

# Natural courtship song variation caused by an intronic retroelement in an ion channel gene

Yun Ding<sup>1</sup>, Augusto Berrocal<sup>1†</sup>, Tomoko Morita<sup>1</sup>, Kit D. Longden<sup>1</sup> & David L. Stern<sup>1</sup>

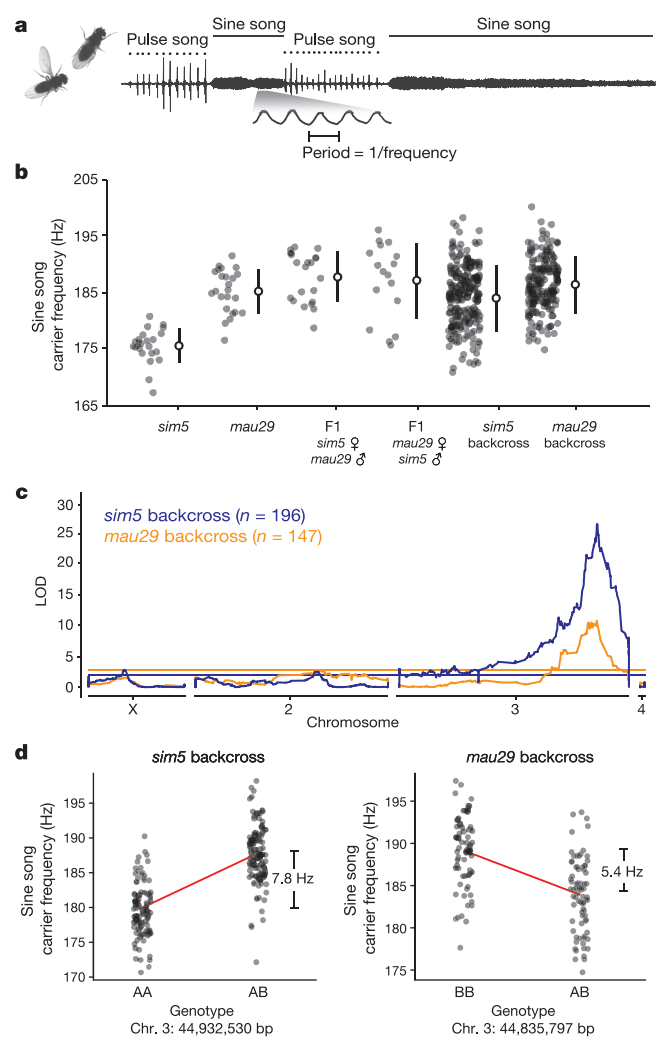
Animal species display enormous variation for innate behaviours, but little is known about how this diversity arose. Here, using an unbiased genetic approach, we map a courtship song difference between wild isolates of *Drosophila simulans* and *Drosophila mauritiana* to a 966 base pair region within the *slowpoke* (*slo*) locus, which encodes a calcium-activated potassium channel<sup>1</sup>. Using the reciprocal hemizyosity test<sup>2</sup>, we confirm that *slo* is the causal locus and resolve the causal mutation to the evolutionarily recent insertion of a retroelement in a *slo* intron within *D. simulans*. Targeted deletion of this retroelement reverts the song phenotype and alters *slo* splicing. Like many ion channel genes, *slo* is expressed widely in the nervous system and influences a variety of behaviours<sup>3,4</sup>; *slo*-null males sing little song with severely disrupted features. By contrast, the natural variant of *slo* alters a specific component of courtship song, illustrating that regulatory evolution of a highly pleiotropic ion channel gene can cause modular changes in behaviour.

During courtship in *Drosophila* species, males vibrate their wings to produce a ‘song’ that attracts females<sup>5</sup>. Courtship song is relatively easy to quantify<sup>6</sup> and varies widely between species<sup>7</sup>, making song an excellent system for genetic studies. However, despite decades of work<sup>8–11</sup>, no causative loci for natural variation in courtship song have been identified definitively<sup>12</sup>. In particular, candidate gene approaches have failed to identify loci contributing to natural variation<sup>13</sup>. We have therefore taken an unbiased, whole-genome approach to identify loci underlying natural variation in courtship song.

