include all the regions of mandatory or

193    optional CDSs but differ by the two mutually exclusive exons "19c/k" and "26d/l," where one

194    carries exons "19c" and "26k," and the other contains exons "19d" and "26l." CDSs of all the

195    transcriptional variants of LOC109421922 were merged via overlaps. Those merged CDSs

196    perfectly matched AaloF1 except for LOC109421922 including exon 16.5 and except for

197    one mutually exclusive exon "26k" whose sequence itself was found to be intact in

198    MNAF02001058.1.

199    The *VGSC* gene in the chromosome level assembly of the *A. aegypti* genome, AaegL5.0

200    (Matthews et al., 2018), was annotated in the same manner. The whole *VGSC* gene is

201    encoded as AAEL023266 (the NCBI *Aedes aegypti* Annotation Release 101) on

202    chromosome 3. AAEL023266 has 13 transcripts, these CDSs were merged via overlaps as

203    in the canu_80X_arrow2.2 assembly of *A. albopictus*. AAEL023266 appeared to be lacking

204    an exon corresponding to exon 16.5, whereas we found its sequence between exons 16

205    and 17. AAEL023266, however, contains an additional exon between exons 11 and 12. The

206    21 bp small exon was assumed to correspond to exon "12" in the *Anopheles gambiae*

207    genome described by Davies et al. (2007) and is situated within the intracellular loop

208    between domains I and II. We found the sequence homologous to this exon also in the two

209    *A. albopictus* genome assemblies, AaloF1 and canu_80X_arrow2.2; thus, which means we

210    had failed to include this exon in the probe design. The complete *VGSC* sequence was also

211    found in scaffold NIGP01000811 and was assumed to be a redundant haplotigs. This

212    scaffold was purged from the assembly by hard-masking.

213    In *C. quinquefasciatus* genome assembly Cpip_J2 (Arensburger et al., 2010), the *VGSC*

214    gene is located in scaffold supercont3.182. The *VGSC* gene in supercont3.182, however,

8

215   contains shorter exon 13 which was truncated by scaffolding gap and lacks the entire exon

216   14. Complete exons 13 and 14 were found in another scaffold, supercont3.1170, which

217   contains an incomplete *VGSC* gene probably as a haplotig. We fused contig

218   AAWU01037504.1 containing exons 13 and 14 of *VGSC* from supercont3.1170 into

219   supercont3.182 to restore the complete coding sequence of the *VGSC* gene, thereby

220   creating supercont3.182_2 (Fig. S1A). The *VGSC* in supercont3.182 (and

221   supercont3.182_2) still contained *kdr* aa substitutions, L932F and I936V, as already

222   reported by Davies et al. (2007) plus unusual frameshifting deletions in exons 26l and 32

223   (Fig. S1B). For these reasons, supercont3.182_2 was further polished by the *consensus*

224   module in BCFtools (Danecek & McCarthy, 2017) with the "-H 1" option using the variant

225   information for Cpip-JNA-01 in the VCF file generated as described below, thereby finally

226   resulting in supercont3.182_3. The latter scaffold was added to the genome assembly, and

227   the original scaffolds supercont3.182 and supercont3.1170 were purged by hard-masking.

228   We were not able to find an exon corresponding to "exon 16.5" in *A. albopictus* and *A.*

229   *aegypti*.

230   **Bioinformatic analysis**

231   The FASTQ data were mapped to the reference using *BWA mem* (v.0.7.17) (Li & Durbin,

232   2009) with default options. The resultant BAM files were sorted by the *sort* program from

233   the SAMtools suite (v.1.9) (Li et al., 2009), and we removed PCR duplicates by the *rmdup*

234   programs from the SAMtools. Variant calling was performed on the resulting BAM files of

235   each species in the FreeBayes software (v.1.2.0) (Garrison & Marth, 2012) with default

236   options. The variant annotation (for aa changes) was conducted with the *csq* program from

237   the BCFtools suite (v.1.9) (Danecek & McCarthy, 2017) with options "*-l -p a*". Finally, the

238   discovered aa changes were projected onto the corresponding position in the *M. domestica*

239   VGSC protein sequence (GenBank accession No.: AAB47604). Those bioinformatic

240   processes (Fig. 1A) were automated in pipeline tools *MoNaS* (Mosquito Na$^+$ channel

241   mutation Search; https://github.com/ItokawaK/MoNaS) written in the Python3 script

242   language.

243   For estimating the level of enrichment, five sets of random data on whole-genome shotgun

9

244　　paired-end reads (150 bp × 2, 300 ± 50 bp insert length, 1 million read pairs) from each

245　　reference genomic assembly were simulated in the *wgsim* software

246　　(https://github.com/lh3/wgsim). The *multicov* program from the Bedtools suite (v.2.27.1)

247　　(Quinlan & Hall, 2010) was applied to calculate the number of reads overlapping with any

248　　targeted CDS regions. Nucleotide identities of exons were calculated using *Muscle* (Edgar,

249　　2004) and BioPython's *Phylo* package (Talevich, Invergo, Cock, & Chapman, 2012). *R*

250　　(v.3.5.1) (R_Development_Core_Team, 2014) and the *ggplot2* package (Wickham, 2016)

251　　were utilized for summarizing and visualizing the data.

252　　　　　　　　　　　　　　　　　　Results

253　　 For the 56 mosquito gDNA samples, 4.9 million demultiplexed read pairs (150 bp PE) in

254　　total were obtained after a single run of Illumina MiniSeq. From those samples*,* 40–170, 50–

255　　120, and 63–100 thousand read-pairs were obtained from each individual mosquito of

256　　species *A. albopictus* and *A. aegypti* and *C. pipiens* complex, respectively. The raw read

257　　data were deposited to DDBJ Sequence Read Archive (DRA) under BioProjectID:

258　　PRJDB7889. In *A. albopictus*, 44% of all reads on average overlapped with any of the *VGSC*

259　　CDSs under study (Reads overlapping a per kilobase exon and per million sequenced reads:

260　　RPKM = 6.5 $\times$ 10$^4$), which was approximately 8.3 × 10$^4$-fold enrichment compared to

261　　simulated random data on whole-genome shotgun sequencing (Fig. 2). Although the probe

262　　set was designed based on the *A. albopictus* genomic sequence, the same probe set

263　　captured *VGSC* CDSs from the gDNA of *A. aegypti* and *C. pipiens* complex at an on-target

264　　rate similar to that of *A. albopictus* (Fig. 2).

