



Minerva Access is the Institutional Repository of The University of Melbourne

Author/s:

Zhang, J;Zhang, F;Tay, WT;Robin, C;Shi, Y;Guan, F;Yang, Y;Wu, Y

Title:

Population genomics provides insights into lineage divergence and local adaptation within the cotton bollworm

Date:

2022-07-01

Citation:

Zhang, J., Zhang, F., Tay, W. T., Robin, C., Shi, Y., Guan, F., Yang, Y. & Wu, Y. (2022). Population genomics provides insights into lineage divergence and local adaptation within the cotton bollworm. *Molecular Ecology Resources*, 22 (5), pp.1875-1891. <https://doi.org/10.1111/1755-0998.13581>.

Persistent Link:

<https://hdl.handle.net/11343/336305>

1

2 DR. WEE TEK TAY (Orcid ID : 0000-0002-8451-0811)

3 DR. YIDONG WU (Orcid ID : 0000-0003-3456-3373)

4

5

6 Article type : Resource Article

7

8

9 **In preparation for:** Molecular Ecology Resources10 **Running title:** Population genomics of *Helicoverpa armigera*

11

12 **Population genomics provides insights into lineage divergence and local**  
13 **adaptation within the cotton bollworm**14 Jianpeng Zhang<sup>1</sup>, Feng Zhang<sup>1</sup>, Wee Tek Tay<sup>2</sup>, Charles Robin<sup>3</sup>, Yu Shi<sup>1</sup>, Fang Guan<sup>1</sup>, Yihua Yang<sup>1</sup>,  
15 Yidong Wu<sup>1,\*</sup>

16

17 <sup>1</sup> College of Plant Protection, Nanjing Agricultural University, Nanjing 210095, China18 <sup>2</sup> CSIRO Black Mountain Laboratories, Clunies Ross Street, ACT 2601, Australia19 <sup>3</sup> School of BioSciences, University of Melbourne, Parkville, VIC 3010, Australia.

20

21 \* Corresponding author. E-mail: wyd@njau.edu.cn

22 **ABSTRACT**23 The cotton bollworm *Helicoverpa armigera* is a cosmopolitan pest and its diverse habitats plausibly  
24 contribute to the formation of diverse lineages. Despite the significant threat it poses to economic  
25 crops worldwide, its evolutionary history and genetic basis of local adaptation are poorly

This is the author manuscript accepted for publication and has undergone full peer review but has not been through the copyediting, typesetting, pagination and proofreading process, which may lead to differences between this version and the [Version of Record](#). Please cite this article as [doi: 10.1111/1755-0998.13581](https://doi.org/10.1111/1755-0998.13581)

This article is protected by copyright. All rights reserved

26 understood. In this study, we *de novo* assembled a high-quality chromosome-level reference  
27 genome of *H. a. armigera* (contig N50 = 7.34 Mb), with 99.13% of the HaSCD2 assembly assigned  
28 into 31 chromosomes (Z-chromosome + 30 autosomes). We constructed an ultra-dense variation  
29 map across 14 cotton bollworm populations and identified a novel lineage in northwestern China.  
30 Historical inference showed that effective population size changes coincided with global  
31 temperature fluctuation. We identified nine differentiated genes in the three *H. armigera* lineages  
32 (*H. a. armigera*, *H. a. conferta*, and the new northwestern Chinese lineage), of which *per* and *clk*  
33 genes are involved in circadian rhythm. Selective sweep analyses identified a series of GO  
34 categories related to climate adaptation, feeding behavior and insecticide tolerance. Our findings  
35 reveal fundamental knowledge of the local adaptation of different cotton bollworm lineages and  
36 will guide the formulation of cotton bollworm management measures at different scales.

37 **Keywords:** cotton bollworm, population genomics, genetic structure, local adaptation

## 38 INTRODUCTION

39 The cotton bollworm, *Helicoverpa armigera* (Lepidoptera: Noctuidae), is a devastating pest. It can  
40 feed on a broad range of agricultural and horticultural plants, including cotton, soybean, maize and  
41 sunflower (Fitt, 1989). It is capable of long-range migration, facultative diapause and high  
42 fecundity, and these features enable the cotton bollworm to successfully colonize and flourish in  
43 various ecological niches, and therefore, it can be seen as a model species for studying local  
44 adaptation in Lepidoptera (Jones, Parry, Tay, Reynolds, & Chapman, 2019; Wu & Guo, 2005).  
45 Until its recent invasion into the Americas (Tay et al., 2013), the cotton bollworm was naturally  
46 distributed throughout Asia, Europe and Africa (known as *H. a. armigera*, *sensu* Hardwick, 1965)  
47 as well as Oceania where the subspecies, *H. a. conferta* is found. A third subspecies confined to the  
48 Canton Island (Republic of Kiribati) in the South Pacific, *H. a. commoni*, was also proposed by  
49 Hardwick (1965) based on body size and coloration differences.

50 Over the last decade, numerous investigations on the cotton bollworm radically reinvigorated our  
51 understanding of the population genomics of this major agricultural pest (Jones et al., 2019; Van  
52 Leeuwen, Dermauw, Mavridis, & Vontas, 2020). In particular, the previous short-read  
53 whole-genome sequencing efforts generated the draft genomes of *H. a. conferta* and its sister  
54 species *H. zea* (Pearce et al., 2017), which greatly promoted our understanding of insecticide

This article is protected by copyright. All rights reserved

55 resistance, biotic incursion, population admixture, and insect-plant interactions (Anderson et al.,  
56 2016; Guo et al., 2020; Pearce et al., 2017; Valencia-Montoya et al., 2020). However, the inherent  
57 deficiencies, e.g., fragmentary sequences and numerous N-gap regions frustrate efforts to fully  
58 characterize genotype-phenotype relationships. Meanwhile, despite of their morphological  
59 similarity, genomic data has supported Hardwick's (1965) taxonomic classification separating *H. a.*  
60 *armigera* and *H. a. conferta* into two allopatric and genetically distinct subspecies (Anderson, Tay,  
61 McGaughran, Gordon, & Walsh, 2016). Thus, it is necessary to account for population structure and  
62 not rely upon genomic data sourced predominantly from a single subspecies to build a  
63 comprehensive management framework. A perfect illustration of this is the desire to understand and  
64 learn from the recent establishment of expanding *H. a. armigera* populations in the South American  
65 continent (Arnemann et al., 2019; Tay et al., 2013, 2017) and the Caribbean Islands (Jones et al.,  
66 2019). Furthermore, recent methodological and technological advances in third generation  
67 sequencing (TGS) have greatly improved the reference genomes of different agricultural pests (Fritz  
68 et al., 2018; Quan et al., 2019; Zhu et al., 2017). Compared to next generation sequencing (NGS)  
69 technology, TGS is able to generate longer reads which can span many repetitive regions and  
70 provide a more comprehensive genomic portrait (Sedlazeck, Lee, Darby, & Schatz, 2018).

71 Elucidating the dynamics, composition and relationships of different populations has practical  
72 benefits for pest management. In particular, unrecognized adaptive introgression may render local  
73 pest management unsuccessful. For example in 2013, it was confirmed that *H. a. armigera* had  
74 invaded Brazil (Czepak, Albernaz, Vivan, Guimarães, & Carvalhais, 2013; Tay et al., 2013).  
75 Invaded individuals were subsequently shown to hybridize with the previously isolated sister  
76 species *H. zea*, resulting in the introgression of an allele conferring resistance to fenvalerate and  
77 cypermethrin, *CYP337B3* (Anderson et al., 2018; Joußen et al., 2012; Valencia-Montoya et al.,  
78 2020; Walsh et al. 2018). Other insecticide resistance alleles and loci that counter diverse  
79 insecticides including *Bacillus thuringiensis* (Bt) toxins have been reported in different cotton  
80 bollworm populations and could also spread through populations, subspecies or species hybrids  
81 (Tabashnik & Carrière, 2017). For example, cotton bollworm in northern China has evolved  
82 high-level resistances to synthetic pyrethroid and organophosphate insecticides since the 1990s,  
83 whereas populations in northwestern China remained sensitive as recently as 2013 (Yang, Li, &  
84 Wu, 2013). A dominant point mutation conferring dominant resistance to Cry1Ac was identified in

85 northern China populations but it has not been detected in other regions (Jin et al., 2018). With  
86 globalization of trade and human movement, regions are less isolated and pests are more likely to  
87 spread (Tay & Gordon, 2019). Consequently, the invasion risks of noctuid pests have greatly  
88 elevated despite of strict quarantine measures worldwide.

89 Here, we combined single-molecule real-time (SMRT) sequencing and genetic map to assemble a  
90 chromosome-level genome of *H. a. armigera*. We re-sequenced the whole genomes of 141 cotton  
91 bollworm from China and downloaded 30 *Helicoverpa* complex data. By the population genomics  
92 approach, we explored the population structure and demographic history within the cotton  
93 bollworm. We also identified various highly divergent genomic regions that suggested climate, host  
94 plants and insecticides have contributed to shaping the cotton bollworm populations. Our research  
95 provides novel insights into adaptation by the cotton bollworm as well as valuable genomic  
96 resources for management of this major agricultural pest.

## 97 **MATERIALS AND METHODS**

### 98 **SCD strain and sequencing for genome assembly**

99 The SCD strain (*H. a. armigera*), used for genome assembly, was originally sampled from Africa in  
100 the 1970s and continuously reared on artificial diets in our laboratory (Nanjing, China) without  
101 exposure to any chemical pesticides or Bt toxins (Jin et al., 2018). For this study, the SCD strain  
102 was bred through three generations of successive sib-crossing to reduce its heterozygosity.  
103 Subsequently, we selected a single SCD male (ZZ) *H. a. armigera* pupa for total genomic DNA  
104 isolation using the phenol/chloroform method. The purity and concentration of the extracted DNA  
105 were measured using OneDrop™ 1000 Spectrophotometer (OneDrop Technologies). The extracted  
106 DNA was used to construct Illumina library with 350-bp insert size and Pacific Biosciences  
107 (PacBio) library with 20-kb insert size following their respective standard protocols. The short-read  
108 library was sequenced on Illumina HiSeq Xten platform using the 150-bp paired-end (PE) strategy,  
109 and seven single-molecule real-time (SMRT) long-read libraries were sequenced on PacBio Sequel  
110 platform.

### 111 **Genome survey and *de novo* genome assembly**

112 We summarized the main processes and their corresponding programs throughout genome assembly

This article is protected by copyright. All rights reserved

113 as illustrated in Figure S1. Qualified Illumina data (Phred score > 15) were used to count *k*-mers (*k*  
114 = 17) by khist.sh built-in BBTools v38.49 (<https://doi.org/jgi.doe.gov/data-and-tools/bbtools/>), then  
115 genome size and heterozygosity rate were estimated using GenomeScope v1.0.0 (Vurture et al.,  
116 2017). The subreads from seven SMRT cells were merged and converted to FASTA format using  
117 SAMtools v0.1.19 (Li et al., 2009). The Canu v1.8 (Koren et al., 2017) and Falcon v1.3.0 (Chin et  
118 al., 2016) were applied to *de novo* assemble draft assemblies, respectively. The advanced options  
119 “batOptions=-dg 3 -db 3 -dr 1 -ca 500 -cp 50” were selected in Canu to maintain haplotype  
120 separation. These two draft assemblies had complementary contiguity, we therefore merged them  
121 using quickmerge v0.3 (Chakraborty, Baldwin-Brown, Long, & Emerson, 2016) to salvage missing  
122 and/or wrong segments. Alternative haplotypes were removed from the merged assembly using  
123 Purge Haplotigs v1.1.0 (Roach, Schmidt, & Borneman, 2018). Non-redundant sequences were  
124 polished for two rounds with short PE reads using Pilon v1.23 (Walker et al., 2014). The vector and  
125 microbe contaminations were searched using HS-BLASTN v1.0.0 (Chen, Ye, Zhang, & Xu, 2015)  
126 against the Nt and UniVec databases, with the alignment threshold set as 1e-5. Finally, we  
127 employed SeqKit v0.14.0 (Shen, Le, Li, & Hu, 2016) to remove identified contaminations (i.e.,  
128 vector and microbe) from the HaSCD2 assembly.

### 129 **Genetic map construction**

130 For genetic map construction, the SCD individual was crossed with an AY individual (collected in  
131 2011 from Anyang city, China) as illustrated in Figure S2. Genomic DNA of 100 BC1 offspring  
132 and their parents was isolated to prepare short-read libraries, and these libraries were sequenced  
133 using the Illumina HiSeq Xten platform. The short reads were mapped to the assembled HaSCD2  
134 genome using BWA-mem v0.7.12, with ‘-M’ parameters (Li & Durbin, 2010). SAMtools v0.1.19  
135 was used to format intermediate files (Li et al., 2009). Genotype and variants calling were  
136 conducted by GATK v3.8 (McKenna et al., 2010), and hard filtering was executed with the  
137 following criteria: QD < 2.0; FS > 60.0; MQ < 30.0; HaplotypeScore > 13.0; MQRankSum < -12.5;  
138 ReadPosRankSum < -8.0. High-quality SNP markers were used to construct the genetic map using  
139 Lep-map3 (Rastas, 2017). Female-informative markers (i.e., heterozygous alleles in the mother and  
140 homozygous in the father) were selected to determine linkage groups (LGs) with LOD score of 15.  
141 Subsequently, male-informative markers (i.e., homozygous alleles in the mother and heterozygous  
142 in the father) were sorted and used to compute genetic distance of pairwise markers. Finally, the

This article is protected by copyright. All rights reserved

143 contigs of HaSCD2 assembly were clustered, ordered and orientated based on the genetic map using  
144 ALLMAPS (Tang et al., 2015).

### 145 **Quality assessment of the HaSCD2 assembly**

146 We evaluated the quality of the HaSCD2 assembly in three aspects, i.e., completeness, continuity,  
147 and accuracy. First, we assessed the completeness using Benchmarking Universal Single-Copy  
148 Orthologs v3.0.2 (BUSCO) analysis against insecta\_odb9 database (n=1,658) (Simão, Waterhouse,  
149 Ioannidis, Kriventseva, & Zdobnov, 2015). Generally, the N50 and L50 values can be affected by  
150 the genome sizes and chromosome numbers. Thus, we used the chromosome-normalized N50  
151 (CN50) and chromosome-normalized L50 (CL50) to evaluate the continuity at the level of  
152 hypothetical chromosomes (Jiao et al., 2017). Moreover, we assessed the accuracy of HaSCD2  
153 assembly by detecting genomic variants between the reference genome and sequencing data used  
154 for *de novo* genome assembly. In theory, heterozygous variants represent the discordant alleles and  
155 homozygous variants represent the misassembled segments. Specifically, Illumina short reads were  
156 used to detect single nucleotide polymorphisms (SNPs; GQ > 30, DP > 50) and short insertions and  
157 deletions (InDels; length < 50 bp, GQ > 30, DP > 50), PacBio long reads were used to detect  
158 structural variants (SVs; length > 50 bp) using Minimap2 v2.17 (Li, 2018) and Sniffle v1.0.12a  
159 (Sedlazeck, Rescheneder, et al., 2018), loci with less than 20 supporting reads or that were flagged  
160 as ‘imprecise’ (i.e., ambiguous breakpoints) and/or ‘unresolved’ (i.e., low quality) were removed.