Male flies in the *D. melanogaster* species subgroup, which includes the species studied here, produce courtship song that often contains two components: trains of continuous, approximately sinusoidal sound at a certain carrier frequency, called ‘sine song’, and a series of pulses separated by a characteristic interval, called ‘pulse song’<sup>5,7</sup> (Fig. 1a and Supplementary Video 1). *D. simulans* and *D. mauritiana* diverged about 240 thousand years ago<sup>14</sup> and many features of their songs have changed. For example, the carrier frequency of sine song differs by 9.7 Hz between two wild-type isolates of these species, *sim5* and *mau29* (Fig. 1b and Supplementary Audios 1 and 2). We performed quantitative trait locus (QTL) mapping of courtship song traits between *sim5* and *mau29* using a high-throughput song phenotyping platform<sup>6</sup> and multiplexed shotgun genotyping<sup>15</sup>. The F1 hybrids and backcross progeny produced sine song with a frequency similar to *mau29*, indicating that the *mau29* allele(s) is largely dominant over the *sim5* allele(s) (Fig. 1b). In both backcrosses, we detected a single significant QTL at about 44.9 Mb on chromosome 3 (Fig. 1c) and the QTL explains most of the difference in sine frequency (Fig. 1d). We also identified one QTL for pulse song carrier frequency (Extended Data Fig. 1a, b) and two QTLs for inter-pulse interval (Extended Data Fig. 1c, d). The QTLs for these traits are located at different positions, indicating that different song features are genetically separable, consistent with the genetic modularity observed for other evolved behaviours<sup>16,17</sup>.

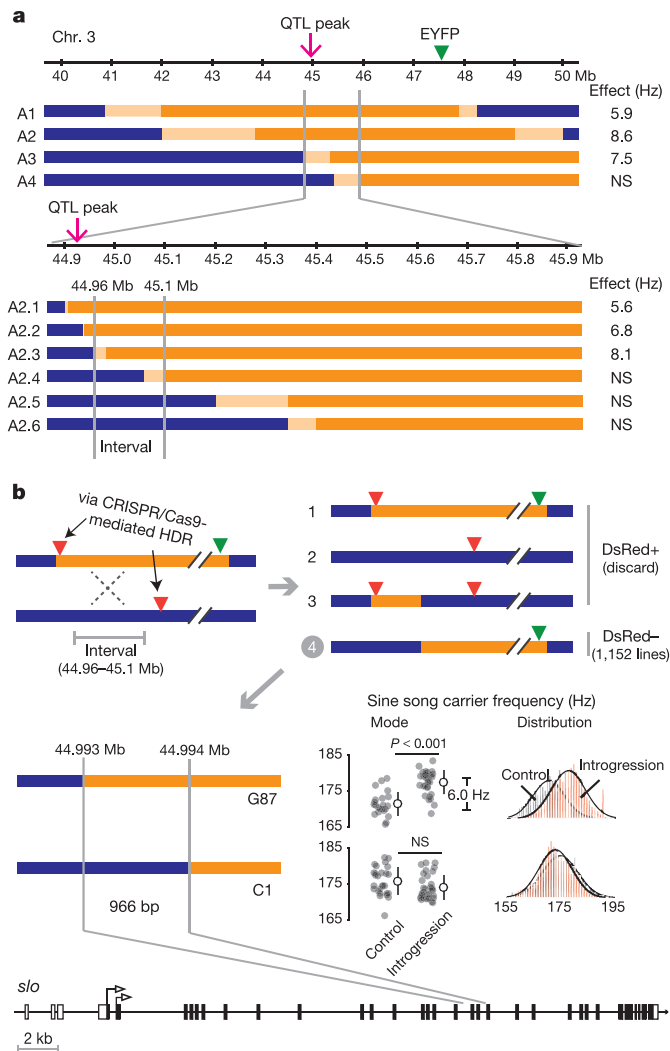
To validate and fine-map the sine frequency QTL, we produced *D. mauritiana white* (*mauW*) strains with randomly inserted transposable

elements carrying 3XP3::EYFP marker to facilitate targeted introgression. Using a marker located at 47.75 Mb on chromosome 3, we introgressed *D. mauritiana* DNA near the QTL into *sim5* (Fig. 2a). Of the many lines screened, several critical lines delimited the causal locus within an interval of about 1 Mb (44.9–45.9 Mb) that includes the



**Figure 1 | QTL analysis of sine song frequency difference between *sim5* and *mau29*.** **a**, Illustration of pulse and sine song. **b**, Sine song frequency in parental strains, F1 hybrids, and backcross males. Mean  $\pm$  s.d. **c**, QTL map of *sim5* (blue) and *mau29* (orange) backcross. LOD, logarithm of the odds. Horizontal lines mark  $P = 0.01$ . **d**, Effect plots of the chromosome 3 QTL from each backcross. A, *sim5* allele; B, *mau29* allele. Red lines connect means for each genotype.