265　　 Fig. 3 shows a distribution of median and minimum sequencing depths within each CDS

266　　in each individual sample after PCR duplicates were removed. In *A. albopictus*, most

267　　nucleotides in all exons were covered deeply with minimum bias in all samples. In *A. aegypti*

268　　and *C. quinquefasciatus*, however, some exons were covered at relatively low depth partly

269　　or entirely. In particular, exons 2 and 16.5 were covered at nearly or absolutely zero depth.

270　　 Fig. 4 presents a distribution of the allele balance in genotypes containing single or multiple

271　　nucleotide variants (SNVs or MNVs). The ratio of the first allele in a heterozygous genotype

272　　was distributed mostly around 50%, which was substantially different from the homozygous

273    genotype (near 100%) except for one SNV or MNV site in exon 32 of the *A. albopictus* gene

274    located in the GGT (Gly) trinucleotide tandem repeats variable in length near the C-terminus

275    of VGSC, where accurate calling of the genotype is difficult.

276    Aa substitutions identified at the end of the MoNaS pipeline are listed in Table 2. Of

277    those, F1534C in Aalb-SP (Shinji Kasai et al., 2019)*,* S989P and V1016G in Aaeg-SP (S

278    Kasai et al., 2014), and L1014F in Cpip-JPP (Hardstone et al., 2007), all previously known

279    to exist in those strains, were recalled correctly. In the Aaeg-Mex strain, V410L, V1016I,

280    and F1534C variants, which are known as *kdr* (Brengues et al. 2003; Haddi et al. 2017;

281    Kawada et al. 2009), were detected. Other aa substitutions detected—C749*Y (*in

282    mosquito aa coordinates because there was no corresponding aa in *M. domestica*),

283    A2023T and G2046E in *A. albopictus*, S723T in *A. aegypti*, and K109R, Y319F, T1632S,

284    E1633D, E1856D, G2051A, and A2055V in *C. pipiens* complex—are not known for their

285    effects on insecticide susceptibility.

## Discussion

287    In this study, we evaluated potential of targeted enrichment sequencing of *VGSC* CDSs

288    from the gDNA of mosquitos using hybridization capture probes. The result of the

289    experiment is quite promising: most nucleotides of *VGSC* CDSs were covered at sufficient

290    read depths even in samples with less than a 30 Mbp (0.1 million read-pairs) sequencing

291    effort.

292    Even though the probe set was designed on the basis of the *A. albopictus* genome

293    sequence only, it successfully enriched *VGSC* CDSs from the gDNA of two other Culicinae

294    mosquito species, *A. aegypti* and *C. pipiens* complex, which are estimated to have diverged

295    71.4 and 179 million years ago, respectively, from *A. albopictus* (X.-G. Chen et al., 2015).

296    Although the estimated entire efficiency of enrichment was lower by approximately 50% in

297    *C. pipiens* samples than in *A. albopictus* samples, the on-target ratio was still comparable

298    or rather higher in *C. pipiens* (Fig. 2). This result can be explained by the much smaller

299    genome size of *C. pipiens* complex (579 Mb in the CpipJ2 assembly) as compared to *A.*

300    *albopictus* (2.25 Gbp in the C6/36 assembly). Applicability of a single probe set to multi

301    species (e.g., the same genus or family) is clearly advantageous because this obviates the

11

302    need to prepare each custom probe sets specific for each single species and may enable

303    capture even in species lacking prior genome information. Nonetheless, the evolutionary

304    distance will limit the range of species that one probe set can be applied. In this study, the

305    mapping results on each exon indicated that capture efficiency decreased in some exons

306    (Fig. 3). The empirical observation suggests that less than 87.5% in identity or less than 60

307    bp in size for the homology track of targets could decrease the efficiency of capture

308    significantly (Fig. 5). Especially, our probe set failed to capture two optionally used exons, 2

309    and 16.5, in *A. aegypti* and *C. pipiens* complex, which are among the smallest exons

310    targeted (Fig. 3). It is assumed that those tiny exons alone do not provide enough

311    thermostability for probe–target DNA duplex during the capture. Because the probes for

312    those small exons contain flanking intronic sequences of the *A. albopictus* genome, those

313    flanking sequences may have provided enough homology region to capture sequences from

314    this species. Although it is straightforward to optimize our probe set further at least for the

315    two other species of mosquito simply by adding species specific probes for those exons and

316    flanking intronic regions, small exons in general will be a major challenge when a probe set

317    is aimed to be used for a group of species rather than specific targets because the homology

318    in an intronic region will decay more rapidly than that in an exonic region during speciation.

319    We also missed another exon corresponding to "exon 12" in *Anopheles gambiae* described

320    by Davies et al. (2007) during probe design (see Materials and Methods). Such tiny and

321    rarely used exons may be difficult to annotate without high-quality high-throughput RNA

322    sequencing data. Nevertheless, in mosquitoes, all such tiny exons are actually situated on

323    the N-terminal intracellular loop or the intracellular loop between domains I and domain II in

324    VGSC, where no resistance-associated mutation has been found so far (Dong et al., 2014).

325    Therefore, it is not clear whether ignoring those small exons of *VGSC* from analysis does

326    pose a serious problem for insecticide resistance research.

327     The process of genotyping *VGSC* carried out in this study was automated in MoNaS. This

328    program sequentially runs tools conducting mapping of NGS reads to a reference, sorting,

329    removal of PCR duplicates, indexing for BAM files, variant calling, variant annotation, and

330    finally integration of these results across multiple samples into a single table with conversion

331    of the aa coordinates to those corresponding to the *M. domestica* VGSC protein. The

332  automation in MoNaS allows researchers to process raw NGS reads of many samples via

333  a simple command line operation without expert knowledge of the bioinformatics field.