### 161 **Repetitive elements annotation**

162 Repetitive elements including tandem and interspersed repetitive sequences were discovered based  
163 on homology and structural characteristics. First, we *de novo* constructed a species-specific  
164 transposable element library for the HaSCD2 assembly using RepeatModeler2 v2.0.1 (Flynn et al.,  
165 2020). Known repeat families extracted from Dfam\_3.2 (Hubley et al., 2016) and  
166 RepBase-20181026 databases (Bao, Kojima, & Kohany, 2015) were combined with  
167 species-specific repeat library, and RepeatMasker v4.1.0 (Tarailo-Graovac & Chen, 2009) was used  
168 to search and classify repetitive elements in the HaSCD2 assembly.

### 169 **Transcriptome sequencing and analyses**

170 Transcriptome of the SCD strain was conducted using individuals at three main developmental  
171 stages, including embryos (~2 days, 40 eggs), larvae (5 instars) and adults (2 male and 2 female  
This article is protected by copyright. All rights reserved

172 individuals). These specimens were snap-frozen in liquid nitrogen for extraction of total RNA using  
173 TRIzol (Invitrogen). The purities and concentrations of total RNA were measured using OneDrop™  
174 1000 Spectrophotometer (OneDrop Technologies). RNA-seq libraries were constructed using  
175 Illumina TruSeq RNA Sample Preparation Kit according to the manufacturer's protocol, and 150-bp  
176 short reads were sequenced on Illumina HiSeq 2500 platform. Low-quality sequences (Phred score  
177 < 10) and adapter sequences were trimmed in 4-bp sliding windows using Trimmomatic v0.36  
178 (Bolger, Lohse, & Usadel, 2014). The clean reads were mapped to the HaSCD2 assembly by  
179 HISAT2 v2.1.0 (Kim, Langmead, & Salzberg, 2015). StringTie v1.3.4 (Pertea et al., 2015) was used  
180 to parse sorted BAM files and to reconstruct transcripts from all transcriptome data. Redundant  
181 transcripts were removed with CD-HIT v4.8.1 (Fu, Niu, Zhu, Wu, & Li, 2012) with an identity  
182 cutoff of 0.9.

### 183 **Repetitive elements mask and gene prediction**

184 Repetitive elements identified above were masked using SeqKit v0.14.0 (Shen, Le, Li, & Hu, 2016)  
185 prior to gene annotation. We conducted *ab initio* gene prediction using BRAKER v2.1.0 (Hoff,  
186 Lange, Lomsadze, Borodovsky, & Stanke, 2016) invoking GeneMark-ET v4.38 (Lomsadze, 2005)  
187 and Augustus v3.3 (Stanke, Steinkamp, Waack, & Morgenstern, 2004). Namely, GeneMark-ET  
188 generated initial gene structures, then AUGUSTUS integrated the predicted genes and RNA-seq  
189 information to generate the final gene predictions. The protein sequences of *Acyrtosiphon pisum*  
190 (GCF\_005508785.1), *Apis mellifera* (GCF\_003254395.2), *Bombyx mori* (GCF\_014905235.1),  
191 *Drosophila melanogaster* (GCF\_000001215.4), *Helicoverpa armigera* (GCF\_002156985.1) and  
192 *Tribolium castaneum* (GCF\_000002335.3) were retrieved from NCBI RefSeq database (O'Leary et  
193 al., 2016). We used MAKER v2.31.10 (Holt & Yandell, 2011) to integrate *ab initio*,  
194 transcriptome-based and homology-protein evidence. The molecular functions of annotated genes  
195 were assigned against the UniProtKB (SwissProt + TrEMBL) database using Diamond v0.9.24  
196 (Buchfink, Xie, & Huson, 2015). Protein domains, Gene Ontology (GO) and pathways (MetaCyc  
197 and Reactome) were searched against Pfam (El-Gebali et al., 2019), PANTHER (Mi, Muruganujan,  
198 Ebert, Huang, & Thomas, 2019), Gene3D (Lewis et al., 2018), Superfamily (Wilson et al., 2009),  
199 and CDD databases (Marchler-Bauer et al., 2017) using InterProScan 5.34-73.0 (Mitchell et al.,  
200 2019). The KEGG pathway assignments were conducted using KofamKOALA (Aramaki et al.,  
201 2020) based on Kofam, a customized HMM database of KEGG Orthologs (KOs).

This article is protected by copyright. All rights reserved

---

## 202 **Collinearity**

203 We utilized BLASTP (Camacho et al., 2009) to perform all-vs-all homologous gene searches  
204 between *H. armigera* and other lepidopteran insects (*Bombyx mori*, *S. frugiperda*, *S. exigua*, *S.*  
205 *litura*) with a cut-off e-value of 1e-5. The top 5 BLASTP hits for a gene were selected to detect  
206 inter-genome collinear blocks using MCSanX toolkit (Wang et al., 2012) with default settings. The  
207 links of putative collinear blocks were visualized by Circos v0.69 (Krzywinski et al., 2009).

## 208 **Genomic comparisons between two subspecies**

209 The csiro4b assembly (Pearce et al., 2017) was retrieved from NCBI with RefSeq accession  
210 GCF\_002156985.1. The csiro4b assembly was assembled based on multiple sequencing libraries  
211 with different insert size and there remained 26,911 N-gaps as padding of unsolved regions,  
212 covering over 56.71 Mb genomic sequences (Figure 2c). We illustrated the process of closing  
213 N-gaps in the csiro4b assembly in Figure S3. Specifically, scaffolds in the csiro4b assembly were  
214 split into contigs at arbitrary Ns to avoid erroneous junctions, then we coalesced these contigs into  
215 pseudomolecules with padding of 100 Ns according to the HaSCD2 assembly using RaGOO v1.11  
216 (Alonge et al., 2019). Subsequently, we located the genomic coordinates of N-gaps in the csiro4b  
217 assembly and extracted additional 500-bp downstream and upstream flanking sequences of the  
218 N-gaps. To obtain corresponding locations and sizes of these unsolved regions in the HaSCD2  
219 assembly, we aligned PE flanking sequences to the HaSCD2 assembly using MUMmer v3.23  
220 (Kurtz et al., 2004). N-gaps were considered as successfully closed when they simultaneously  
221 satisfied the following criteria: (1) matched bases in PE flanking sequences were  $\geq 200$  bp long; (2)  
222 PE flanking sequences were assigned to the same chromosome.

## 223 **Field cotton bollworm collection and variant calling**

224 During 2016-2018, we used light traps to collect 141 adult cotton bollworm from 13 locations in  
225 mainland China, consisted of Yellow River Region (YRR) of Ac (n=7), Ay (n=10), Gy (n=5), Hm  
226 (n=10), Kf (n=10), Np (n=5), Xj (n=6); Changjiang River Region (CRR) of Aq (n=8), Df (n=19), Jz  
227 (n=20); and Northwestern Region (NR) of Hmi (n=12), Kel (n=12) and Sc (n=17) (Figure 3a, Table  
228 S8). All collected samples were preserved in absolute ethanol for genomic DNA extraction and  
229 library preparation, followed by sequencing on the Illumina HiSeq Xten platform. Partial  
230 mitochondrial cytochrome oxidase subunit I (COI) sequences were retrieved from raw sequencing

This article is protected by copyright. All rights reserved

231 reads by MitoZ (Meng, Li, Yang, & Liu, 2019) to confirm the *H. armigera* species status (Table  
232 S8). In addition, published re-sequencing resources, including 16 *H. a. conferta* from Australia, 7  
233 *H. punctigera* from Australia and 7 *H. zea* from Brazil (n=3) and America (n=4) were also  
234 integrated for downstream analyses (Anderson et al., 2018; Pearce et al., 2017). The sequence  
235 adapters and low-quality reads (Phred score < 20) were trimmed in 4-bp sliding windows using  
236 Trimmomatic v0.36 (Bolger et al., 2014). The clean re-sequencing reads were mapped to the  
237 HaSCD2 assembly using BWA v0.7.12 with default parameters (Li & Durbin, 2010). We used  
238 SAMtools v0.1.19 (Li et al., 2009) to remove duplicated reads. GATK v4.1.4.1 (McKenna et al.,  
239 2010) was used for variant calling and for filtering out low-quality variants as described above.  
240 Finally, the high-quality SNPs were annotated using SnpEff v5.0d (Cingolani et al., 2012).

### 241 **Population diversity and structure**

242 The pairwise nucleotide diversity ( $\theta\pi$ ), Watterson's estimator ( $\theta_w$ ), expected heterozygosity ( $H_E$ )  
243 and Tajima's  $D$ , were calculated using VCFtools v0.1.16 (Danecek et al., 2011). Loci with minor  
244 allele frequency (MAF) < 0.01 and genotyping rate < 90% were further removed for downstream  
245 population structure analyses. We employed VCF2Dis v1.42  
246 (<https://doi.org/github.com/BGI-shenzhen/VCF2Dis>) to calculate pairwise distance ( $p$ -dist) matrices  
247 between individuals, and constructed neighbor-joining (NJ) trees using the 'neighbor' module in  
248 PHYLIP v3.695 (<https://doi.org/evolution.genetics.washington.edu/phylip.html>). Phylogenetic trees  
249 were visualized using FigTree (<http://doi.org/tree.bio.ed.ac.uk/software/figtree/>). We also  
250 performed principal component analysis (PCA) using the smartPCA module built-in EIGENSOFT  
251 v7.2.1 with no automatic outlier removal allowed (Patterson, Price, & Reich, 2006). The  
252 Tracy-Widom test was applied to demonstrate statistical significance of different eigenvectors.  
253 Finally, individual ancestries were quantified using fastStructure (Raj, Stephens, & Pritchard,  
254 2014).

### 255 **Inference of demographic history**

256 We selected two representative individuals with high sequencing depth (> 25-fold, Table S8) to  
257 infer the historical effective population size of CRR & YRR and NR lineages using Pairwise  
258 Sequentially Markovian Coalescence (PSMC) v0.6.5 with 100 bootstrap replications (-b 100) (Li &  
259 Durbin, 2011). Other parameters were set as following: "-N25 -t15 -r5 -p 4+25\*2+4+6". The

260 SMC++ v1.15.4 which was shown to have a better resolution in inferring demographic history of  
261 the recent past than the PSMC program (Terhorst, Kamm, & Song, 2017), was also used to infer  
262 more recent demographic history of CRR & YRR and NR lineages based on multiple individuals,  
263 and outputs from these two modelers were scaled with a generation time of 0.25 year and a  
264 mutation rate of  $2.9 \times 10^{-9}$  (Keightley et al., 2015).

### 265 **Detection of divergent patterns and adaptive footprints**

266 All female individuals (i.e., individuals with the heterogametic sex chromosome) were removed to  
267 avoid adverse effects on divergence computation. The male individuals with diploid  
268 Z-chromosomes have comparable sequence depth coverage between Z-chromosome and autosomes.  
269 We identified 90 male cotton bollworms for downstream analyses, including 47 individuals from  
270 CRR & YRR, 33 individuals from NR and 10 *H. a. conferta* individuals (Table S8 and Table S9).  
271 We split the HaSCD2 assembled 31 chromosomes into 10 kb sliding windows with a step size of 5  
272 kb. Windows with more than 10 SNPs were retained for calculating inter-populations fixation index  
273 ( $F_{ST}$ ) using VCFtools v0.1.16 (Danecek et al., 2011) with negative  $F_{ST}$  values treated as zero. We  
274 performed Mantel test to evaluate the relationship between genetic distance ( $F_{ST}$ ) and geographic  
275 distance (km) using the ‘vegan’ package  
276 (<https://doi.org/cran.r-project.org/web/packages/vegan/index.html>). The significance of the  
277 correlation was estimated based on 9,999 permutations.

278 Complementary methods can be used to detect selective signatures. We calculated fixation index  
279 ( $F_{ST}$ ) and nucleotide diversity ratio ( $\theta_{\pi A}/\theta_{\pi B}$ ) for each sliding window (5-kb windows with the step  
280 size of 2 kb). The  $F_{ST}$  values were Z-transformed and the nucleotide diversity ratios were  
281  $\log_2$ -transformed. Selective signatures with the highest 5%  $Z(F_{ST})$  and the highest 5%  $\log_2(\theta_{\pi A}/\theta_{\pi B})$   
282 were considered to be candidate selection targets of population B, whereas the lowest 5%  
283  $\log_2(\theta_{\pi A}/\theta_{\pi B})$  were population A. Genes that overlapped more than 50% in length with candidate  
284 selection targets were defined as putative selective genes. The Gene Ontology (GO) terms of  
285 putative selective genes were retrieved from the HaSCD2 annotation, and GO enrichment analysis  
286 was carried out to identify gene sets that were significantly over-represented relative to common  
287 genes using the R package clusterProfiler (Yu, Wang, Han, & He, 2012).

## 288 **RESULTS**

This article is protected by copyright. All rights reserved

---

## 289 **Chromosome-scale assembly of the cotton bollworm**

290 Our sequencing generated 5.48 million long reads (~137-fold coverage; N50 size of 12.89 kb) using  
291 single-molecule real-time (SMRT) sequencing, and 132.95 million PE reads (~112-fold coverage)  
292 using the Illumina HiSeq Xten platform (Table S1). The genome size and heterozygosity rate of  
293 cotton bollworm were estimated to be ~352.30 Mb and ~0.84% respectively (Figure S4). The two  
294 genome assemblers, i.e., Falcon and Canu, produced two draft assemblies with total sizes of 398.01  
295 Mb (N50 size of 4.81 Mb) and 454.70 Mb (N50 size of 4.17 Mb), respectively. After merging  
296 assemblies and discarding alternative haplotypes, we obtained a 356.67 Mb primary assembly,  
297 named as HaSCD2, that consisted of 106 contigs with a N50 length of 7.34 Mb (Table 1). Our  
298 assembly is 5.81% larger than the published *H. a. conferta* assembly (csiro4b, 337.07 Mb; Pearce et  
299 al., 2017). Importantly, our assembly is much more continuous, i.e., contig N50 of 7.34 Mb vs.  
300 23.48 kb, representing an approximately 313-fold increase in average contig length (Figures 1a and  
301 1b, Table 1).