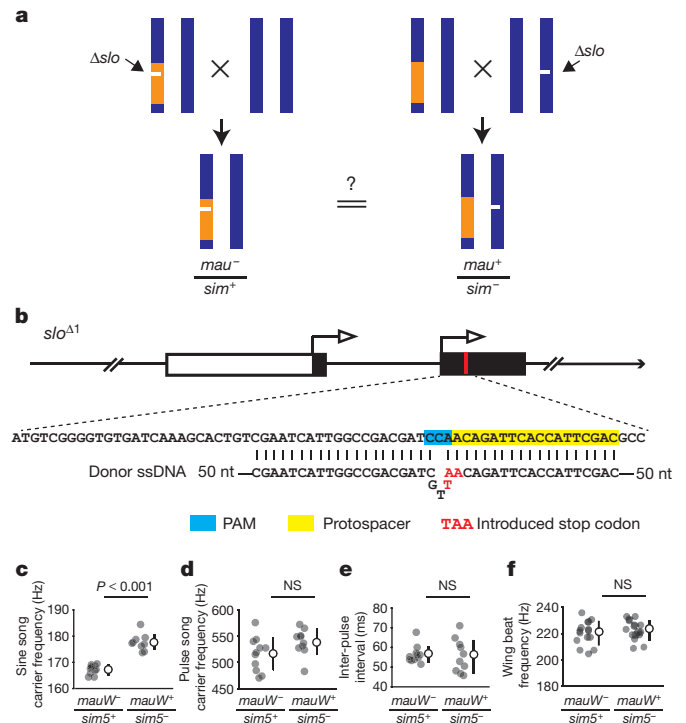
<sup>1</sup>Janelia Research Campus, Howard Hughes Medical Institute, Ashburn, Virginia 20147, USA. <sup>†</sup>Present address: Department of Molecular and Cell Biology, University of California, Berkeley, California 94720-3200, USA.



**Figure 2 | Fine-scale mapping identifies *slo* as the candidate causal locus.** **a**, Introgression of a *mauW* EYFP marker (green triangle) into a *sim5* background identifies a 140 kb region that causes a sine song frequency difference. Bar colour denotes species identity of DNA: blue, *sim5* DNA; orange, *mauW* DNA; light orange, breakpoint region. **b**, High-resolution mapping resolves the causal mutations to a 966 bp region of the *slo* locus. Upper panel, schematic of the targeted mapping strategy. Green and red triangles represent EYFP and DsRed markers. The products of either the absence (1, 2) or presence (3, 4) of recombination between DsRed markers are indicated. Middle panel, sine song frequency phenotypes (mean  $\pm$  s.d.) of the two recombinant lines defining the minimal interval.  $P$  value by one-way ANOVA; NS, non-significant. Lower panel, diagram of *slo* gene structure showing minimal mapping interval. Black and open boxes indicate coding and non-coding exons, respectively, and open triangles mark the two alternative start codons. Genotyping and phenotyping data provided in Extended Data Fig. 2.

QTL peak (Fig. 2a and Extended Data Fig. 2a). This result validated our QTL map and suggests that *mau29* and *mauW* share the same genetic cause for higher sine frequency.

The *D. simulans* and *D. mauritiana* genomes differ at approximately one in every hundred base pairs and most differences are presumably irrelevant to song evolution. Therefore, unlike mapping of laboratory-induced mutations, we required high resolution to localize the causal nucleotides. To gain further resolution, we generated 500 additional introgression lines, defining a causal region of about 140 kb (44.96–45.1 Mb) (Fig. 2a and Extended Data Fig. 2a), which was still too large to identify a candidate gene. We therefore designed a targeted mapping strategy to resolve the precise location of the causal locus. We employed CRISPR/Cas9-mediated homology-directed repair