334  MoNaS can be run locally with appropriate genome reference data. Also, a web-service of

335  MoNaS implemented with JBrowse alignment viewer (Buels et al., 2016) is provided by NIID

336  Pathogen Genomics Center's severer (https://gph.niid.go.jp/monas) (Fig. 1B).

337

**Acknowledgements**

344

**Authors' contributions**

346  KI and SK designed the study. KI and OK designed the capture probes. KI and TT refined

347  the annotation of the *VGSC* gene. KI conducted all the experiments. KI, TSe, KY, and MK

348  contributed to the bioinformatic analysis and development of MoNaS. SK, YM, MS ans TSa

349  contributed maintaining the laboratory colony or collected specimens of mosquitoes in the

350  field. KI drafted this manuscript. All the coauthors critically revised the manuscript and

351  approved it for publication.

352

**Data accessibility**

354  Raw NGS reads obtained in this study were deposited to DDBJ Sequence Read Archive

355  (BioProjectID: PRJDB7889, see Table 1). Annotation information for *VGSC* and new

356  reference sequence of *VGSC* gene in *C. pipens* complex used in this study

357  (supercont3.182_3) are provided in S3 Appendix.zip file. Web service and source codes of

358    MoNaS are hosted in https://gph.niid.go.jp/monas and https://github.com/ItokawaK/MoNaS,

359    respectively.

360

361    Reference

362    Altschul, S. F., Gish, W., Miller, W., Myers, E. W., & Lipman, D. J. (1990). Basic local alignment

363        search tool. *Journal of Molecular Biology*, *215*(3), 403–410. doi:10.1016/S0022-

364        2836(05)80360-2

365    Arensburger, P., Megy, K., Waterhouse, R. M., Abrudan, J., Amedeo, P., Antelo, B., …

366        Atkinson, P. W. (2010). Sequencing of Culex quinquefasciatus establishes a platform for

367        mosquito comparative genomics. *Science*, *330*(6000), 86–88.

368        doi:10.1126/science.1191864

369    Brengues, C., Hawkes, N. J., Chandre, F., Mccarroll, L., Duchon, S., Guillet, P., … Hemingway,

370        J. (2003). Pyrethroid and DDT cross-resistance in Aedes aegypti is correlated with novel

371        mutations in the voltage-gated sodium channel gene. *Medical and Veterinary Entomology*,

372        *17*(1), 87–94. doi:10.1046/j.1365-2915.2003.00412.x

373    Buels, R., Yao, E., Diesh, C. M., Hayes, R. D., Munoz-Torres, M., Helt, G., … Holmes, I. H.

374        (2016). JBrowse: a dynamic web platform for genome visualization and analysis. *Genome*

375        *Biology*, *17*(1), 66. doi:10.1186/s13059-016-0924-1

376    Busvine, J. R. (1951). Mechanism of Resistance to Insecticide in Houseflies. *Nature*, *168*(4266),

377        193–195. doi:10.1038/168193a0

378    Chen, H., Li, K., Wang, X., Yang, X., Lin, Y., Cai, F., … Ma, Y. (2016). First identification of kdr

379        allele F1534S in VGSC gene and its association with resistance to pyrethroid insecticides

380        in Aedes albopictus populations from Haikou City, Hainan Island, China. *Infectious*

381        *Diseases of Poverty*, *5*(1), 31. doi:10.1186/s40249-016-0125-x

382    Chen, X.-G., Jiang, X., Gu, J., Xu, M., Wu, Y., Deng, Y., … James, A. A. (2015). Genome

383        sequence of the Asian Tiger mosquito, *Aedes albopictus* , reveals insights into its biology,

384        genetics, and evolution. *Proceedings of the National Academy of Sciences*, *112*(44),

385        E5907–E5915. doi:10.1073/pnas.1516410112

386    Clarkson, C. S., Miles, A., Harding, N. J., Weetman, D., Kwiatkowski, D., Donnelly, M., &

14

387      Consortium, T. &lt;I&gt;Anopheles gambiae&lt;/I&gt; 1000 G. (2018). The genetic

388        architecture of target-site resistance to pyrethroid insecticides in the African malaria

389        vectors Anopheles gambiae and Anopheles coluzzii. *BioRxiv*, 323980. doi:10.1101/323980

390    Danecek, P., & McCarthy, S. A. (2017). BCFtools/csq: haplotype-aware variant consequences.

391        *Bioinformatics*, *33*(13), 2037–2039. doi:10.1093/bioinformatics/btx100

392    Davies, T. G. E., Field, L. M., Usherwood, P. N. R., & Williamson, M. S. (2007). A comparative

393        study of voltage-gated sodium channels in the Insecta: implications for pyrethroid

394        resistance in Anopheline and other Neopteran species. *Insect Molecular Biology*, *16*(3),

395        361–75. doi:10.1111/j.1365-2583.2007.00733.x

396    Dong, K., Du, Y., Rinkevich, F., Nomura, Y., Xu, P., Wang, L., … Zhorov, B. S. (2014, July).

397        Molecular biology of insect sodium channels and pyrethroid resistance. *Insect*

398        *Biochemistry and Molecular Biology*. NIH Public Access. doi:10.1016/j.ibmb.2014.03.012

399    Edgar, R. C. (2004). MUSCLE: multiple sequence alignment with high accuracy and high

400        throughput. *Nucleic Acids Res*, *32*(5), 1792–1797. doi:10.1093/nar/gkh340

401    Garrison, E., & Marth, G. (2012). Haplotype-based variant detection from short-read

402        sequencing. Retrieved from http://arxiv.org/abs/1207.3907

403    Giraldo-Calderón, G. I., Emrich, S. J., MacCallum, R. M., Maslen, G., Dialynas, E., Topalis, P.,

404        … Lawson, D. (2015). VectorBase: an updated bioinformatics resource for invertebrate

405        vectors and other organisms related with human diseases. *Nucleic Acids Research*,

406        *43*(D1), D707–D713. doi:10.1093/nar/gku1117

407    Gnirke, A., Melnikov, A., Maguire, J., Rogov, P., LeProust, E. M., Brockman, W., … Nusbaum,

408        C. (2009). Solution hybrid selection with ultra-long oligonucleotides for massively parallel

409        targeted sequencing. *Nature Biotechnology*, *27*(2), 182–189. doi:10.1038/nbt.1523

410    Haddi, K., Tomé, H. V. V., Du, Y., Valbon, W. R., Nomura, Y., Martins, G. F., … Oliveira, E. E.

411        (2017). Detection of a new pyrethroid resistance mutation (V410L) in the sodium channel

412        of Aedes aegypti: a potential challenge for mosquito control. *Scientific Reports*, *7*(1),

413        46549. doi:10.1038/srep46549

414    Hardstone, M., Leichter, C., Harrington, L., Kasai, S., Tomita, T., & Scott, J. (2007). Cytochrome

415        P450 monooxygenase-mediated permethrin resistance confers limited and larval specific

416        cross-resistance in the southern house mosquito, Culex pipiens quinquefasciatus.