302 Clean sequence data from 100 BC1 cotton bollworm offspring and their parents (~5.88-fold average  
303 coverage for offspring; ~18.83-fold average coverage for parents, Table S2) were mapped to the  
304 HaSCD2 assembly. This yielded 425,788 female-informative markers that were assigned into 31  
305 linkage groups (LGs), and 756,597 male-informative markers that spanned a genetic distance of  
306 1545.4 cM (Table S3). Finally, 94 of the 106 contigs (353.56 Mb) were anchored into 31 LGs, of  
307 which 351.05 Mb sequences were further oriented (Figure 1c, Table S3). Remarkably, eight  
308 chromosomes are represented by a single contig, including the Z-chromosome, the longest  
309 chromosome in HaSCD2 assembly (Table S3).

## 310 **Quality assessment of HaSCD2 assembly**

311 We confirmed that the HaSCD2 assembly was highly complete through the BUSCO assessment  
312 (97.7% completeness) (Table S4). The CN50 and CL50 values of the HaSCD2 assembly were 7.34  
313 Mb and 1 (Table S5), indicating that nearly all chromosomes were assembled by a few contigs only.  
314 Irrespective of whether the NGS data were sourced from *H. a. armigera* or *H. a. conferta*  
315 (Anderson et al., 2018), the mapping rate and coverage breadth of the HaSCD2 assembly were  
316 significantly higher than those of the csiro4b assembly (Figures 2a and 2b), which resulted from  
317 longer length of most chromosomes (Figure 2c). The conserved synteny between cotton bollworm

318 and other moths confirmed the reliability of our chromosome assignment on the basis of the genetic  
319 map (Figure S5). Moreover, we only detected 129, 273 and 2,322 homozygous SNPs, InDels and  
320 SVs (covering ~0.52 Mb) and thus estimated the sequence accuracy of the HaSCD2 assembly to be  
321 ~99.85%. In particular, most heterozygous variants were restricted to hotspots in the genome, which  
322 have failed to become homozygous possibly due to closely linked lethal alleles segregating in  
323 repulsion (Figure S6). Overall these results confirm our HaSCD2 genome is more comprehensive  
324 and will be of greater utility for conducting genome-wide analyses of the cotton bollworm.

### 325 **Repeat and gene annotation**

326 In total, 96.44 Mb sequences (27.04%) were annotated as repetitive elements in the HaSCD2  
327 assembly, of which 74.56% of repeats were classified into various repetitive families that included  
328 long interspersed nuclear elements (LINEs, 21.04%), rolling-circle transposons (RCs, 18.23%),  
329 long terminal repeats (LTRs, 14.46%), short interspersed nuclear elements (SINEs, 9.58%) and  
330 DNA transposon (8.32%) (Table S6). Compared to the csiro4b assembly (14.60%, 49.21 Mb;  
331 Pearce et al., 2017), repetitive elements annotated in HaSCD2 were increased by 95.98% in length,  
332 owing to the better performance of long reads sequencing technology in resolving repetitive  
333 sequences. In addition, the total proportion of repetitive elements in the HaSCD2 assembly was  
334 lower than the other noctuid pests with larger genome sizes, indicating that repetitive elements are  
335 an important factor affecting differences in the size of noctuid genomes (Table S7).

336 A total of 18,668 protein-coding genes were annotated in HaSCD2 assembly, and the mean gene  
337 length was 6,776.03 bp (Table 1). The BUSCO analysis revealed 95.0% of single-copy genes in  
338 insecta\_odb9 database were identified as complete in our HaSCD2 assembly (Table S4). The gene  
339 number in the HaSCD2 assembly was slightly more than the csiro4b assembly (17,086), which  
340 suggested that the HaSCD2 offered opportunities to unmask protein-coding genes in those  
341 previously unsolved regions in the csiro4b genome (Table 1). In addition, approximately 95.37% of  
342 genes were functionally annotated, of which 9,067 and 6,689 genes were assigned with GO and  
343 KEGG accessions, respectively.

### 344 **Comparisons of HaSCD2 and csiro4b assemblies**

345 A total of 12,115 N-gaps (45.02%) were successfully localized to unique locations based on their  
346 flanking sequences matches, and corresponded to greater than 22.45 Mb in the HaSCD2 assembly

This article is protected by copyright. All rights reserved

(Figure 2d). Noticeably, 44.45% of these assigned N-gaps (9.98 Mb) overlapped with various repetitive families, including long interspersed nuclear elements (LINEs, 2.07 Mb), rolling-circles (RCs, 1.65 Mb), long terminal repeat (LTRs, 1.64 Mb), short interspersed nuclear elements (SINEs, 0.89 Mb) and DNA transposons (0.84 Mb) (Figure 2e). The total proportion of repetitive sequences in N-gaps was higher than the whole genome, indicating that the performance of NGS in genome assembly was greatly hindered by repetitive sequences. On the other hand, the proportions of different repetitive families in N-gaps regions were very close to that in the whole HaSCD2 genome, revealing that N-gaps regions have no obvious bias to any repeat category (Figure 2e). Meanwhile, these assigned N-gaps overlapped with exonic (2.71 Mb), intronic (5.67 Mb) and intergenic (13.80 Mb) regions, and ~2,987 genes were improved due to at least one missing fragments in their exons (Figure S6). For example, we found a completely missing and a partially fragmentary odorant binding proteins (OBPs) in an assigned N-gap located in HaChr18 of HaSCD2 assembly (Figure 2f).

### 360 **Unique lineage of cotton bollworm in northwestern China**

361 We conducted a population genomic study of 141 individuals of cotton bollworm, which were  
362 collected from 13 locations in three cotton-producing regions of China, i.e., the Yellow River  
363 Region (YRR), the Changjiang River Region (CRR) and the Northwestern Region (NR) (Wu &  
364 Guo, 2005) (Figure 3a, Table S8). We generated a total of 746.74 Gb PE sequencing data, 88.25%  
365 of which were mapped to the HaSCD2 assembly with a mean depth of 10.21-fold (Table S8). By  
366 integrating publicly available data of 30 genomes from the *Helicoverpa* complex (see MATERIALS  
367 AND METHODS, Table S9), we identified 28,714,909 bi-allelic SNPs, including 2,891,882  
368 non-synonymous variants (Table S10), and they covered more than 87.10% of the HaSCD2  
369 assembly (Table S11).

370 After stringent quality filtering, we refined 5,227,071 high-quality SNPs to infer their evolutionary  
371 attributes. The pairwise *p*-dist matrices based on genome-wide, Z-chromosome, and non-coding  
372 SNPs, were calculated to construct neighbor-joining (NJ) trees, respectively (Figure 3b, Figure S7).

373 All phylogenetic trees consistently showed that the cotton bollworm was clustered into three  
374 distinct groups, that is, *H. a. conferta* sampled in Australia, *H. a. armigera* sampled in CRR and  
375 YRR, and *H. armigera* sampled in NR. The result was consistent with the sampling locations of  
376 these specimens. Similarly, principal component analysis (PCA) demonstrated that the first two  
This article is protected by copyright. All rights reserved

377 eigenvectors jointly separated all individuals into three clusters, and they explained ~4.72%  
378 variation present across the 13 Chinese populations as well as a *H. a. conferta* population (Figure  
379 3c). We further used fastStructure (Raj et al., 2014) to gain insights into the ancestry of each cotton  
380 bollworm individual (Figure 3d). When  $K = 2$ , NR lineage, rather than *H. a. conferta* subspecies,  
381 was initially separated from the other groups, which indicated their unique genetic composition and  
382 perhaps warranted them being considered their own subspecies. When  $K = 3$ , three distinct genetic  
383 clusters were in agreement with the phylogenetic tree and PCA analysis. Furthermore, asymmetric  
384 genetic admixture events seem to exist from CRR & YRR to the NR lineage (Figure 3d).

### 385 **Population diversity and demographic history**

386 As shown in Table S12, the genome-wide genetic diversity of different populations in the NR  
387 lineage was similar to that in CRR and YRR lineages. The mean Tajima's  $D$  values of different  
388 populations ranged from -1.22 to -0.67 (Table S12), meaning that different cotton bollworm  
389 populations have an excess of low frequency alleles. At the lineage levels, the high proportions of  
390 negative Tajima's  $D$  values across the whole genome (97.26% for CRR & YRR and 97.49% for  
391 NR) suggest that population expansion was the most likely cause (Biswas & Akey, 2006) (Figure  
392 4a).

393 PSMC analysis revealed that CRR & YRR and NR lineages have maintained a similar and constant  
394 effective population sizes for up to *ca.* 0.2 million years ago (Mya), after which, their effective  
395 population sizes began to simultaneously increase but at different growth rates, coinciding with the  
396 most recent rise of global temperature (Figure 4b). At marine isotope stage 5 (MIS5), the effective  
397 population size of CRR & YRR exceeded that of the NR lineage. During the last glaciation (LG),  
398 the global temperature began to gradually fall and effective population sizes of the two lineages  
399 reached a plateau. Then when the temperature decreased to its lowest point the NR lineage (but not  
400 the CRR & YRR lineage) had a decrease in population size. Subsequently, the CRR & YRR and the  
401 NR lineages started to expand again as the global temperature rose. In general, historical global  
402 temperature profoundly affected the effective population sizes of the CRR & YRR and the NR  
403 lineages and may therefore likely have facilitated their divergence. In addition, the PSMC result  
404 further confirmed the assumption that the genome-wide negative Tajima's  $D$  values were the result  
405 of population expansions. SMC++ analysis further revealed the dramatic population expansion of  
406 the CRR & YRR lineage between ~5,000 – 8,000 years ago (Figure S8).

This article is protected by copyright. All rights reserved

---

407 **Genomic divergence among different lineages**

408 The fixation indices ( $F_{ST}$ ) in pairwise cotton bollworm populations ranged from 0 to 0.076 (Figure  
409 S9). Several populations from YRR showed no divergence (i.e.,  $F_{ST} = 0$ ) which indicated frequent  
410 gene exchange underpinned their homogenization. A Mantel test showed that the divergence levels  
411 of different populations was highly correlated with their geographical distances (Figure S10,  $P =$   
412  $1e-04$ ). This suggested that spatial distance has been highly influential in generating the observed  
413 genetic variation. We also compared the differentiation levels among NR, CRR & YRR, and *H. a.*  
414 *conferta*. We found *H. a. conferta* vs. NR exhibited the highest divergence levels ( $F_{ST} = 0.069$ ),  
415 followed by *H. a. conferta* vs. CRR & YRR ( $F_{ST} = 0.041$ ), while CRR & YRR vs. NR was lowest  
416 ( $F_{ST} = 0.025$ ), which was still higher than the divergence level between the corn and rice strains of  
417 *Spodoptera frugiperda* ( $F_{ST} = 0.019$ ; Gouin et al., 2017). Overall, our genomic analyses further  
418 suggested that the NR is likely to be a unique cotton bollworm lineage.

419 Among these three *H. armigera* lineages, the divergence levels of Z-chromosome in the male  
420 individuals were 5.38-, 2.67-, and 4.44-fold higher than autosomes in the comparisons of CRR &  
421 YRR vs. NR, CRR & YRR vs. *H. a. conferta*, and NR vs. *H. a. conferta*, respectively (Figure 5a).  
422 The genomic regions with the highest 1%  $F_{ST}$  value in the above comparisons, 63.55%, 34.12%,  
423 and 61.33% regions were located in the Z-chromosome. Thus, we deduced that the divergence rate  
424 of the Z-chromosome was faster than autosomes of cotton bollworm either at the level of the whole  
425 chromosomes or at outlier loci. Surprisingly, discordant trends of divergence levels were detected in  
426 Z-chromosome and autosomes among different pairwise lineages, i.e., the Z-chromosome  
427 divergence level of CRR & YRR vs. *H. a. conferta* was lower than that of CRR & YRR vs. NR  
428 (mean  $F_{ST} = 0.085$  and  $0.091$ , Wilcoxon test,  $W=5,687,304$ ,  $P = 2.43e-05$ ), whereas the former had  
429 obvious elevation in autosomal divergence (mean  $F_{ST} = 0.032$  and  $0.017$ , Wilcoxon test,  
430  $W=731,287,952$ ,  $P < 2.2e-16$ ).

431 A total of 25 windows were identified with the highest 1%  $F_{ST}$  values among all three comparisons  
432 (Figure 5b). These windows covered 195 kb genomic regions, of which 76.92% were located on the  
433 Z-chromosome. From these highly divergent regions, nine protein-coding genes were annotated  
434 (four genes located in the Z-chromosome), including two clock genes (*per* and *clk*), a G-protein  
435 coupled receptor and Acyl-CoA desaturase (Table S13). Spatiotemporal transcriptome analysis  
436 showed that *per* and *clk* genes broadly expressed in most tissues and developmental stages (Figure  
This article is protected by copyright. All rights reserved

437 5c). Notably, the *per* and *clk* genes act in the negative feedback loop of the core circadian system  
438 and they have been shown to be associated with nocturnal flight, feeding behavior and endocrine  
439 activity (Tomioka & Matsumoto, 2015). We identified 17 and 2 non-synonymous SNPs in *per* and  
440 *clk* genes, respectively, which could contribute to the differences in facultative diapause and annual  
441 generations observed among different cotton bollworm lineages (Table S14). Expression levels of  
442 two circadian genes (*per* and *clk*) in adults were also markedly higher than in larvae and pupae  
443 (Figure 5c), which could highlight their critical roles in behavioral rhythms during adulthood, such  
444 as nocturnality and mating activity. In addition, expression levels of the two circadian genes in  
445 peripheral tissues were higher than other adult tissues, such as in antennae, abdomen and thorax  
446 (Figure 5c), emphasizing the divergence of circadian genes may be related to orientation and  
447 reproduction (Merlin, Gegear, & Reppert, 2009; Tomioka & Matsumoto, 2015).

#### 448 **Selection signatures on climate adaptation, feeding behavior and insecticide tolerance**

449 Loci that feature in different cotton bollworm lineages may have roles responding to environmental  
450 factors, such as low temperature and deficient diets in NR populations, or the different pest  
451 management strategies in CRR & YRR and in *H. a. conferta* (Table S14). Using a combination of  
452  $\log_2$ -transformed  $\theta_{\pi A}/\theta_{\pi B}$  and  $Z$ -transformed  $F_{ST}$  values (see MATERIALS AND METHODS), we  
453 detected 310/349 selective windows for the NR lineage when compared to CRR & YRR and *H. a.*  
454 *conferta*, 185/192 selective windows for CRR & YRR when compared to NR and *H. a. conferta*,  
455 and 145/269 selective windows for *H. a. conferta* when compared to NR and CRR & YRR (Figure  
456 6a, Figure S11).