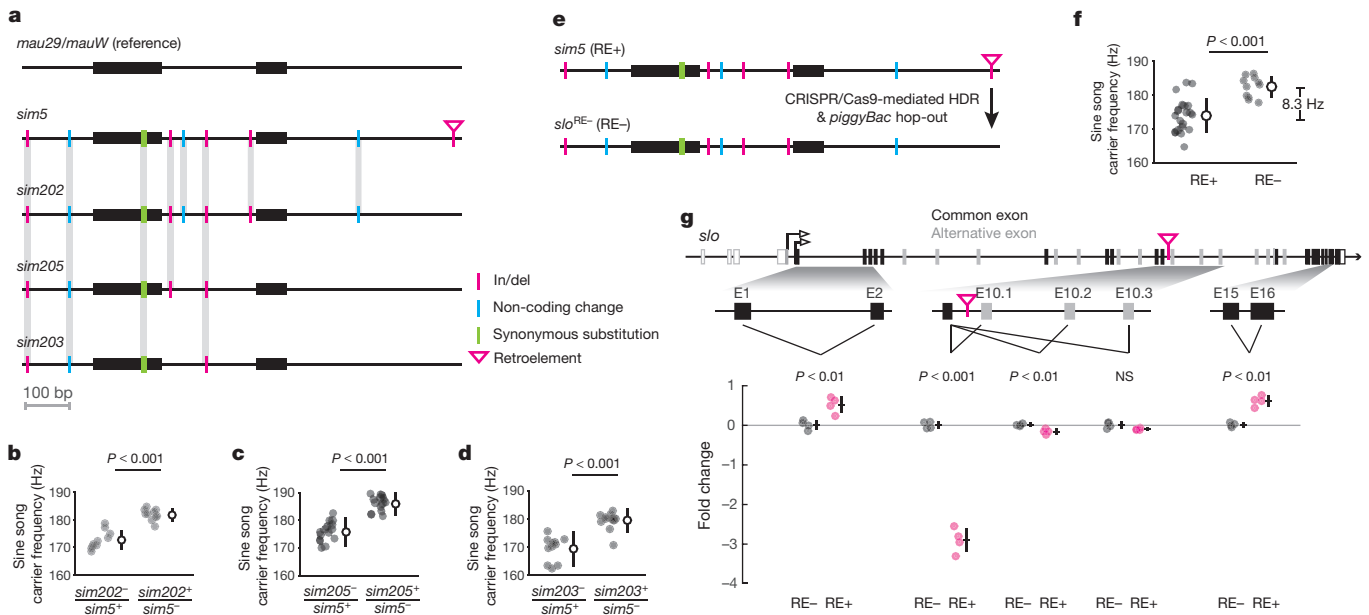


**Figure 3 | Evolution of *slo* causes sine song frequency variation.**

**a**, Schematic of reciprocal hemizyosity test between *D. simulans* (*sim*) and *D. mauritiana* (*mau*) *slo* alleles in the genetic background of the introgression line A2.3. Blue and orange bars indicate *sim5* and *mauW* DNA, respectively.  $\Delta slo$ , *slo*-null allele. **b**, Design of *slo*-null allele *slo* <sup>$\Delta 1$</sup>  induced by CRISPR/Cas9-mediated HDR using single-stranded DNA (ssDNA) as homology donor. PAM, protospacer adjacent motif; nt, nucleotides. **c–f**, Reciprocal hemizyosity test of behavioural phenotypes (mean  $\pm$  s.d.) between *sim5* and *mauW*: sine song frequency (**c**); pulse song frequency (**d**); inter-pulse interval (**e**); and wing beat frequency (**f**). +, *slo* wild-type allele; –, *slo*-null allele *slo* <sup>$\Delta 1$</sup> .  $P$  values by one-way ANOVA; NS, non-significant.

(HDR)<sup>18</sup> to insert 3XP3::DsRed markers at specific sites flanking the 140 kb interval (Extended Data Fig. 3) and identified recombinants in the interval by screening for eye colour (Fig. 2b). We generated 1,152 recombinants within a 170 kb region, providing mapping resolution of one recombination event per 148 bp, on average. All recombinants were genotyped and a subset were phenotyped when the mapped breakpoints could potentially further reduce the causal interval (Extended Data Fig. 2b, c). This effort identified a causal region of 966 bp within the *slo* locus (Fig. 2b), which encodes a calcium-activated potassium channel required to shape the excitability and firing pattern of neurons<sup>1,19</sup>.