15

417      *Pesticide Biochemistry and Physiology*, *89*(3), 175–184. doi:10.1016/j.pestbp.2007.06.006

418    Kasai, S, Komagata, O., Itokawa, K., Shono, T., Ng, L. C., Kobayashi, M., & Tomita, T. (2014).

419      Mechanisms of pyrethroid resistance in the dengue mosquito vector, Aedes aegypti: target

420      site insensitivity, penetration, and metabolism. *PLoS Negl Trop Dis*, *8*(6), e2948.

421      doi:10.1371/journal.pntd.0002948

422    Kasai, S, Ng, L. C., Lam-Phua, S. G., Tang, C. S., Itokawa, K., Komagata, O., … Tomita, T.

423      (2011). First detection of a putative knockdown resistance gene in major mosquito vector,

424      Aedes albopictus. *Jpn J Infect Dis*, *64*(3), 217–221. Retrieved from

425      http://www.ncbi.nlm.nih.gov/pubmed/21617306

426    Kasai, Shinji, Caputo, B., Tsunoda, T., Cuong, T. C., Maekawa, Y., Lam-Phua, S. G., … Tomita,

427      T. (2019). First detection of a Vssc allele V1016G conferring a high level of insecticide

428      resistance in Aedes albopictus collected from Europe (Italy) and Asia (Vietnam), 2016: a

429      new emerging threat to controlling arboviral diseases. *Eurosurveillance*, *24*(5), 1700847.

430      doi:10.2807/1560-7917.ES.2019.24.5.1700847

431    Kawada, H., Higa, Y., Komagata, O., Kasai, S., Tomita, T., Thi Yen, N., … Takagi, M. (2009).

432      Widespread Distribution of a Newly Found Point Mutation in Voltage-Gated Sodium

433      Channel in Pyrethroid-Resistant Aedes aegypti Populations in Vietnam. *PLoS Negl Trop*

434      *Dis*, *3*(10), e527. doi:10.1371/journal.pntd.0000527

435    Kim, D., Langmead, B., & Salzberg, S. L. (2015). HISAT: A fast spliced aligner with low memory

436      requirements. *Nature Methods*. doi:10.1038/nmeth.3317

437    Kutsuna, S., Kato, Y., Moi, M. L., Kotaki, A., Ota, M., Shinohara, K., … Ohmagari, N. (2015).

438      Autochthonous Dengue Fever, Tokyo, Japan, 2014. *Emerging Infectious Diseases*, *21*(3),

439      517–520. doi:10.3201/eid2103.141662

440    Li, H., & Durbin, R. (2009). Fast and accurate short read alignment with Burrows-Wheeler

441      transform. *Bioinformatics*, *25*(14), 1754–1760. doi:10.1093/bioinformatics/btp324

442    Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., … Durbin, R. (2009). The

443      Sequence Alignment/Map format and SAMtools. *Bioinformatics*, *25*(16), 2078–2079.

444      doi:10.1093/bioinformatics/btp352

445    Lund, A. E., & Narahashi, T. (1983). Kinetics of sodium channel modification as the basis for the

446      variation in the nerve membrane effects of pyrethroids and DDT analogs. *Pesticide*

16

447       *Biochemistry and Physiology*, *20*(2), 203–216. doi:10.1016/0048-3575(83)90025-1

448    Marcombe, S., Farajollahi, A., Healy, S. P., Clark, G. G., & Fonseca, D. M. (2014). Insecticide

449       Resistance Status of United States Populations of Aedes albopictus and Mechanisms

450       Involved. *PLoS ONE*, *9*(7), e101992. doi:10.1371/journal.pone.0101992

451    Matthews, B. J., Dudchenko, O., Kingan, S. B., Koren, S., Antoshechkin, I., Crawford, J. E., …

452       Vosshall, L. B. (2018). Improved reference genome of Aedes aegypti informs arbovirus

453       vector control. *Nature*, *563*(7732), 501–507. doi:10.1038/s41586-018-0692-z

454    Miller, J. R., Koren, S., Dilley, K. A., Puri, V., Brown, D. M., Harkins, D. M., … Shabman, R. S.

455       (2018). Analysis of the Aedes albopictus C6/36 genome provides insight into cell line utility

456       for viral propagation. *GigaScience*, *7*(3). doi:10.1093/gigascience/gix135

457    Miyazaki, M., Ohyama, K., Dunlap, D. Y., & Matsumura, F. (1996). Cloning and sequencing of

458       thepara-type sodium channel gene from susceptible andkdr-resistant German

459       cockroaches (Blattella germanica) and house fly (Musca domestica). *MGG Molecular &*

460       *General Genetics*, *252*(1–2), 61–68. doi:10.1007/BF02173205

461    Quinlan, A. R., & Hall, I. M. (2010). BEDTools: a flexible suite of utilities for comparing genomic

462       features. *Bioinformatics*, *26*(6), 841–842. doi:10.1093/bioinformatics/btq033

463    R_Development_Core_Team. (2014). R: A language and environment for statistical computing.

464       *R Foundation for Statistical Computing, Vienna, Austria.* Retrieved from http://www.r-

465       project.org

466    Reiter, P., Fontenille, D., & Paupy, C. (2006). Aedes albopictus as an epidemic vector of

467       chikungunya virus: another emerging problem? *The Lancet Infectious Diseases*, *6*(8),

468       463–464. doi:10.1016/S1473-3099(06)70531-X

469    Rinkevich, F. D., Du, Y., & Dong, K. (2013). Diversity and convergence of sodium channel

470       mutations involved in resistance to pyrethroids. *Pesticide Biochemistry and Physiology*,

471       *106*(3), 93–100. doi:10.1016/j.pestbp.2013.02.007

472    Saavedra-Rodriguez, K., Urdaneta-Marquez, L., Rajatileka, S., Moulton, M., Flores, A. E.,

473       Fernandez-Salas, I., … Black, W. C. th. (2007). A mutation in the voltage-gated sodium

474       channel gene associated with pyrethroid resistance in Latin American Aedes aegypti.