457 For the NR lineage, two GO categories were commonly enriched irrespective of whether *H. a.*  
458 *conferta* or CRR & YRR was set as the control lineage (Table S15), i.e., oxidation-reduction  
459 process (GO:0055114) and ATPase activity (GO:0016887). Another three GO categories related to  
460 energy metabolic processes were enriched in a single comparison (Table S15). These intense  
461 selection signatures were consistent with expected greater energy consumption for the NR lineage  
462 to adapt to the extreme winter, especially the process of resisting severe cold and looking for host  
463 plants in the next spring (Arrese & Soulages, 2010). A previous study of *Drosophila melanogaster*  
464 populations has shown that fluctuating temperature could induce genes involving in “ATP synthesis  
465 coupled proton transport” category (Hsu, Belmouaden, Nolte, & Schlötterer, 2020). Meanwhile, the  
466 down-regulation of genes involved in “oxidation-reduction process” and “ATP synthesis coupled  
This article is protected by copyright. All rights reserved

467 proton transport” is also considered to be the response of Antarctic midges to extreme desert  
468 habitats (Teets et al., 2012). These results suggest that diverse insects have evolved similar climate  
469 adaptation strategies.

470 For *H. a. conferta*, we annotated 203 and 289 genes in these selective regions when compared to  
471 NR and CRR & YRR lineages, and four GO categories were significantly overrepresented (Table  
472 S16). Two enriched GO categories, gustatory and olfactory receptors, are identified, of which  
473 gustatory receptors are the only chemosensory protein family to display radiations in lepidopteran  
474 insects with narrow host ranges (Pearce et al., 2017), which presumably indicated their important  
475 roles in insect-plant interactions (Xu, 2020). Previous research has revealed that these  
476 chemosensory genes were differentially expressed when cotton bollworm were fed on different host  
477 plants (Pearce et al., 2017). In addition, olfactory receptors are required by the insect to detect and  
478 respond to innumerable volatile odors in the field, and they are relevant to foraging, mate  
479 recognition, and oviposition site selection (Montagné, de Fouchier, Newcomb, & Jacquin-Joly,  
480 2015). When NR was set as the control lineage, the GO term named “structural constituent of  
481 cuticle” was uniquely enriched. For lepidopteran pests, the cuticle is the first barrier to slow down  
482 the penetration of insecticide molecules within their bodies (Balabanidou, Grigoraki, & Vontas,  
483 2018), thus natural selection on genes encoding cuticle may effectively decrease the efficiency of  
484 insecticides in the field.

485 We further detected a specific divergent signature between CRR & YRR and NR lineages spanning  
486 a 120 kb window on chromosome 7 (Figure 6a). In this genomic region, the Tajima’s *D* values in  
487 the NR lineage were significantly higher than the CRR & YRR lineage, indicating positive selection  
488 could lead to local decrease of Tajima’s *D* values (Figure 6b). The nucleotide diversity in the CRR  
489 & YRR lineage also decreased, indicating that genes in this regions experienced intense positive  
490 selection (Figure 6c). Six glutathione S-transferase (GST) genes were annotated in this region  
491 (Figure 6d, Table S17). These GSTs that belonged to the Delta clan have well-established roles in  
492 xenobiotic degradation, including insecticide and plant secondary metabolites (Cui, Rui, Yang,  
493 Wang, & Yuan, 2017; Low et al., 2010; Pearce et al., 2017; Rane et al., 2019). Our analysis of  
494 public transcriptome data (Jin et al., 2019) revealed *GSTD1h*, *GSTD1k* and *GSTD1l* were  
495 constitutively expressed in insecticide treatment and control groups. Moreover, none of these GSTs  
496 are differentially expressed in response to different insecticides (Figure 6f). In this region, we also

This article is protected by copyright. All rights reserved

497 detected two linkage disequilibrium (LD) blocks (Figure 6e), which may suggest that at least two  
498 GST genes were targeted by the natural selection rather than a single selection with extended  
499 genetic hitchhiking. We identified 75 missense SNPs within these divergent signatures and further  
500 identified 15 loci whose alleles were specific to the CRR & YRR or the NR lineage (Table S18).  
501 Thus, these loci may provide critical cues for understanding host adaptation.

## 502 DISCUSSION

503 In this study, we generated a chromosome-level reference genome of *H. armigera* based on the  
504 long-read sequencing technology and an ultra-dense genetic map. This new genome offers not only  
505 significant advancements on the preceding work but also represents a near complete reference  
506 genome resource for this species. The HaSCD2 assembly had significant advances in a series of  
507 assessments, especially the CN50 value which was even superior to other recently published  
508 reference genomes of lepidopteran pests (Table S7). In addition, the number of different repetitive  
509 DNA families annotated in HaSCD2 assembly were dramatically elevated, and this provides  
510 fundamental information critical to resolve complex biological and evolutionary questions related to  
511 repetitive sequences (Klai et al., 2020; Tay, Behere, Batterham, & Heckel, 2010). In the genomic  
512 regions that failed to be assembled in the csiro4b draft genome, we found numerous genes including  
513 two OBP genes that could be biologically significant to receive odorants in the olfactory system  
514 (Figure 2f). Therefore, our assembly represents a significantly improved genome resource of this  
515 global pest, offering opportunities for more comprehensive genome-wide studies owing to its  
516 excellent features in completeness and continuity.

517 With applications of different molecular markers, deductions on the genetic structure of the cotton  
518 bollworm have been revised several times at different scales (Anderson et al., 2016; Behere et al.,  
519 2007; Behere, Tay, Russell, Kranthi, & Batterham, 2013; Nibouche, Buès, Toubon, & Poitout,  
520 1998; Weeks et al., 2010). In particular, several studies based on whole-genome SNPs (Anderson et  
521 al., 2016; Anderson et al., 2018; Pearce et al., 2017) consistently revealed the subspecies boundary  
522 between *H. a. armigera* and *H. a. conferta*, and the presence of gene flow between *H. a. armigera*  
523 and *H. zea*, which emphasized the biological risks of allopatric lineages. We jointly analyzed 14  
524 cotton bollworm populations (13 Chinese populations and 1 Australian population) to investigate  
525 genetic structure and diversity within *H. armigera* complex. We found the NR individuals collected

526 from northwestern China formed a unique cotton bollworm lineage, which was genetically distinct  
527 from the known *H. a. conferta* and *H. a. armigera* lineages. This conclusion may well explain why  
528 NR and CRR & YRR lineages have heterogeneous resistance levels to chemical and Bt insecticides  
529 (Yang et al., 2013). Noticeably, we found that three outliers in the PCA were sampled from Sc  
530 population that lies on the western region of Xinjiang. Previous research has demonstrated that  
531 European corn borer, *Ostrinia nubilalis*, invaded into Yili area and made secondary contact with  
532 native Asian corn borer, *O. furnacalis* (Wang et al., 2017). Thus, these three outliers may provide  
533 an important cue for the likely presence of other genetically distinct *H. armigera* lineages in the  
534 India subcontinent, Europe, and/or Africa. Variations in genomic diversity in several other insects  
535 have been attributed to factors, such as climate stress, artificial domestication and biological  
536 invasion (Ding et al., 2018; Montero-Mendieta et al., 2019; Xia et al., 2009; Wu et al., 2019).  
537 Interestingly, although suffering from severe coldness and deficient diets, NR lineage maintains a  
538 similar level of genetic diversity to the CRR & YRR lineage (Table S14). Therefore, the NR lineage  
539 exhibits remarkable genetic plasticity and capacities to adapt to intensive planting of transgenic Bt  
540 cotton and has the potential to generate new mutations conferring Bt resistances such as that  
541 observed in the northern region of China. The explicit population relationships we have identified  
542 form a basis to monitor and manage this pest complex across spatial and temporal scales.

543 Contemporary studies in a wide range of taxa are using genomic variation to track the biological  
544 processes of speciation and population differentiation on a landscape scale (Wolf & Ellegren,  
545 2017). Aided by an ultra-dense variation map, we investigated the patterns of divergence in three  
546 distinct cotton bollworm lineages. We found that Z-chromosomes had higher divergence levels than  
547 autosomes in any pairwise lineage comparisons, reflecting their faster cumulative rates of  
548 differentiated loci. Similar phenomena have been reported during the speciation processes of the  
549 sibling species of birds, fish, and other moths (Presgraves, 2018; Van Belleghem et al., 2018). The  
550 Z-chromosome is less exchangeable between species, and thus more refractory to gene  
551 introgression than autosomes (Presgraves, 2018). We infer gene introgression from CRR & YRR to  
552 NR effectively reduced their interspecific divergence level in the autosomes, but the divergence  
553 level in Z-chromosome remained similar across *H. armigera* subspecies and populations. We also  
554 identified nine universally divergent genes in inter-lineage comparisons, including two circadian  
555 genes (*per* and *clk*). It is worth noting that the *per* gene was among the genes that exhibited elevated

556  $F_{ST}$  in a less well-powered but independent study of Chinese and Australian *H. armigera*  
557 populations (Song et al., 2018). Previous studies developed mitochondrial COI markers or used full  
558 mitochondrial genomes to distinguish species in the *Helicoverpa* complex (Behere et al., 2008;  
559 Behere et al., 2007; Walsh et al. 2019) and track the invasion routes of *H. a. armigera* (Tay et al.  
560 2017; Arnemann et al., 2019). However, mitochondrial markers are maternally transmitted and  
561 non-recombining compared to nuclear markers, and they have limited resolution to distinguish  
562 different lineages within *H. armigera* (Behere et al., 2007). These nine genes we listed in Table S13  
563 have the potential to be developed as lineage-specific markers to accurately identify the different  
564 evolutionary lineages within the *H. armigera* subspecies.

565 The patterns of differentiation across the genomes measured by  $F_{ST}$  imply local environments could  
566 shape genomic divergence. Our analyses indicated that genes that were involved in climate  
567 adaptation, feeding behaviors, and insecticide tolerance, were also under strong positive selection.  
568 Previous research also demonstrated that different host plants could induce these GST genes to  
569 become highly expressed in detoxification and digestive tissues, such as mid-gut, Malpighian  
570 tubules and fat body (Pearce et al., 2017), which is consistent with its recent duplication and  
571 diversification evolution. The Delta clan of GSTs highlighted in this study have been inferred to be  
572 involved in the adaptation of fly species to brassica (Gloss et al., 2014). Cotton bollworm residing  
573 in the northwestern China have a unique evolutionary ancestry and local habitat (Figure 3 and Table  
574 S14), and they are still sensitive to most chemical insecticides and Bt toxins. Meanwhile, a  
575 dominant point mutation underlying resistance to Cry1Ac was only detected in northern China  
576 populations (Jin et al., 2018), and highlighted the need for concerted management practice to  
577 prevent the introgression of the dominant Cry1Ac resistance gene from the northern populations  
578 into the NR populations of *H. armigera* in China.

## 579 **ACKNOWLEDGEMENTS**

580 This study was supported by grants from the National Natural Science Foundation of China (grant  
581 no. 31930093 to YW) and the SAFEA of China (grant no. BP0719029 to YW). We thank Yanhui  
582 Lu for providing cotton bollworm samples from northwestern China, and Shuai Zhan for critical  
583 reading of the manuscript. WTT was supported by CSIRO Health & Biosecurity.

## 584 **CONFLICT OF INTEREST**

This article is protected by copyright. All rights reserved

---

585 The authors declare that there is no conflict of interest.

## 586 AUTHOR CONTRIBUTIONS

587 Yidong Wu, Jianpeng Zhang, Feng Zhang and Yihua Yang conceived the ideas and designed the  
588 methodology; Jianpeng Zhang, Yu Shi and Fang Guan conducted the molecular experiments;  
589 Jianpeng Zhang, Feng Zhang, Wee Tek Tay, and Charles Robin analysed the data; and Jianpeng  
590 Zhang and Yidong Wu led the writing of the manuscript. All authors contributed critically to the  
591 drafts and gave final approval for publication.

592

## 593 DATA AVAILABILITY STATEMENT

594 The HaSCD2 assembly and corresponding annotation file were uploaded at Figshare:  
595 <https://doi.org/10.6084/m9.figshare.14913858.v1>. The PacBio and Illumina sequencing reads for  
596 genome assembly and gene annotation were deposited in the NCBI database with the BioProject  
597 Number of PRJNA623133. The whole genome resequencing data of the pedigreed individuals was  
598 deposited with the BioProject Number of PRJNA730914. The raw sequencing data of field cotton  
599 bollworm was deposited with the BioProject Number of PRJNA731848, and the high-quality SNPs  
600 was at [https://figshare.com/articles/dataset/population\\_genomics\\_of\\_cotton\\_bollworm/17008922](https://figshare.com/articles/dataset/population_genomics_of_cotton_bollworm/17008922).

## 601 REFERENCES

- 602 Alonge, M., Soyk, S., Ramakrishnan, S., Wang, X., Goodwin, S., Sedlazeck, F. J., ... Schatz, M. C. (2019).  
603 RaGOO: fast and accurate reference-guided scaffolding of draft genomes. *Genome Biology*, 20(1), 224.  
604 <https://doi.org/10.1186/s13059-019-1829-6>
- 605 Anderson, C. J., Tay, W. T., McGaughran, A., Gordon, K. H. J., & Walsh, T. K. (2016). Population structure  
606 and gene flow in the global pest, *Helicoverpa armigera*. *Molecular Ecology*, 25(21), 5296–5311.  
607 <https://doi.org/10.1111/mec.13841>
- 608 Anderson, C. J., Oakeshott, J. G., Tay, W. T., Gordon, K. H. J., Zwick, A., & Walsh, T. K. (2018).  
609 Hybridization and gene flow in the mega-pest lineage of moth, *Helicoverpa*. *Proceedings of the*  
610 *National Academy of Sciences*, 115(19), 5034–5039. <https://doi.org/10.1073/pnas.1718831115>
- 611 Aramaki, T., Blanc-Mathieu, R., Endo, H., Ohkubo, K., Kanehisa, M., Goto, S., & Ogata, H. (2020).  
This article is protected by copyright. All rights reserved