Although *slo* is a good candidate gene for courtship song variation, the causal mutation(s) in the mapped interval could, in principle, affect neighbouring genes and not *slo* itself. To directly test whether the evolved change(s) acts on the *slo* locus to alter sine frequency, we performed a reciprocal hemizyosity test, which is considered genetic proof for identifying causal genes underlying quantitative variation<sup>2</sup>. The test is implemented by generating null alleles in each of two strains and by crossing the mutant strains to the reciprocal non-mutated strains (Fig. 3a). The test, therefore, reveals the effects of alternative wild-type alleles in the same genetic background. We used CRISPR/Cas9-mediated HDR<sup>18</sup> to introduce a stop codon in the first coding exon shared by all *slo* splice isoforms in both the *sim5* and *mauW* alleles of the introgression line A2.3 (Fig. 3b and Extended Data Fig. 4). Hemizygous males carrying the *mauW slo* allele sang sine song at 10.0 Hz higher frequency than hemizygous males carrying the *sim5 slo* allele (Fig. 3c). Therefore, variation acting on *slo* causes a sine frequency difference between the two strains.



**Figure 4 | Intronic insertion of a retroelement at the *slo* locus is the causal mutation.** **a**, Candidate mutations in *D. simulans* strains *sim5*, *sim202*, *sim205*, and *sim203*, using *mau29* and *mauW* sequences as reference. Shared sites are highlighted by grey bars. **b–d**, Reciprocal hemizygosity test of sine frequency (mean  $\pm$  s.d.) between *sim5* and *sim202* (**b**), between *sim5* and *sim205* (**c**), and between *sim5* and *sim203* (**d**). +, *slo* wild-type allele; –, *slo*-null allele *slo* <sup>$\Delta$ 2</sup>. *P* values by one-way ANOVA. **e**, Targeted deletion of the retroelement insertion in *sim5* to generate *slo*<sup>RE-</sup>. Experimental details provided in Extended Data Fig. 8.

The *slo* gene is expressed broadly in the fly nervous system<sup>3</sup> (Extended Data Fig. 5a, c, d) and influences many locomotor behaviours<sup>4</sup>. We found that *slo*-null males of *sim5* produced little song with disrupted pulse and sine events (Extended Data Fig. 5b, e–h). In contrast, in the reciprocal hemizygosity test, the evolved allele of *slo* altered only sine frequency, but not pulse song frequency (Fig. 3d), inter-pulse interval (Fig. 3e), wing beat frequency in tethered flight (Fig. 3f), or any other song traits we measured (Extended Data Fig. 6). Therefore, while *slo* has pleiotropic roles, the natural variant of *slo* alters a specific component of song.

There are ten candidate differences in the minimal mapping interval (Fig. 4a): two synonymous coding changes, three non-coding single nucleotide changes, four small non-coding insertions/deletions (in/del), and one 6.7 kb intronic retroelement. To identify the causal mutation(s), we exploited natural variation. We analysed songs from 12 *D. simulans* and 12 *D. mauritiana* wild-type isolates. On average, the sine song frequency of *D. simulans* is 7.8 Hz lower than *D. mauritiana*, although each species contains extensive variation for sine frequency (Extended Data Fig. 7). We sequenced the 966 bp minimal interval in all these strains and found that many of the candidate differences are polymorphic within species, potentially allowing reciprocal hemizygosity testing to narrow down the causal mutation(s).

We generated new *slo*-null alleles in three *D. simulans* strains, *sim202*, *sim205*, and *sim203* (Extended Data Fig. 5b), which harbour nine, six, and five of the ten candidate differences, respectively (Fig. 4a). We performed reciprocal hemizygosity tests between *sim5* and each strain. If the causal mutation is shared by two *D. simulans* strains, then we expect no sine frequency difference between the reciprocal hemizygotes; if the causal mutation is specific to *sim5*, then we expect to observe a difference. For each comparison, the hemizygote with a single *sim5* copy produced song with a significantly lower sine frequency (Fig. 4b–d). The 6.7 kb retroelement insertion is the only polymorphism that differs between *sim5* and all three other strains, making it the best candidate.