475       *Insect Mol Biol*, *16*(6), 785–798. doi:IMB774 [pii]10.1111/j.1365-2583.2007.00774.x

476    Srisawat, R., Komalamisra, N., Eshita, Y., Zheng, M., Ono, K., Itoh, T. Q., … Rongsriyam, Y.

477        (2010). Point mutations in domain II of the voltage-gated sodium channel gene in

478        deltamethrin-resistant Aedes aegypti (Diptera: Culicidae). *Applied Entomology and*

479        *Zoology*, *45*(2), 275–282. doi:10.1303/aez.2010.275

480    Talevich, E., Invergo, B. M., Cock, P. J., & Chapman, B. A. (2012). Bio.Phylo: A unified toolkit

481        for processing, analyzing and visualizing phylogenetic trees in Biopython. *BMC*

482        *Bioinformatics*, *13*(1), 209. doi:10.1186/1471-2105-13-209

483    Wickham, H. (2016). *ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York.

484        Retrieved from http://ggplot2.org

485    Williamson, M. S., Martinez-Torres, D., Hick, C. A., & Devonshire, A. L. (1996). Identification of

486        mutations in the houseflypara-type sodium channel gene associated with knockdown

487        resistance (kdr) to pyrethroid insecticides. *MGG Molecular & General Genetics*, *252*(1–2),

488        51–60. doi:10.1007/BF02173204

489    Xu, J., Bonizzoni, M., Zhong, D., Zhou, G., Cai, S., Li, Y., … Chen, X.-G. (2016). Multi-country

490        Survey Revealed Prevalent and Novel F1534S Mutation in Voltage-Gated Sodium

491        Channel (VGSC) Gene in Aedes albopictus. *PLOS Neglected Tropical Diseases*, *10*(5),

492        e0004696. doi:10.1371/journal.pntd.0004696

493

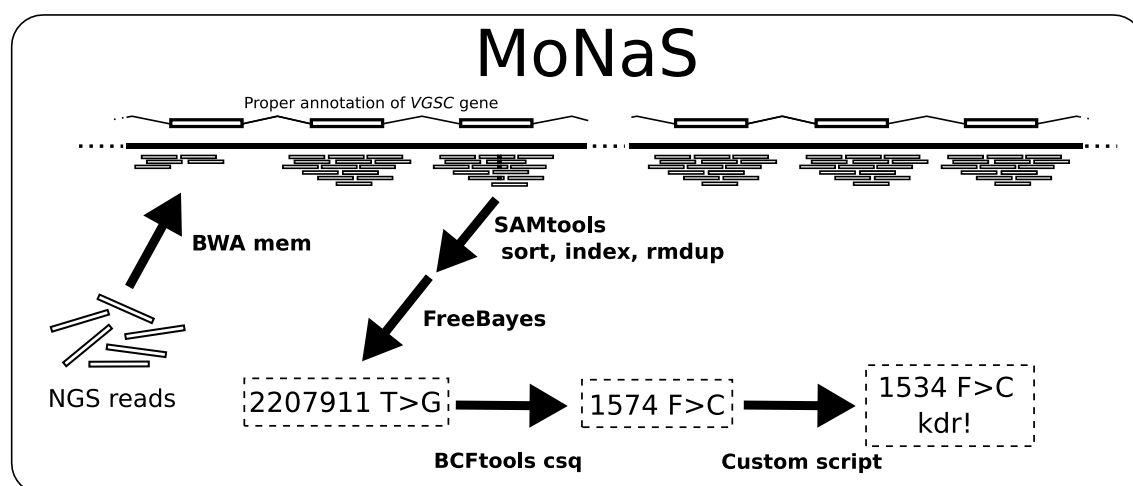494

Table 1 Mosquito samples used in this study

| ID | Species | Lab. Colony / Wild | num | Description | DDBJ Accession no. |
|---|---|---|---|---|---|
| Aalb-SP | *A. albopictus* | Lab. Colony | 8 | Originated from Singapore in 2016. This strain is known to possess 1534C *kdr* variant (Kasai et al., 2019). | DRR167925–932 |
| Aalb-Viet | *A. albopictus* | Lab. Colony | 2 | Originated from Hanoi, Viet Nam in 2016. | DRR167935–936 |
| Aalb-Okayama | *A. albopictus* | Lab. Colony | 2 | Originated from Okayama, Japan in 2015. | DRR167923–924 |
| Aalb-Toyama | *A. albopictus* | Lab. Colony | 2 | Originated from Tokyo, Japan in 2015. | DRR167933–934 |
| Aalb-Ishigaki | *A. albopictus* | Lab. Colony | 2 | Originated from Ishigaki-zima, Okinawa, Japan in 2016. | DRR167921–922 |
| Aalb-Yona | *A. albopictus* | Wild | 8 | Caught wild from Yonaguni-zima, Okinawa, Japan in 2017. | DRR167937–944 |
| Aaeg-Mex | *A. aegypti* | Lab. Colony | 16 | Originated from Monterrey, Mexico in 2008. | DRR167897–912 |
| Aaeg-SP | *A. aegypti* | Lab. Colony | 8 | Originated from Singapore in 2009. This strain is known to possesses 989P–1016G kdr haplotype (Kasai 2014). | DRR167913–920 |
| Cpip-JNA | *C. quinquefasciatus* | Lab. Colony | 2 | This strain was selected for *CYP9M10* genotype (Itokawa et al., 2011) from JHB strain originated from Johannesburg, South Africa in 2001 (Arensburger et al., 2010). | DRR167945–946 |
| Cpip-JPP | *C. quinquefasciatus* | Lab. Colony | 2 | Originated from Saudi Arabia, selected by permethrin for 20 generations (Amin et al., 1989). This strain is known to possesses 1014F *kdr* variant (Hardstone et al., 2007) . | DRR167947–948 |
| Cpip-Ryo | *C. pipiens pallens* | Lab. Colony | 2 | Originated from Kanagawa, Japan in 2015. | DRR167951–952 |
| Cpip-JP_mix | *C. pipiens* form *molestus* | Lab. Colony | 2 | Mixed from several lab. colonies originated from different places of Japan in 2003–2004. | DRR167949–950 |