- 612 KofamKOALA: KEGG Ortholog assignment based on profile HMM and adaptive score threshold.  
613 *Bioinformatics*, 36(7), 2251–2252. <https://doi.org/10.1093/bioinformatics/btz859>
- 614 Arnemann, J. A., Roxburgh, S., Walsh, T. K., Guedes, J., Gordon, K. H. J., Smagghe, G., & Tay, W. T.  
615 (2019). Multiple incursion pathways for *Helicoverpa armigera* in Brazil show its genetic diversity  
616 spreading in a connected world. *Scientific Reports*, 9(1), 19380.  
617 <https://doi.org/10.1038/s41598-019-55919-9>
- 618 Arrese, E. L., & Soulages, J. L. (2010). Insect Fat Body: Energy, Metabolism, and Regulation. *Annual*  
619 *Review of Entomology*, 55(1), 207–225. <https://doi.org/10.1146/annurev-ento-112408-085356>
- 620 Balabanidou, V., Grigoraki, L., & Vontas, J. (2018). Insect cuticle: a critical determinant of insecticide  
621 resistance. *Current Opinion in Insect Science*, 27, 68–74. <https://doi.org/10.1016/j.cois.2018.03.001>
- 622 Bao, W., Kojima, K. K., & Kohany, O. (2015). Repbase Update, a database of repetitive elements in  
623 eukaryotic genomes. *Mobile DNA*, 6(1), 4–9. <https://doi.org/10.1186/s13100-015-0041-9>
- 624 Behere, G. T., Tay, W. T., Russell, D. A., & Batterham, P. (2008). Molecular markers to discriminate among  
625 four pest species of *Helicoverpa* (Lepidoptera: Noctuidae). *Bulletin of Entomological Research*, 98(6),  
626 599–603. <https://doi.org/10.1017/S0007485308005956>
- 627 Behere, G. T., Tay, W. T., Russell, D. A., Heckel, D. G., Appleton, B. R., Kranthi, K. R., & Batterham, P.  
628 (2007). Mitochondrial DNA analysis of field populations of *Helicoverpa armigera* (Lepidoptera:  
629 Noctuidae) and of its relationship to *H. zea*. *BMC Evolutionary Biology*, 7(1), 117.  
630 <https://doi.org/10.1186/1471-2148-7-117>
- 631 Behere, G. T., Tay, W. T., Russell, D. A., Kranthi, K. R., & Batterham, P. (2013). Population genetic  
632 structure of the cotton bollworm *Helicoverpa armigera* (Hübner) (Lepidoptera: Noctuidae) in India as  
633 inferred from EPIC-PCR DNA markers. *PLoS ONE*, 8(1), e53448.  
634 <https://doi.org/10.1371/journal.pone.0053448>
- 635 Biswas, S., & Akey, J. M. (2006). Genomic insights into positive selection. *Trends in Genetics*, 22(8), 437–  
636 446. <https://doi.org/10.1016/j.tig.2006.06.005>
- 637 Bolger, A. M., Lohse, M., & Usadel, B. (2014). Trimmomatic: A flexible trimmer for Illumina sequence data.  
638 *Bioinformatics*, 30(15), 2114–2120. <https://doi.org/10.1093/bioinformatics/btu170>
- 639 Buchfink, B., Xie, C., & Huson, D. H. (2015). Fast and sensitive protein alignment using DIAMOND.  
640 *Nature Methods*, 12(1). <https://doi.org/10.1038/nmeth.3176>

- 641 Camacho, C., Coulouris, G., Avagyan, V., Ma, N., Papadopoulos, J., Bealer, K., & Madden, T. L. (2009).  
642 BLAST+: Architecture and applications. *BMC Bioinformatics*, *10*, 1–9.  
643 <https://doi.org/10.1186/1471-2105-10-421>
- 644 Chakraborty, M., Baldwin-Brown, J. G., Long, A. D., & Emerson, J. J. (2016). Contiguous and accurate *de*  
645 *novo* assembly of metazoan genomes with modest long read coverage. *Nucleic Acids Research*, *44*(19),  
646 1–12. <https://doi.org/10.1093/nar/gkw654>
- 647 Chen, Y., Ye, W., Zhang, Y., & Xu, Y. (2015). High speed BLASTN: an accelerated MegaBLAST search  
648 tool. *Nucleic Acids Research*, *43*(16), 7762–7768. <https://doi.org/10.1093/nar/gkv784>
- 649 Chin, C.-S., Peluso, P., Sedlazeck, F. J., Nattestad, M., Concepcion, G. T., Clum, A., ... Schatz, M. C.  
650 (2016). Phased diploid genome assembly with single-molecule real-time sequencing. *Nature Methods*,  
651 *13*(12), 1050–1054. <https://doi.org/10.1038/nmeth.4035>
- 652 Cingolani, P., Platts, A., Wang, L. L., Coon, M., Nguyen, T., Wang, L., ... Ruden, D. M. (2012). A program  
653 for annotating and predicting the effects of single nucleotide polymorphisms, SnpEff. *Fly*, *6*(2), 80–92.  
654 <https://doi.org/10.4161/fly.19695>
- 655 Cui, L., Rui, C., Yang, D., Wang, Z., & Yuan, H. (2017). *De novo* transcriptome and expression profile  
656 analyses of the Asian corn borer (*Ostrinia furnacalis*) reveals relevant flubendiamide response genes.  
657 *BMC Genomics*, *18*(1), 20. <https://doi.org/10.1186/s12864-016-3431-6>
- 658 Czepak, C., Albernaz, K. C., Vivan, L. M., Guimarães, H. O., & Carvalhais, T. (2013). First reported  
659 occurrence of *Helicoverpa armigera* (Hübner) (Lepidoptera: Noctuidae) in Brazil. *Pesquisa*  
660 *Agropecuária Tropical*, *43*(1), 110–113. <https://doi.org/10.1590/S1983-40632013000100015>
- 661 Danecek, P., Auton, A., Abecasis, G., Albers, C. A., Banks, E., DePristo, M. A., ... Durbin, R. (2011). The  
662 variant call format and VCFtools. *Bioinformatics*, *27*(15), 2156–2158.  
663 <https://doi.org/10.1093/bioinformatics/btr330>
- 664 Ding, D., Liu, G., Hou, L., Gui, W., Chen, B., & Kang, L. (2018). Genetic variation in PTPN1 contributes to  
665 metabolic adaptation to high-altitude hypoxia in Tibetan migratory locusts. *Nature Communications*,  
666 *9*(1), 4991. <https://doi.org/10.1038/s41467-018-07529-8>
- 667 El-Gebali, S., Mistry, J., Bateman, A., Eddy, S. R., Luciani, A., Potter, S. C., ... Finn, R. D. (2019). The  
668 Pfam protein families database in 2019. *Nucleic Acids Research*, *47*(D1), D427–D432.  
669 <https://doi.org/10.1093/nar/gky995>

- 670 Fitt, G. (1989). The ecology of *Heliothis* species in relation to agroecosystems. *Annual Review of*  
671 *Entomology*, 34(1), 17–52. <https://doi.org/10.1146/annurev.ento.34.1.17>
- 672 Flynn, J. M., Hubley, R., Goubert, C., Rosen, J., Clark, A. G., Feschotte, C., & Smit, A. F. (2020).  
673 RepeatModeler2 for automated genomic discovery of transposable element families. *Proceedings of the*  
674 *National Academy of Sciences*, 117(17), 9451–9457. <https://doi.org/10.1073/pnas.1921046117>
- 675 Fritz, M. L., DeYonke, A. M., Papanicolaou, A., Micinski, S., Westbrook, J., & Gould, F. (2018).  
676 Contemporary evolution of a Lepidopteran species, *Heliothis virescens*, in response to modern  
677 agricultural practices. *Molecular Ecology*, 27(1), 167–181. <https://doi.org/10.1111/mec.14430>
- 678 Fu, L., Niu, B., Zhu, Z., Wu, S., & Li, W. (2012). CD-HIT: accelerated for clustering the next-generation  
679 sequencing data. *Bioinformatics*, 28(23), 3150–3152. <https://doi.org/10.1093/bioinformatics/bts565>
- 680 Gloss, A. D., Vassão, D. G., Hailey, A. L., Nelson Dittrich, A. C., Schramm, K., Reichelt, M., ... Whiteman,  
681 N. K. (2014). Evolution in an ancient detoxification pathway is coupled with a transition to herbivory  
682 in the drosophilidae. *Molecular Biology and Evolution*, 31(9), 2441–2456.  
683 <https://doi.org/10.1093/molbev/msu201>
- 684 Gouin, A., Bretaudeau, A., Nam, K., Gimenez, S., Aury, J.-M., Duvic, B., ... Fournier, P. (2017). Two  
685 genomes of highly polyphagous lepidopteran pests (*Spodoptera frugiperda*, Noctuidae) with different  
686 host-plant ranges. *Scientific Reports*, 7(1), 11816. <https://doi.org/10.1038/s41598-017-10461-4>
- 687 Guo, M., Du, L., Chen, Q., Feng, Y., Zhang, J., Zhang, X., ... Liu, Y. (2021). Odorant receptors for detecting  
688 flowering plant cues are functionally conserved across moths and butterflies. *Molecular Biology and*  
689 *Evolution*, 38(4), 1413–1427. <https://doi.org/10.1093/molbev/msaa300>
- 690 Hardwick, D. F. (1965). The corn earworm complex. *Memoirs of the Entomological Society of Canada*,  
691 97(S40), 5–247. <https://doi.org/10.4039/entm9740fv>
- 692 Hoff, K. J., Lange, S., Lomsadze, A., Borodovsky, M., & Stanke, M. (2016). BRAKER1: Unsupervised  
693 RNA-Seq-based genome annotation with GeneMark-ET and AUGUSTUS. *Bioinformatics*, 32(5), 767–  
694 769. <https://doi.org/10.1093/bioinformatics/btv661>
- 695 Holt, C., & Yandell, M. (2011). MAKER2: an annotation pipeline and genome-database management tool  
696 for second-generation genome projects. *BMC Bioinformatics*, 12(1), 491.  
697 <https://doi.org/10.1186/1471-2105-12-491>
- 698 Hsu, S. K., Belmouaden, C., Nolte, V., & Schlötterer, C. (2020). Parallel gene expression evolution in  
This article is protected by copyright. All rights reserved

- 699 natural and laboratory evolved populations. *Molecular Ecology*, 30(4), 1–11.  
700 <https://doi.org/10.1111/mec.15649>
- 701 Hubley, R., Finn, R. D., Clements, J., Eddy, S. R., Jones, T. A., Bao, W., ... Wheeler, T. J. (2016). The  
702 Dfam database of repetitive DNA families. *Nucleic Acids Research*, 44(D1), D81–D89.  
703 <https://doi.org/10.1093/nar/gkv1272>
- 704 Jiao, W.-B., Accinelli, G. G., Hartwig, B., Kiefer, C., Baker, D., Severing, E., ... Schneeberger, K. (2017).  
705 Improving and correcting the contiguity of long-read genome assemblies of three plant species using  
706 optical mapping and chromosome conformation capture data. *Genome Research*, 27(5), 778–786.  
707 <https://doi.org/10.1101/gr.213652.116>
- 708 Jin, L., Wang, J., Guan, F., Zhang, J. P., Yu, S., Liu, S. Y., ... Wu, Y. D. (2018). Dominant point mutation in  
709 a tetraspanin gene associated with field-evolved resistance of cotton bollworm to transgenic Bt cotton.  
710 *Proceedings of the National Academy of Sciences*, 115(46), 11760–11765.  
711 <https://doi.org/10.1073/pnas.1812138115>
- 712 Jin, M., Liao, C., Chakrabarty, S., Zheng, W., Wu, K. M., & Xiao, Y. T. (2019). Transcriptional response of  
713 ATP-binding cassette (ABC) transporters to insecticides in the cotton bollworm, *Helicoverpa armigera*.  
714 *Pesticide Biochemistry and Physiology*, 154, 46–59. <https://doi.org/10.1016/j.pestbp.2018.12.007>
- 715 Jones, C. M., Parry, H., Tay, W. T., Reynolds, D. R., & Chapman, J. W. (2019). Movement ecology of pest  
716 *Helicoverpa*: implications for ongoing spread. *Annual Review of Entomology*, 64(1), 1–19.  
717 <https://doi.org/10.1146/annurev-ento-011118-111959>
- 718 Joußen, N., Agnolet, S., Lorenz, S., Schöne, S. E., Ellinger, R., Schneider, B., & Heckel, D. G. (2012).  
719 Resistance of Australian *Helicoverpa armigera* to fenvalerate is due to the chimeric P450 enzyme  
720 *CYP337B3*. *Proceedings of the National Academy of Sciences of the United States of America*, 109(38),  
721 15206–15211. <https://doi.org/10.1073/pnas.1202047109>
- 722 Keightley, P. D., Pinharanda, A., Ness, R. W., Simpson, F., Dasmahapatra, K. K., Mallet, J., ... Jiggins, C. D.  
723 (2015). Estimation of the spontaneous mutation rate in *Heliconius melpomene*. *Molecular Biology and*  
724 *Evolution*, 32(1), 239–243. <https://doi.org/10.1093/molbev/msu302>
- 725 Kim, D., Langmead, B., & Salzberg, S. L. (2015). HISAT: A fast spliced aligner with low memory  
726 requirements. *Nature Methods*, 12(4), 357–360. <https://doi.org/10.1038/nmeth.3317>
- 727 Klai, K., Chenais, B., Zidi, M., Djebbi, S., Caruso, A., Denis, F., ... Mezghani K. M. (2020). Screening of

- 728 *Helicoverpa armigera* mobilome revealed transposable element insertions in insecticide resistance  
729 genes. *Insects*, 11(12), 879. <https://doi.org/10.3390/insects11120879>
- 730 Koren, S., Walenz, B. P., Berlin, K., Miller, J. R., Bergman, N. H., & Phillippy, A. M. (2017). Canu:  
731 Scalable and accurate long-read assembly via adaptive k-mer weighting and repeat separation. *Genome*  
732 *Research*, 27(5), 722–736. <https://doi.org/10.1101/gr.215087.116>
- 733 Krzywinski, M., Schein, J., Birol, I., Connors, J., Gascoyne, R., Horsman, D., ... Marra, M. A. (2009).  
734 Circos: An information aesthetic for comparative genomics. *Genome Research*, 19(9), 1639–1645.  
735 <https://doi.org/10.1101/gr.092759.109>
- 736 Kurtz, S., Phillippy, A., Delcher, A. L., Smoot, M., Shumway, M., Antonescu, C., & Salzberg, S. L. (2004).  
737 Versatile and open software for comparing large genomes. *Genome Biology*, 5(2). R12.  
738 <https://doi.org/10.1186/gb-2004-5-2-r12>
- 739 Lewis, T. E., Sillitoe, I., Dawson, N., Lam, S. D., Clarke, T., Lee, D., ... Lees, J. (2018). Gene3D: Extensive  
740 prediction of globular domains in proteins. *Nucleic Acids Research*, 46(D1), D435–D439.  
741 <https://doi.org/10.1093/nar/gkx1069>
- 742 Li, H. (2018). Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*, 34(18), 3094–3100.  
743 <https://doi.org/10.1093/bioinformatics/bty191>
- 744 Li, H., & Durbin, R. (2010). Fast and accurate long-read alignment with Burrows-Wheeler transform.  
745 *Bioinformatics*, 26(5), 589–595. <https://doi.org/10.1093/bioinformatics/btp698>
- 746 Li, H., & Durbin, R. (2011). Inference of human population history from individual whole-genome  
747 sequences. *Nature*, 475(7357), 493–496. <https://doi.org/10.1038/nature10231>
- 748 Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., ... Durbin, R. (2009). The Sequence  
749 Alignment/Map format and SAMtools. *Bioinformatics*, 25(16), 2078–2079.  
750 <https://doi.org/10.1093/bioinformatics/btp352>
- 751 Lomsadze, A. (2005). Gene identification in novel eukaryotic genomes by self-training algorithm. *Nucleic*  
752 *Acids Research*, 33(20), 6494–6506. <https://doi.org/10.1093/nar/gki937>
- 753 Low, W. Y., Feil, S. C., Ng, H. L., Gorman, M. A., Morton, C. J., Pyke, J., ... Batterham, P. (2010).  
754 Recognition and Detoxification of the Insecticide DDT by *Drosophila melanogaster* Glutathione  
755 S-Transferase D1. *Journal of Molecular Biology*, 399(3), 358–366.  
756 <https://doi.org/10.1016/j.jmb.2010.04.020>