**f**, Comparison of sine song frequency (mean  $\pm$  s.d.) between *sim5* (RE+) and *slo*<sup>RE-</sup> (RE-). *P* value by one-way ANOVA. **g**, Expression differences of five *slo* exon junctions between *sim5* (RE+, magenta) and *slo*<sup>RE-</sup> (RE-, grey), assayed by quantitative reverse-transcription PCR (RT-qPCR). The alternative exons E10.1, E10.2 and E10.3 are mutually exclusive in full-length transcripts. Each exon junction was assayed in four biological replicates with two technical replicates. Mean  $\pm$  s.d. *P* values by two-sided *t*-test; NS, non-significant.

We then directly tested the effect of the retroelement by targeted deletion. We first replaced the retroelement with a 3XP3::DsRed marker via CRISPR/Cas9-mediated HDR and then removed the marker (Fig. 4e and Extended Data Fig. 8). The resultant flies, *slo*<sup>RE-</sup>, sang sine song at 8.3 Hz higher frequency than wild-type *sim5* (Fig. 4f). Therefore, the intronic retroelement in the *slo* locus is the causal mutation.

Among the 12 surveyed *D. simulans* strains, this intronic retroelement insertion was detected only in *sim5*. Therefore, it represents a newly derived, rare variant within *D. simulans*. The retroelement resembles a retrovirus (Extended Data Fig. 9a), although it is distinct from any previously characterized retroelement clades (Extended Data Fig. 10). We therefore named it *Shellder*. We identified many polymorphic putative *Shellder* insertions in wild-type strains of *D. simulans* and *D. mauritiana* using TagMap<sup>20</sup> (Extended Data Fig. 9b), suggesting that *Shellder* is probably propagating actively in *Drosophila* populations.

Like many ion channel genes, *slo* exhibits complex patterns of tissue-specific transcription<sup>21</sup> and alternative splicing<sup>22</sup>, potentially producing channels with distinctive properties that could be exploited during evolution. Using the *Shellder* excision allele, we found that the *Shellder* insertion causes a 2.9-fold decrease specifically in the usage of its flanking exon junction but also a slight increase in the overall expression level of *slo* (Fig. 4g). Thus, the insertion does not appear to cause mRNA decay. Considering that lower *slo* level is associated with lower sine frequency (Extended Data Fig. 5i), the lower frequency phenotype caused by the *Shellder* insertion most likely results from the splicing changes of *slo*.

Previous studies have identified genetic variation influencing behaviour<sup>23–28</sup>, including how variation in sensory systems alters the probability of particular behaviours<sup>23–26</sup>. Our study is the first, to our knowledge, to identify the genetic cause for variation in a motor pattern. Our study illustrates that specific behaviour changes can result from certain kinds of regulatory changes in a highly pleiotropic ion channel, in this case most likely through modifications of alternative splicing.

**Online Content** Methods, along with any additional Extended Data display items and Source Data, are available in the online version of the paper; references unique to these sections appear only in the online paper.