Table 2 Detected amino-acid substitutions

| Species | Population | n | Amino-acid substitutions (num. of homozygous, heterozygous individuals) |
|---------|-----------|---|------------------------------------------------------------------------|
| *A. albopictus* | Aalb-SP | 8 | F1534C**(8,0); A2023T(8,0) |
| | Aalb-Viet | 2 | A2023T(0,1) |
| | Aalb-Okayama | 2 | C749*Y(0,1); A2023T(0,1); G2046E(0,1) |
| | Aalb-Toyama | 2 | A2023T(0,1) |
| | Aalb-Ishigaki | 2 | A2023T(1,0) |
| | Aalb-Yona | 8 | A2023T(0,2) |
| *A. aegypti* | Aaeg-Mex | 16 | V410L**(1,7); S723T(1,7); V1016I**(1,7); F1534C**(16,0) |
| | Aaeg-SP | 8 | S989P**(8,0); V1016G**(8,0) |
| *C. quinquefasciatus* | Cpip-JNA | 2 | K109R(0,1); T1632S(0,1); E1633D(0,1); G2051A(0,1); A2055V(0,1) |
| *C. quinquefasciatus* | Cpip-JPP | 2 | R261K(2,0); L1014F**(2,0) |
| *C. pipiens pallens* | Cpip-Ryo | 2 | Y319F(2,0); T1632S(0,2); E1633D(0,2) |
| *C. pipiens* form *molestus* | Cpip-JP_mix | 2 | Y319F(2,0); L1014F**(2,0); T1632S(2,0); T1633D(2,0); E1856D(2,0) |

The amino-acid coordination is *M. domestica* except C749Y with asterisk (*) since there was no corresponding amino-acid in *M. domestica* (genbank id: AAB47604). Double asterisks (**) indicate known *kdr* substitutions conferring pyrethroid resistance. Variants seen on the Gly repeats near the C-terminal are omitted.
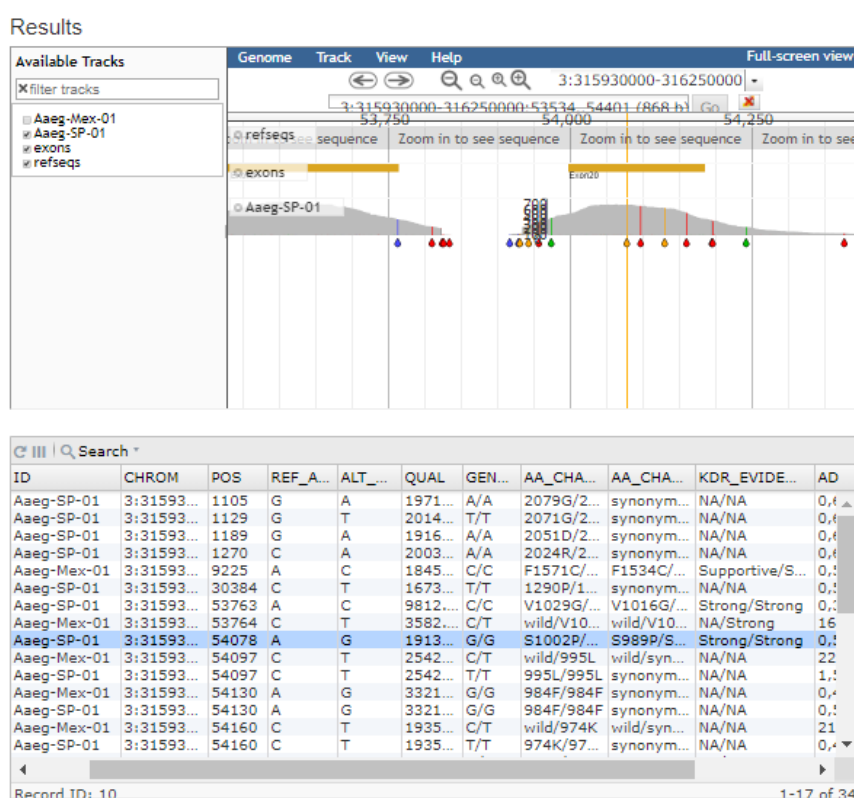
(A)



(B)



Fig. 1 Analytical pipeline MoNaS

(A) A diagram for analytical pipeline MoNaS. MoNaS executes several bioinformatic tools to call variants and aa substitutions. Finally, a custom script converts species aa coordinates to the standard housefly aa coordinates, tells whether each aa substitution is among the known listed *kdr* substitutions and creates a human-readable table from Variant Call Format (VCF). (B) Image of the result output page from MoNaS web-service.
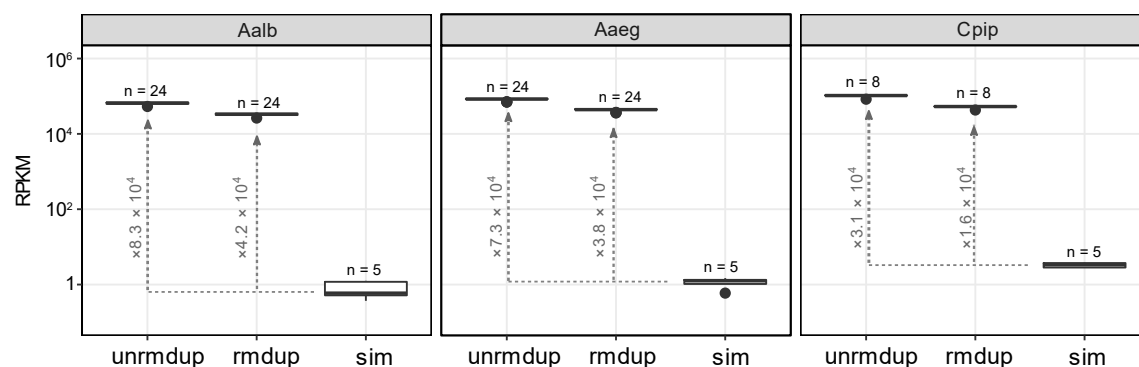
Fig. 2 The NGS reads are enriched in targeted *VGSC* exons by capture