This article is protected by copyright. All rights reserved

- 757 Marchler-Bauer, A., Bo, Y., Han, L., He, J., Lanczycki, C. J., Lu, S., ... Bryant, S. H. (2017).  
758 CDD/SPARCLE: Functional classification of proteins via subfamily domain architectures. *Nucleic*  
759 *Acids Research*, 45(D1), D200–D203. <https://doi.org/10.1093/nar/gkw1129>
- 760 McKenna, A., Hanna, M., Banks, E., Sivachenko, A., Cibulskis, K., Kernytsky, A., ... DePristo, M. A.  
761 (2010). The Genome Analysis Toolkit: A MapReduce framework for analyzing next-generation DNA  
762 sequencing data. *Genome Research*, 20(9), 1297–1303. <https://doi.org/10.1101/gr.107524.110>
- 763 Meng, G., Li, Y., Yang, C., & Liu, S. (2019). MitoZ: a toolkit for animal mitochondrial genome assembly,  
764 annotation and visualization. *Nucleic Acids Research*, 47(11), e63–e63.  
765 <https://doi.org/10.1093/nar/gkz173>
- 766 Merlin, C., Gegeer, R. J., & Reppert, S. M. (2009). Antennal circadian clocks coordinate sun compass  
767 orientation in migratory Monarch butterflies. *Science*, 325(5948), 1700–1704.  
768 <https://doi.org/10.1126/science.1176221>
- 769 Mi, H., Muruganujan, A., Ebert, D., Huang, X., & Thomas, P. D. (2019). PANTHER version 14: More  
770 genomes, a new PANTHER GO-slim and improvements in enrichment analysis tools. *Nucleic Acids*  
771 *Research*, 47(D1), D419–D426. <https://doi.org/10.1093/nar/gky1038>
- 772 Mitchell, A. L., Attwood, T. K., Babbitt, P. C., Blum, M., Bork, P., Bridge, A., ... Finn, R. D. (2019).  
773 InterPro in 2019: Improving coverage, classification and access to protein sequence annotations.  
774 *Nucleic Acids Research*, 47(D1), D351–D360. <https://doi.org/10.1093/nar/gky1100>
- 775 Montagné, N., de Fouchier, A., Newcomb, R. D., & Jacquin-Joly, E. (2015). Advances in the identification  
776 and characterization of olfactory receptors in insects. *Progress in Molecular Biology and Translational*  
777 *Science*, 130, 55–80. <https://doi.org/10.1016/bs.pmbts.2014.11.003>
- 778 Montero-Mendieta, S., Tan, K., Christmas, M. J., Olsson, A., Vilà, C., Wallberg, A., & Webster, M. T.  
779 (2019). The genomic basis of adaptation to high-altitude habitats in the eastern honey bee (*Apis cerana*).  
780 *Molecular Ecology*, 28(4), 746–760. <https://doi.org/10.1111/mec.14986>
- 781 Nibouche, S., Buès, R., Toubon, J.-F., & Poitout, S. (1998). Allozyme polymorphism in the cotton bollworm  
782 *Helicoverpa armigera* (Lepidoptera: Noctuidae): comparison of African and European populations.  
783 *Heredity*, 80(4), 438–445. <https://doi.org/10.1046/j.1365-2540.1998.00273.x>
- 784 O’Leary, N. A., Wright, M. W., Brister, J. R., Ciufu, S., Haddad, D., McVeigh, R., ... Pruitt, K. D. (2016).  
785 Reference sequence (RefSeq) database at NCBI: current status, taxonomic expansion, and functional

- 786 annotation. *Nucleic Acids Research*, 44(D1), D733–D745. <https://doi.org/10.1093/nar/gkv1189>
- 787 Patterson, N., Price, A. L., & Reich, D. (2006). Population structure and eigenanalysis. *PLoS Genetics*, 2(12),  
788 2074–2093. <https://doi.org/10.1371/journal.pgen.0020190>
- 789 Pearce, S. L., Clarke, D. F., East, P. D., Elfekih, S., Gordon, K. H. J., Jermin, L. S., ... Wu, Y. D. (2017).  
790 Genomic innovations, transcriptional plasticity and gene loss underlying the evolution and divergence  
791 of two highly polyphagous and invasive *Helicoverpa* pest species. *BMC Biology*, 15(1), 63.  
792 <https://doi.org/10.1186/s12915-017-0402-6>
- 793 Perteu, M., Perteu, G. M., Antonescu, C. M., Chang, T.-C., Mendell, J. T., & Salzberg, S. L. (2015).  
794 StringTie enables improved reconstruction of a transcriptome from RNA-seq reads. *Nature*  
795 *Biotechnology*, 33(3), 290–295. <https://doi.org/10.1038/nbt.3122>
- 796 Presgraves, D. C. (2018). Evaluating genomic signatures of “the large X-effect” during complex speciation.  
797 *Molecular Ecology*, 27(19), 3822–3830. <https://doi.org/10.1111/mec.14777>
- 798 Quan, Q., Hu, X., Pan, B., Zeng, B., Wu, N., Fang, G., ... Zhan, S. (2019). Draft genome of the cotton aphid  
799 *Aphis gossypii*. *Insect Biochemistry and Molecular Biology*, 105, 25–32.  
800 <https://doi.org/10.1016/j.ibmb.2018.12.007>
- 801 Raj, A., Stephens, M., & Pritchard, J. K. (2014). FastSTRUCTURE: Variational inference of population  
802 structure in large SNP data sets. *Genetics*, 197(2), 573–589.  
803 <https://doi.org/10.1534/genetics.114.164350>
- 804 Rane, R. V., Ghodke, A. B., Hoffmann, A. A., Edwards, O. R., Walsh, T. K., & Oakeshott, J. G. (2019).  
805 Detoxifying enzyme complements and host use phenotypes in 160 insect species. *Current Opinion in*  
806 *Insect Science*, 31, 131–138. <https://doi.org/10.1016/j.cois.2018.12.008>
- 807 Rastas, P. (2017). Lep-MAP3: Robust linkage mapping even for low-coverage whole genome sequencing  
808 data. *Bioinformatics*, 33(23), 3726–3732. <https://doi.org/10.1093/bioinformatics/btx494>
- 809 Roach, M. J., Schmidt, S. A., & Borneman, A. R. (2018). Purge Haplotigs: allelic contig reassignment for  
810 third-gen diploid genome assemblies. *BMC Bioinformatics*, 19(1), 460.  
811 <https://doi.org/10.1186/s12859-018-2485-7>
- 812 Sedlazeck, F. J., Lee, H., Darby, C. A., & Schatz, M. C. (2018). Piercing the dark matter: bioinformatics of  
813 long-range sequencing and mapping. *Nature Reviews Genetics*, 19(6), 329–346.  
814 <https://doi.org/10.1038/s41576-018-0003-4>

- 815 Sedlazeck, F. J., Rescheneder, P., Smolka, M., Fang, H., Nattestad, M., von Haeseler, A., & Schatz, M. C.  
816 (2018). Accurate detection of complex structural variations using single-molecule sequencing. *Nature*  
817 *Methods*, 15(6), 461–468. <https://doi.org/10.1038/s41592-018-0001-7>
- 818 Shen, W., Le, S., Li, Y., & Hu, F. (2016). SeqKit: A cross-platform and ultrafast toolkit for FASTA/Q file  
819 manipulation. *PLoS ONE*, 11(10), e0163962. <https://doi.org/10.1371/journal.pone.0163962>
- 820 Simão, F. A., Waterhouse, R. M., Ioannidis, P., Kriventseva, E. V., & Zdobnov, E. M. (2015). BUSCO:  
821 Assessing genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics*,  
822 31(19), 3210–3212. <https://doi.org/10.1093/bioinformatics/btv351>
- 823 Song, S. V., Anderson, C. J., Good, R. T., Leslie, S., Wu, Y. D., Oakeshott, J. G., & Robin, C. (2018).  
824 Population differentiation between Australian and Chinese *Helicoverpa armigera* occurs in distinct  
825 blocks on the Z-chromosome. *Bulletin of Entomological Research*, 108(6), 817–830.  
826 <https://doi.org/10.1017/S0007485318000081>
- 827 Stanke, M., Steinkamp, R., Waack, S., & Morgenstern, B. (2004). AUGUSTUS: a web server for gene  
828 finding in eukaryotes. *Nucleic Acids Research*, 32, W309–W312. <https://doi.org/10.1093/nar/gkh379>
- 829 Tabashnik, B. E., & Carrière, Y. (2017). Surge in insect resistance to transgenic crops and prospects for  
830 sustainability. *Nature Biotechnology*, 35(10), 926–935. <https://doi.org/10.1038/nbt.3974>
- 831 Tang, H. B., Zhang, X., Miao, C., Zhang, J., Ming, R., Schnable, J. C., ... Lu, J. (2015). ALLMAPS: robust  
832 scaffold ordering based on multiple maps. *Genome Biology*, 16(1), 3.  
833 <https://doi.org/10.1186/s13059-014-0573-1>
- 834 Tarailo-Graovac, M., & Chen, N. (2009). Using RepeatMasker to identify repetitive elements in genomic  
835 sequences. *Current Protocols in Bioinformatics*, 25(1), 1–14.  
836 <https://doi.org/10.1002/0471250953.bi0410s25>
- 837 Tay, W. T., Behere, G. T., Batterham, P., & Heckel, D. G. (2010). Generation of microsatellite repeat  
838 families by RTE retrotransposons in lepidopteran genomes. *BMC Evolutionary Biology*, 10(1), 144.  
839 <https://doi.org/10.1186/1471-2148-10-144>
- 840 Tay, W. T., & Gordon, K. H. J. (2019). Going global – genomic insights into insect invasions. *Current*  
841 *Opinion in Insect Science*, 31, 123–130. <https://doi.org/10.1016/j.cois.2018.12.002>
- 842 Tay, W. T., Soria, M. F., Walsh, T., Thomazoni, D., Silvie, P., Behere, G. T., ... Downes, S. (2013). A brave  
843 New World for an Old World pest: *Helicoverpa armigera* (Lepidoptera: Noctuidae) in Brazil. *PLoS*

- 844 *ONE*, 8(11), e80134. <https://doi.org/10.1371/journal.pone.0080134>
- 845 Tay, W. T., Walsh, T. K., Downes, S., Anderson, C. J., Jermin, L. S., Wong, T. K. F., ... Gordon, K. H. J.  
846 (2017). Mitochondrial DNA and trade data support multiple origins of *Helicoverpa armigera*  
847 (Lepidoptera, Noctuidae) in Brazil. *Scientific Reports*, 7(1), 45302. <https://doi.org/10.1038/srep45302>
- 848 Teets, N. M., Peyton, J. T., Colinet, H., Renault, D., Kelley, J. L., Kawarasaki, Y., ... Denlinger, D. L.  
849 (2012). Gene expression changes governing extreme dehydration tolerance in an Antarctic insect.  
850 *Proceedings of the National Academy of Sciences*, 109(50), 20744–20749.  
851 <https://doi.org/10.1073/pnas.1218661109>
- 852 Terhorst, J., Kamm, J. A., & Song, Y. S. (2017). Robust and scalable inference of population history from  
853 hundreds of unphased whole genomes. *Nature Genetics*, 49(2), 303–309.  
854 <https://doi.org/10.1038/ng.3748>
- 855 Tomioka, K., & Matsumoto, A. (2015). Circadian molecular clockworks in non-model insects. *Current*  
856 *Opinion in Insect Science*, 7, 58–64. <https://doi.org/10.1016/j.cois.2014.12.006>
- 857 Valencia-Montoya, W. A., Elfekih, S., North, H. L., Meier, J. I., Warren, I. A., Tay, W. T., ... Jiggins, C. D.  
858 (2020). Adaptive introgression across semipermeable species boundaries between local *Helicoverpa*  
859 *zea* and invasive *Helicoverpa armigera* moths. *Molecular Biology and Evolution*, 37(9), 2568–2583.  
860 <https://doi.org/10.1093/molbev/msaa108>
- 861 Van't Hof, A. E., Nguyen, P., Dalíková, M., Edmonds, N., Marec, F., & Saccheri, I. J. (2013). Linkage map  
862 of the peppered moth, *Biston betularia* (Lepidoptera, Geometridae): A model of industrial melanism.  
863 *Heredity*, 110(3), 283–295. <https://doi.org/10.1038/hdy.2012.84>
- 864 Van Leeuwen, T., Dermauw, W., Mavridis, K., & Vontas, J. (2020). Significance and interpretation of  
865 molecular diagnostics for insecticide resistance management of agricultural pests. *Current Opinion in*  
866 *Insect Science*, 39, 69–76. <https://doi.org/10.1016/j.cois.2020.03.006>
- 867 Vurture, G. W., Sedlazeck, F. J., Nattestad, M., Underwood, C. J., Fang, H., Gurtowski, J., & Schatz, M. C.  
868 (2017). GenomeScope: fast reference-free genome profiling from short reads. *Bioinformatics*, 33(14),  
869 2202–2204. <https://doi.org/10.1093/bioinformatics/btx153>
- 870 Walker, B. J., Abeel, T., Shea, T., Priest, M., Abouelliel, A., Sakthikumar, S., ... Earl, A. M. (2014). Pilon:  
871 An integrated tool for comprehensive microbial variant detection and genome assembly improvement.  
872 *PLoS ONE*, 9(11). <https://doi.org/10.1371/journal.pone.0112963>

- 873 Walsh, T. K., Joussem, N., Tian, K., McGaughran, A., Anderson, C. J., Qiu, X.H., ... Heckel, D. G. (2018).  
874 Multiple recombination events between two cytochrome P450 loci contribute to global pyrethroid  
875 resistance in *Helicoverpa armigera*. *PLoS ONE*, 13(11), e0197760.  
876 <https://doi.org/10.1371/journal.pone.0197760>
- 877 Walsh, T. K., Perera, O., Anderson, C.J., Gordon, K., Czapak, C., McGaughran, A., ... Tay, W. T. (2019).  
878 Mitochondrial DNA genomes of five major *Helicoverpa* pest species from the Old and New Worlds  
879 (Lepidoptera: Noctuidae). *Ecology and Evolution*, 9(5), 2933–2944. <https://doi.org/10.1002/ece3.4971>
- 880 Wang, Y., Kim, K. S., Guo, W., Li, Q., Zhang, Y., Wang, Z., & Coates, B. S. (2017). Introgression between  
881 divergent corn borer species in a region of sympatry: Implications on the evolution and adaptation of  
882 pest arthropods. *Molecular Ecology*, 26(24), 6892–6907. <https://doi.org/10.1111/mec.14387>
- 883 Wang, Y. P., Tang, H. B., DeBarry, J. D., Tan, X., Li, J., Wang, X. Y., ... Paterson, A. H. (2012). MCSanX:  
884 a toolkit for detection and evolutionary analysis of gene synteny and collinearity. *Nucleic Acids  
885 Research*, 40(7). <https://doi.org/10.1093/nar/gkr1293>
- 886 Weeks, A. R., Endersby, N. M., Lange, C. L., Lowe, A., Zalucki, M. P., & Hoffmann, A. A. (2010). Genetic  
887 variation among *Helicoverpa armigera* populations as assessed by microsatellites: A cautionary tale  
888 about accurate allele scoring. *Bulletin of Entomological Research*, 100(4), 445–450.  
889 <https://doi.org/10.1017/S0007485309990460>
- 890 Wilson, D., Pethica, R., Zhou, Y., Talbot, C., Vogel, C., Madera, M., ... Gough, J. (2009). SUPERFAMILY  
891 – Sophisticated comparative genomics, data mining, visualization and phylogeny. *Nucleic Acids  
892 Research*, 37, 380–386. <https://doi.org/10.1093/nar/gkn762>
- 893 Wolf, J. B. W., & Ellegren, H. (2017). Making sense of genomic islands of differentiation in light of  
894 speciation. *Nature Reviews Genetics*, 18(2), 87–100. <https://doi.org/10.1038/nrg.2016.133>
- 895 Wu, K. M., & Guo, Y. Y. (2005). The evolution of cotton pest management practices in China. *Annual  
896 Review of Entomology*, 50(1), 31–52. <https://doi.org/10.1146/annurev.ento.50.071803.130349>
- 897 Wu, N., Zhang, S., Li, X., Cao, Y., Liu, X., Wang, Q., ... Zhan, S. (2019). Fall webworm genomes yield  
898 insights into rapid adaptation of invasive species. *Nature Ecology & Evolution*, 3(1), 105–115.  
899 <https://doi.org/10.1038/s41559-018-0746-5>
- 900 Xia, Q. Y., Guo, Y., Zhang, Z., Li, D., Xuan, Z., Li, Z., ... Wang, J. (2009). Complete resequencing of 40  
901 genomes reveals domestication events and genes in silkworm (*Bombyx*). *Science*, 326(5951), 433–436.

- 902 <https://doi.org/10.1126/science.1176620>
- 903 Xu, W. (2020). How do moth and butterfly taste?—Molecular basis of gustatory receptors in Lepidoptera.  
904 *Insect Science*, 27(6), 1148–1157. <https://doi.org/10.1111/1744-7917.12718>
- 905 Yamamoto, K., Nohata, J., Kadono-Okuda, K., Narukawa, J., Sasanuma, M., Sasanuma, S., ... Mita, K.  
906 (2008). A BAC-based integrated linkage map of the silkworm *Bombyx mori*. *Genome Biology*, 9(1),  
907 R21. <https://doi.org/10.1186/gb-2008-9-1-r21>
- 908 Yang, Y. H., Li, Y., & Wu, Y. D. (2013). Current status of insecticide resistance in *Helicoverpa armigera*  
909 after 15 years of Bt cotton planting in China. *Journal of Economic Entomology*, 106(1), 375–381.  
910 <https://doi.org/10.1603/EC12286>
- 911 Yu, G. C., Wang, L.-G., Han, Y., & He, Q.-Y. (2012). clusterProfiler: an R package for comparing biological  
912 themes among gene clusters. *OMICS: A Journal of Integrative Biology*, 16(5), 284–287.  
913 <https://doi.org/10.1089/omi.2011.0118>
- 914 Zhu, J., Jiang, F., Wang, X., Yang, P., Bao, Y., Zhao, W., ... Cui, F. (2017). Genome sequence of the small  
915 brown planthopper, *Laodelphax striatellus*. *GigaScience*, 6(12), 1–12. [https://doi:](https://doi.org/10.1093/gigascience/gix109)  
916 [10.1093/gigascience/gix109](https://doi.org/10.1093/gigascience/gix109)

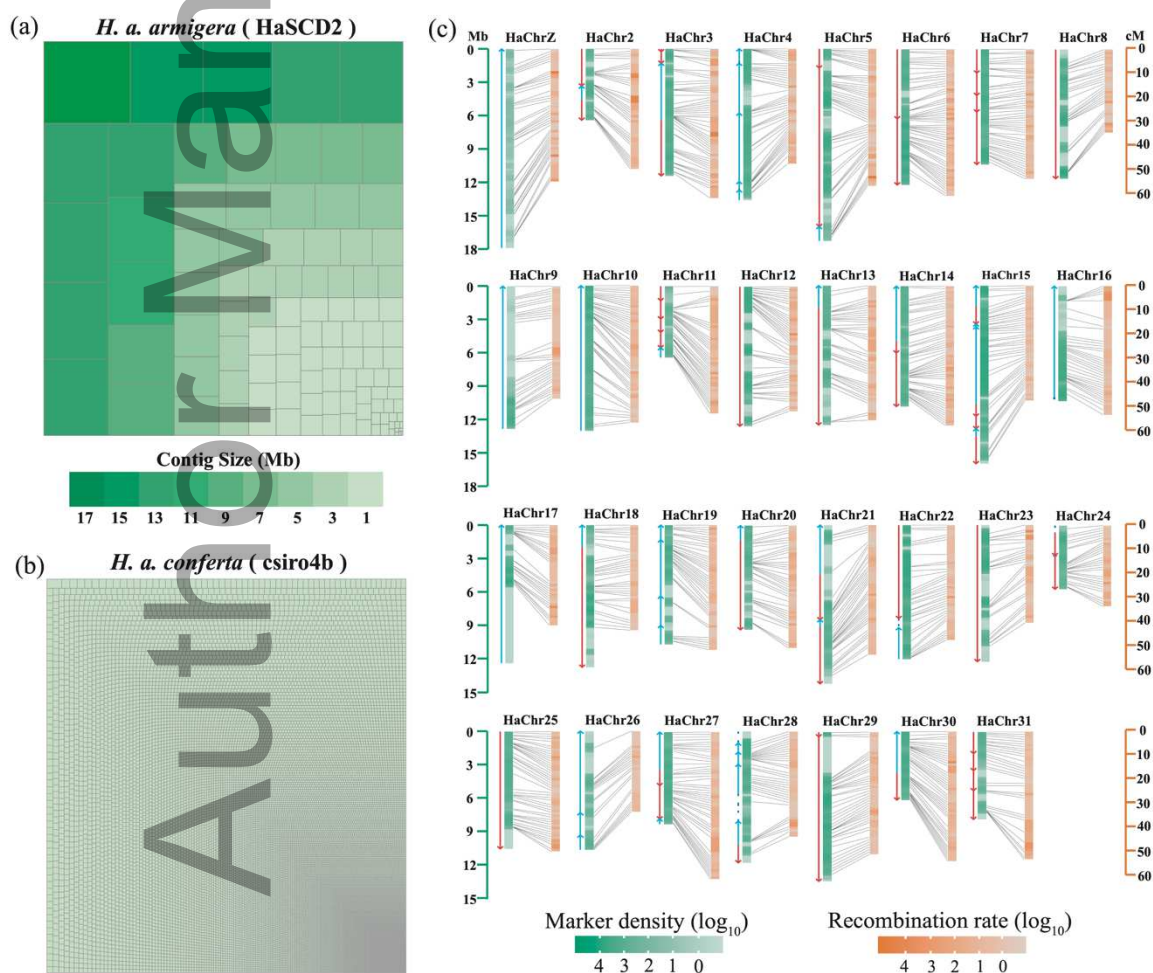
917

918 **Table 1.** Summary statistics of HaSCD2 and csiro4b assemblies as well as their corresponding gene  
919 annotation.

Species	<i>H. a. armigera</i>	<i>H. a. conferta</i>
<b>Genome assembly</b>	HaSCD2	csiro4b
Assembly size (Mb)	356.67	337.07
Number of scaffolds	41	997
Number of contigs	106	24,547
Longest scaffold (Mb)	17.88	6.15
Longest contig (Mb)	17.88	0.31
N50 scaffolds length (Mb)	12.27	1.00
N50 contig length	7.34 Mb	23.48 kb
N90 scaffolds length (Mb)	7.86	0.18

N90 contig length	1.65 Mb	6.21 kb
GC (%)	36.53	36.13
<b>Gene annotation</b>		
Protein-coding genes	18,668	17,086
Mean protein length (aa)	475.55	442.79
Mean gene length (bp)	6,776.03	9,364.52
Exons per gene	5.97	6.05
Exon (%)	7.48	6.75
Mean exon length	239.40	219.96
Intron (%)	23.93	35.48
Mean intron length	919.69	1,368.88

920



921

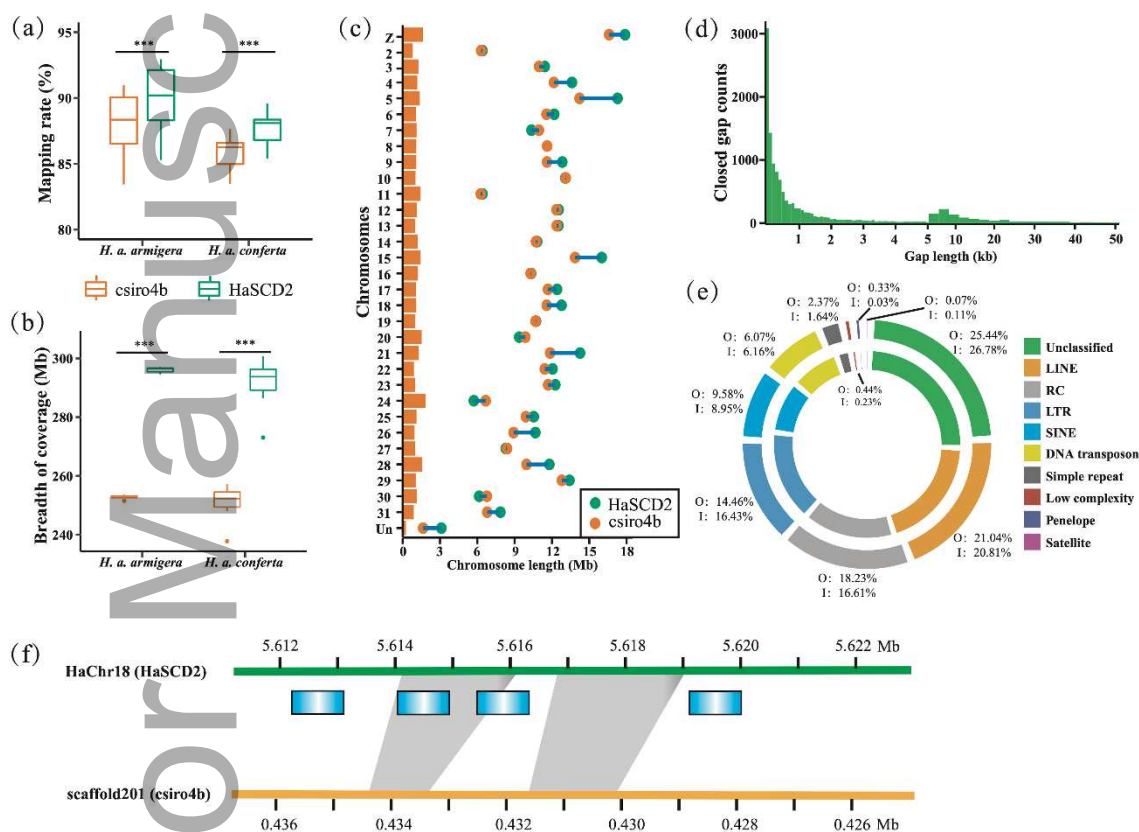
922

923

**Figure 1.** Genome assembly and scaffolding of the HaSCD2. (a) - (b) The differences of contigs length between HaSCD2 and csiro4b assemblies. The proportions and levels of contig lengths are represented by rectangular areas. This article is protected by copyright. All rights reserved

924 and filled colors. (c) Assigning contigs into chromosomes with the guide of high-resolution genetic map. Markers  
 925 with identical recombination distance are collapsed into bins to succinctly show their congruent relationship, and  
 926 the first markers were linked with the corresponding recombination distance. Graduated colors filled in two  
 927 rectangles are determined with log<sub>10</sub>-transformed marker density and recombination distances along 31  
 928 chromosomes. The 94 anchored contigs are represented by red and blue lines and their arrows indicate their  
 929 orientation, and dark blue dotted lines indicate the contigs not oriented.