Distributions of RPKM (number of sequencing reads overlapping to the targeted *VGSC* exons per 1 kbp total exon length and one million reads) for *A. albopictus* (Aalb), *A. aegypti* (Aaeg) and *C. pipiens* complex (Cpip). Labels "unrmdup" and "rmdup" indicate before and after removal of PCR duplicates, respectively. Label "sim" indicates simulated whole genome shotgun (WGS) reads randomly drawn from genome of each species (replicated five times in each species). The values associated with dotted line with arrowhead indicate sizes of fold-change (levels of enrichment) compared to simulated WGS data.
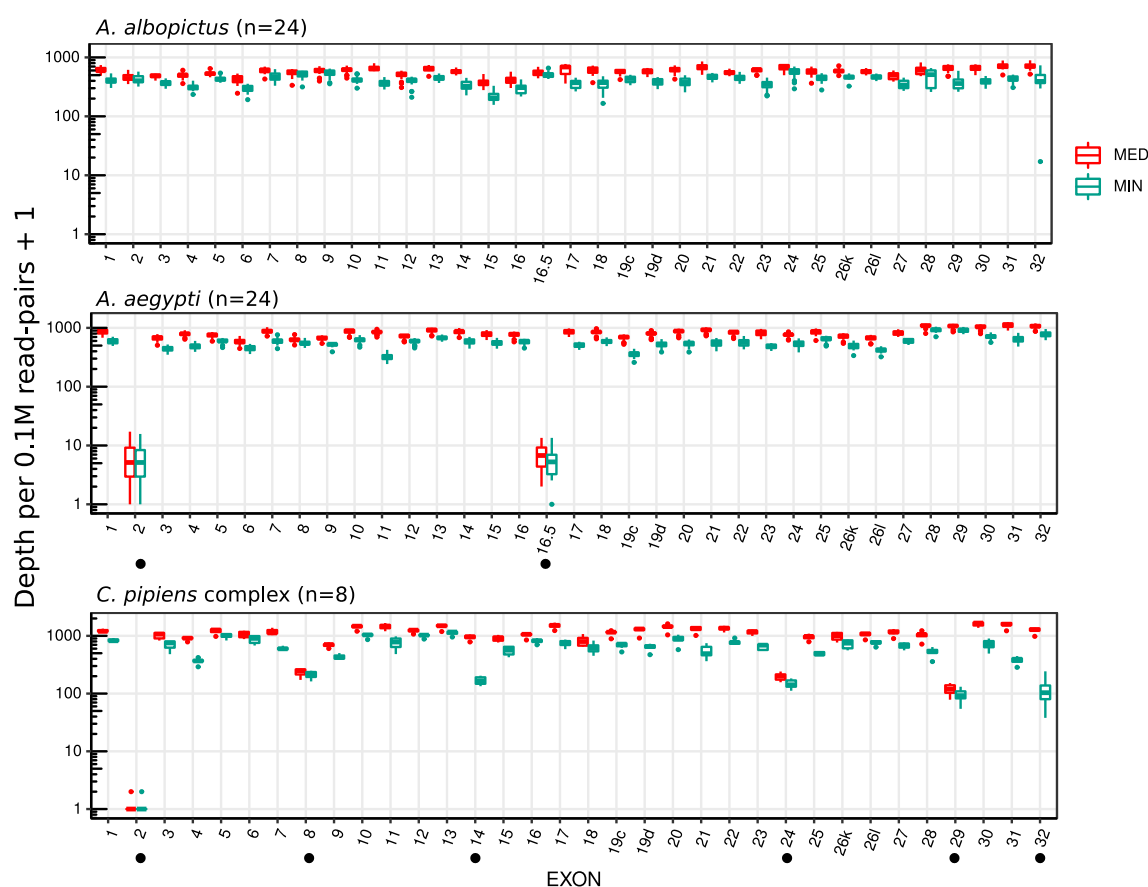
**Fig. 3 Coverage of targeted *VGSC* exons**

Distribution of median (MED) and minimum (MIN) depths per nucleotide (on a logarithmic scale) within each exon and each individual sample after PCR duplicates were removed. Exons labeled with "●" contained nucleotide sites with relatively low coverage.
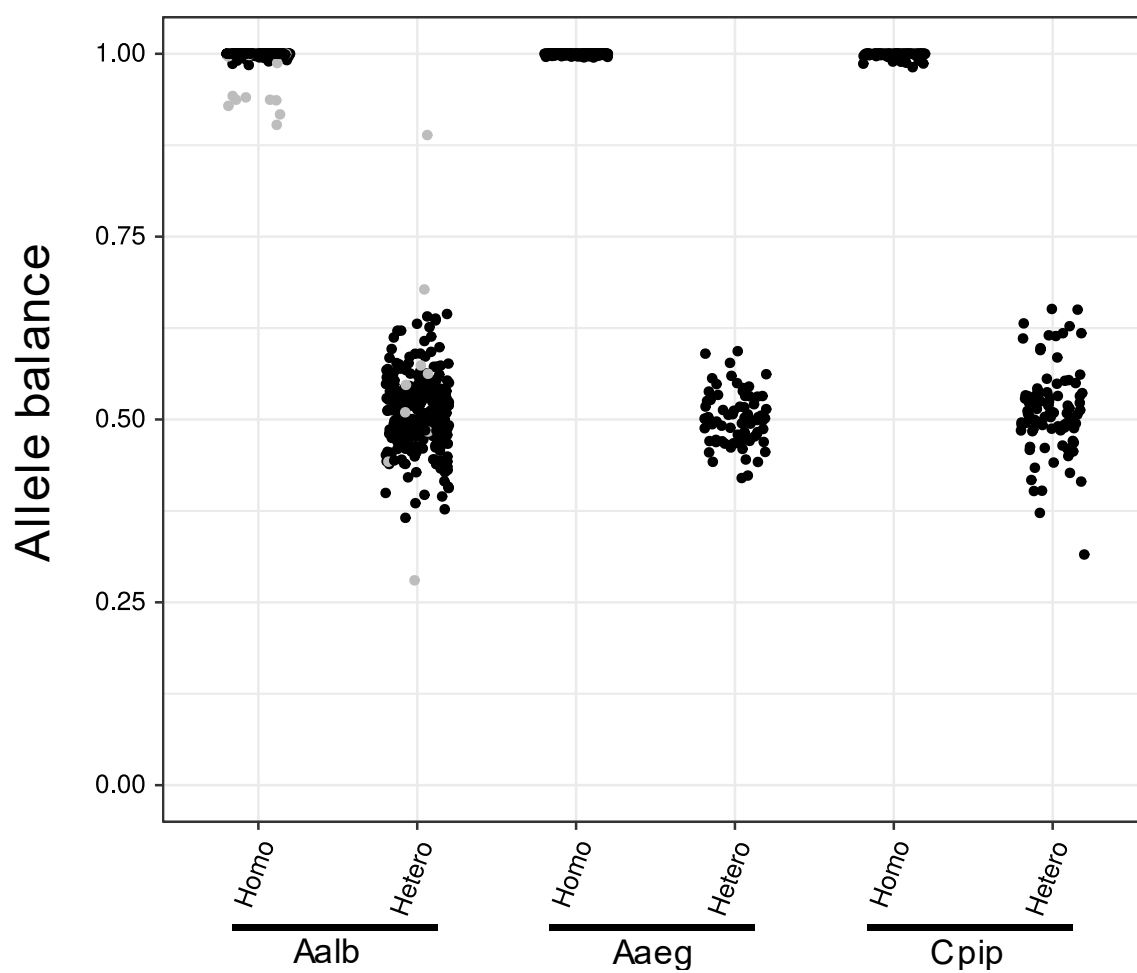
Fig. 4 The allele balance in genotypes containing a variant

The distribution of allele balance in read depth for each allele at heterozygous or homozygous genotypes containing alternative allele(s). The balance was calculated as [read depth of the first allele in "GT" info] / [total depth] in the VCF format. Gray points in *A. albopictus* are at the genotype at GGT (Gly) trinucleotide repeats in exon32.
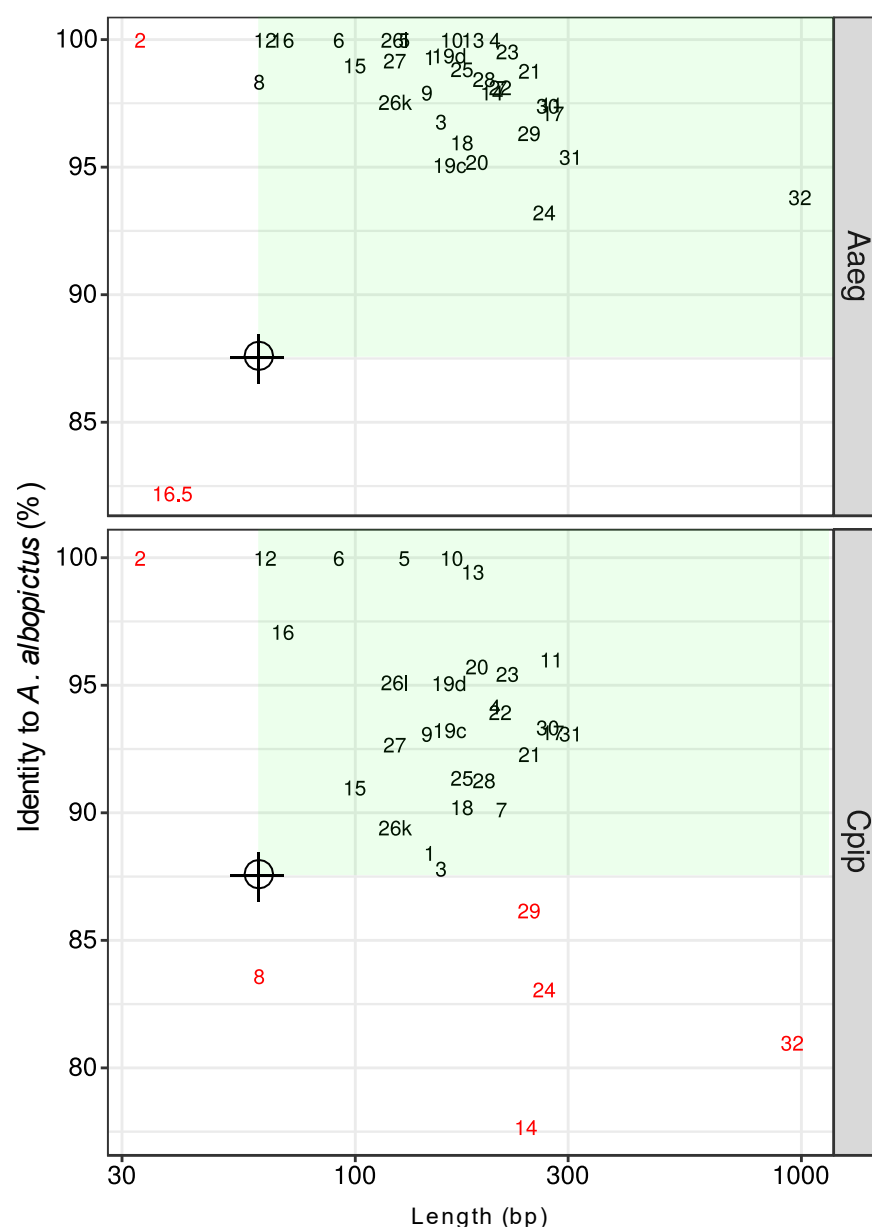
Fig. 5 Length and conservation of the *VGSC* exons (CDSs) targeted

Length (on a logarithmic scale) and percentage identity to *A. albopictus* of each exon in *A. aegypti* and *C. pipiens* complex. Red exon names are those with low-coverage nucleotides in Fig. 3. The green area represents >60 bp length and >87.5% identity.