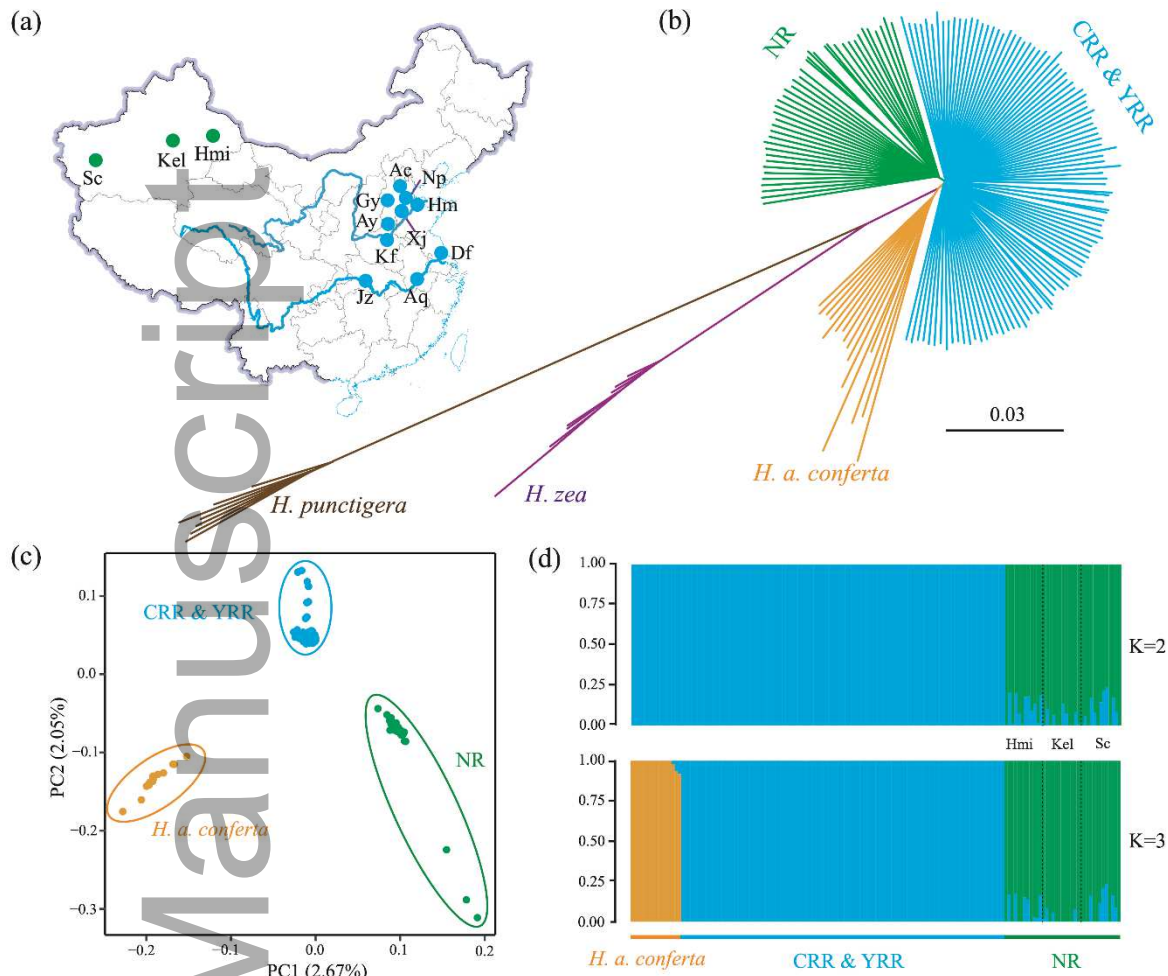
930



931

932 **Figure 2.** Genome comparisons between HaSCD2 and csiro4b assemblies. (a) - (b) Genomic mapping rate and  
 933 breadth of coverage (more than 2-fold coverage) against HaSCD2 and csiro4b assemblies. (c) Comparisons of  
 934 chromosomes size between HaSCD2 and csiro4b assemblies. (d) The distribution of the sizes of closed gaps in the  
 935 csiro4b assembly. (e) Proportions of major repeat categories distributed in the HaSCD2 assembly. The outer circle  
 936 represents the repeat in whole HaSCD2 genome and the inner circle represents the repeat in assigned N-gap  
 937 regions. LINE: long interspersed nuclear elements; RC: rolling-circle transposon; LTR: long terminal repeats;  
 938 SINE: short interspersed nuclear elements. (f) A case of assigned N-gap affecting the completeness of two  
 939 odorant-binding proteins (OBPs). Blue rectangles represent odorant binding proteins, grey shadows represent  
 940 mapping relationships between two assemblies.

This article is protected by copyright. All rights reserved



942

943

944

945

946

947

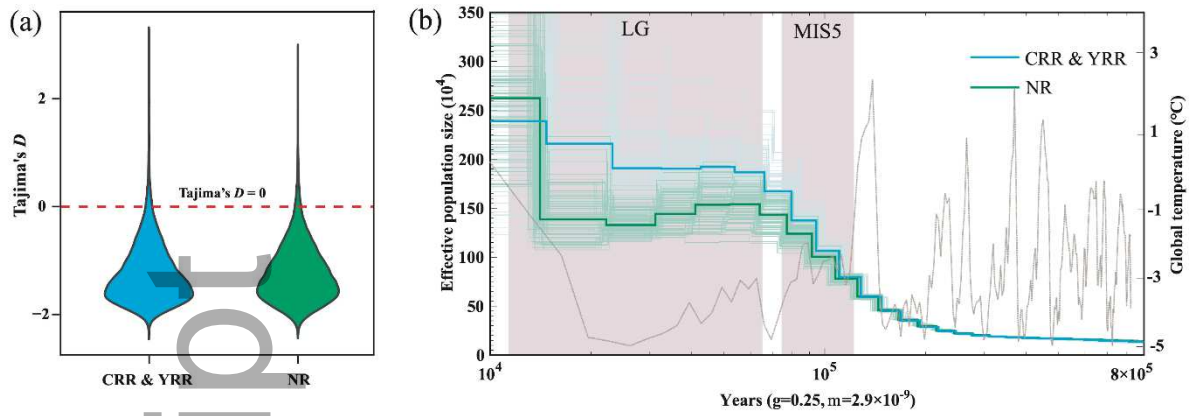
948

949

950

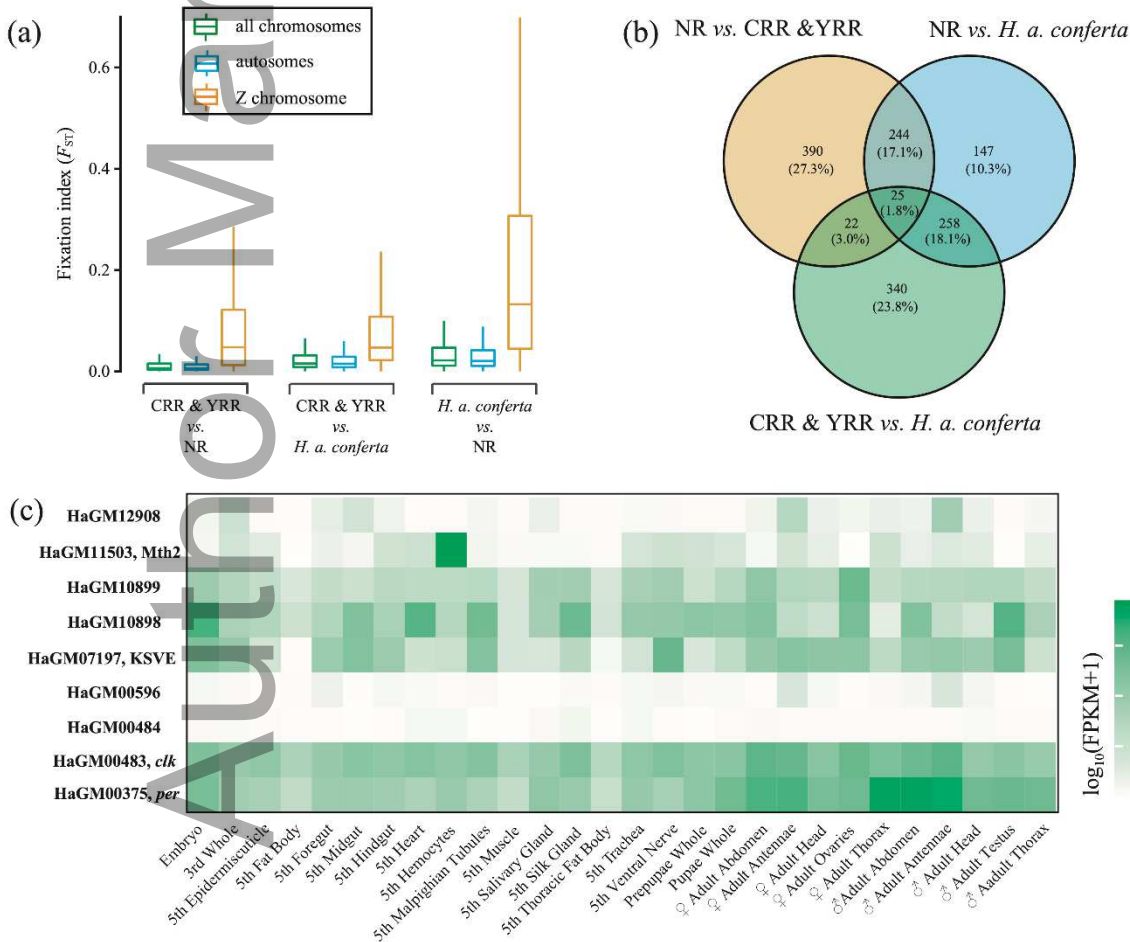
950

**Figure 3.** Geographic distribution and genetic structure of collected cotton bollworm. (a) Schematic map denoting the sampling localities in mainland China. Blue solid circles represent localities in the Yellow River Region (YRR) and the Changjiang River Region (CRR), green solid circles represent localities in the Northwestern Region (NR), two blue curves represent Yellow River (top) and Changjiang River (bottom), respectively. (b) Neighbor-joining (NJ) tree of cotton bollworm inferred from whole-genome SNPs using *H. punctigera* as outgroup. (c) Principle components analysis, PC1=2.67% and PC2=2.05% (Tracy-Widom statistics,  $P < 4.91e-291$ ). (d) Population admixture plot with K equals to 2 (top panel) and 3 (bottom panel).



951  
 952 **Figure 4.** Inference of demographic history of cotton bollworm. (a) Violin plot shows the genome-wide  
 953 distribution of Tajima's  $D$  values in CRR & YRR and NR lineages. The red dotted line indicates Tajima's  $D$   
 954 value of zero. (b) The historical effective population size of representative individuals in CRR & YRR and NR lineages.  
 955 The light blue and green lines represent 100 bootstraps. The dotted grey line represents the historical global  
 956 temperature. LG, last glaciation; MIS5, marine isotope stage 5.

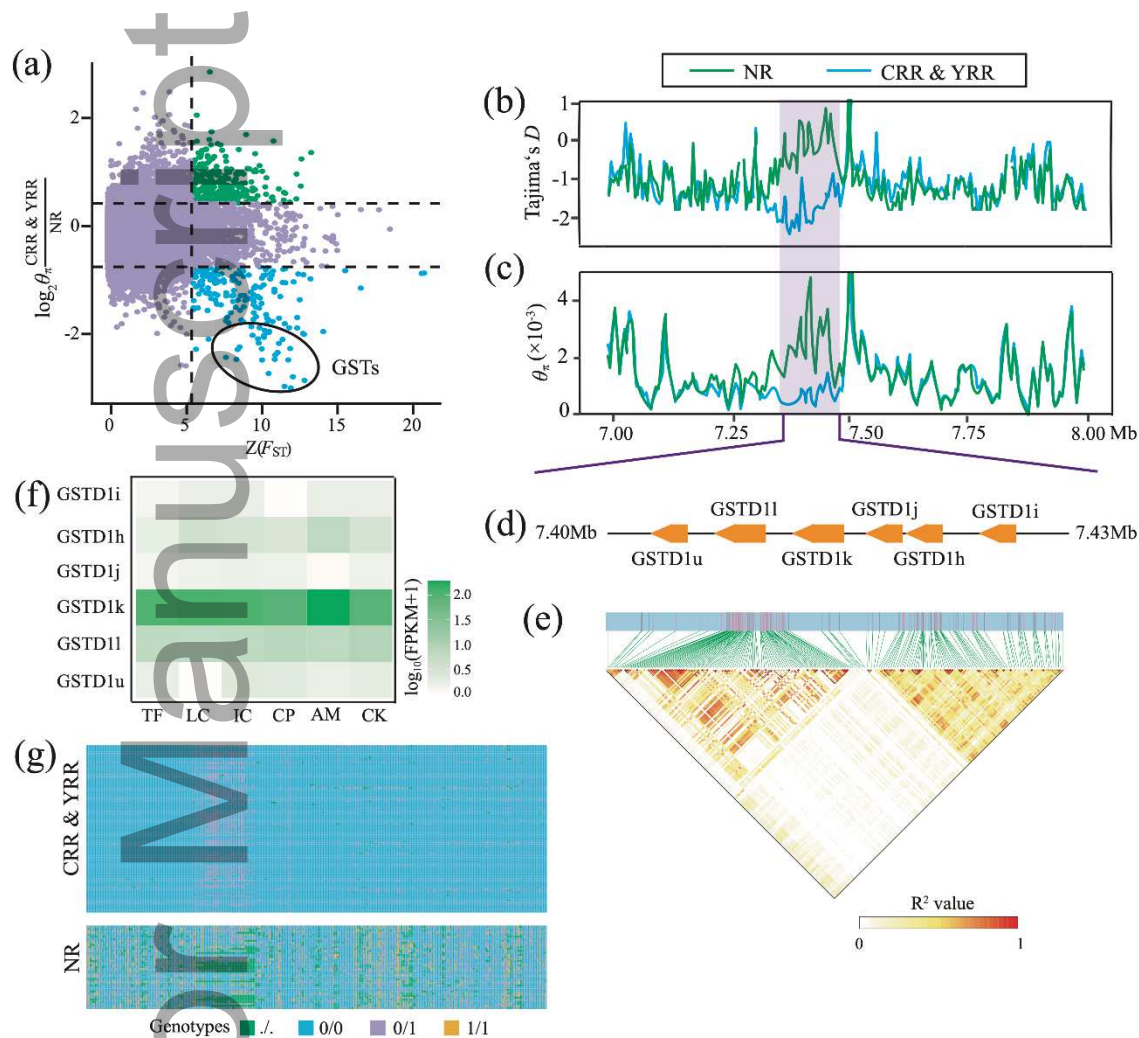
957



958  
 959 **Figure 5.** Divergence patterns among different cotton bollworm lineages. (a)  $F_{ST}$  values in the sex chromosome  
 960 and autosomes; (b) Venn diagram shows the number of shared genes with highest 1%  $F_{ST}$  values in different  
 This article is protected by copyright. All rights reserved

961 comparisons. (c) Expression levels of nine divergent genes in different tissues and developmental stages of cotton  
 962 bollworm.

963



964

965 **Figure 6.** Selection signatures between CRR & YRR and NR lineages. (a) Distribution of  $\log_2$ -transformed  $\theta_{\pi}$   
 966 ratio (Y axis) and Z-transformed  $F_{ST}$  (X axis) between CRR & YRR and NR lineages. Black lines represent their  
 967 highest 1% thresholds. Green dots represent genomic regions under selection for the NR lineage and blue dots for  
 968 the CRR & YRR lineage. (b) - (c) Tajima's  $D$  and  $\theta_{\pi}$  values around selection signatures in HaChr7. Green line  
 969 represents NR lineage and blue line represents CRR & YRR lineage. (d) Six glutathione S-transferase (GST)  
 970 genes annotated in the selection signatures of HaChr7. (e) Patterns of LD blocks around these six GST genes in  
 971 the CRR & YRR lineage. (f) Expressions of six GST genes in *H. armigera* larvae fed on artificial diet  
 972 supplemented with different insecticides. AM: abamectin; CP: chlorpyrifos; IC: indoxacarb; LC:  
 973 lambda-cyhalothrin; TF: tebufenozide; CK: control. (g) Genotypes of loci distributed around these six GST genes.

This article is protected by copyright. All rights reserved