**Received 21 October 2015; accepted 8 July 2016.**

**Published online 10 August; corrected online 17 August 2016**

**(see full-text HTML version for details).**

- Atkinson, N. S., Robertson, G. A. & Ganetzky, B. A component of calcium-activated potassium channels encoded by the *Drosophila slo* locus. *Science* **253**, 551–555 (1991).
- Stern, D. L. Identification of loci that cause phenotypic variation in diverse species with the reciprocal hemizyosity test. *Trends Genet.* **30**, 547–554 (2014).
- Becker, M. N., Brenner, R. & Atkinson, N. S. Tissue-specific expression of a *Drosophila* calcium-activated potassium channel. *J. Neurosci.* **15**, 6250–6259 (1995).
- Atkinson, N. S. *et al.* Molecular separation of two behavioral phenotypes by a mutation affecting the promoters of a Ca-activated K channel. *J. Neurosci.* **20**, 2988–2993 (2000).
- Bennet-Clark, H. C. & Ewing, A. W. The courtship songs of *Drosophila*. *Behaviour* **31**, 288–301 (1968).
- Arthur, B. J., Sunayama-Morita, T., Coen, P., Murthy, M. & Stern, D. L. Multi-channel acoustic recording and automated analysis of *Drosophila* courtship songs. *BMC Biol.* **11**, 11 (2013).
- Greenspan, R. J. & Ferveur, J. F. Courtship in *Drosophila*. *Annu. Rev. Genet.* **34**, 205–232 (2000).
- Huttunen, S., Aspi, J., Hoikkala, A. & Schlötterer, C. QTL analysis of variation in male courtship song characters in *Drosophila virilis*. *Heredity* **92**, 263–269 (2004).
- Gleason, J. M., Nuzhdin, S. V. & Ritchie, M. G. Quantitative trait loci affecting a courtship signal in *Drosophila melanogaster*. *Heredity* **89**, 1–6 (2002).
- Turner, T. L., Miller, P. M. & Cochrane, V. A. Combining genome-wide methods to investigate the genetic complexity of courtship song variation in *Drosophila melanogaster*. *Mol. Biol. Evol.* **30**, 2113–2120 (2013).
- Williams, M. A., Blouin, A. G. & Noor, M. A. F. Courtship songs of *Drosophila pseudoobscura* and *D. persimilis*. II. Genetics of species differences. *Heredity* **86**, 68–77 (2001).
- Stern, D. L. Reported *Drosophila* courtship song rhythms are artifacts of data analysis. *BMC Biol.* **12**, 38 (2014).
- Cande, J., Stern, D. L., Morita, T., Prud'homme, B. & Gompel, N. Looking under the lamp post: neither *fruitless* nor *doublesex* has evolved to generate divergent male courtship in *Drosophila*. *Cell Reports* **8**, 363–370 (2014).
- Garrigan, D. *et al.* Genome sequencing reveals complex speciation in the *Drosophila simulans* clade. *Genome Res.* **22**, 1499–1511 (2012).
- Andolfatto, P. *et al.* Multiplexed shotgun genotyping for rapid and efficient genetic mapping. *Genome Res.* **21**, 610–617 (2011).
- Greenwood, A. K., Wark, A. R., Yoshida, K. & Peichel, C. L. Genetic and neural modularity underlie the evolution of schooling behavior in threespine sticklebacks. *Curr. Biol.* **23**, 1884–1888 (2013).
- Weber, J. N., Peterson, B. K. & Hoekstra, H. E. Discrete genetic modules are responsible for complex burrow evolution in *Peromyscus* mice. *Nature* **493**, 402–405 (2013).
- Bassett, A. R. & Liu, J.-L. CRISPR/Cas9 and genome editing in *Drosophila*. *J. Genet. Genomics* **41**, 7–19 (2014).
- Vergara, C., Latorre, R., Marrion, N. V. & Adelman, J. P. Calcium-activated potassium channels. *Curr. Opin. Neurobiol.* **8**, 321–329 (1998).
- Stern, D. L. Tagmentation-based mapping (TagMap) of mobile DNA genomic insertion sites. Preprint at <http://dx.doi.org/10.1101/037762> (2016).
- Brenner, R., Thomas, T. O., Becker, M. N. & Atkinson, N. S. Tissue-specific expression of a Ca<sup>2+</sup>-activated K<sup>+</sup> channel is controlled by multiple upstream regulatory elements. *J. Neurosci.* **16**, 1827–1835 (1996).
- Yu, J. Y., Upadhyaya, A. B. & Atkinson, N. S. Tissue-specific alternative splicing of BK channel transcripts in *Drosophila*. *Genes Brain Behav.* **5**, 329–339 (2006).
- Bendesky, A., Tsunozaki, M., Rockman, M. V., Kruglyak, L. & Bargmann, C. I. Catecholamine receptor polymorphisms affect decision-making in *C. elegans*. *Nature* **472**, 313–318 (2011).
- Bendesky, A. *et al.* Long-range regulatory polymorphisms affecting a GABA receptor constitute a quantitative trait locus (QTL) for social behavior in *Caenorhabditis elegans*. *PLoS Genet.* **8**, e1003157 (2012).
- de Bono, M. & Bargmann, C. I. Natural variation in a neuropeptide Y receptor homolog modifies social behavior and food response in *C. elegans*. *Cell* **94**, 679–689 (1998).
- McGrath, P. T. *et al.* Quantitative mapping of a digenic behavioral trait implicates globin variation in *C. elegans* sensory behaviors. *Neuron* **61**, 692–699 (2009).
- Osborne, K. A. *et al.* Natural behavior polymorphism due to a cGMP-dependent protein kinase of *Drosophila*. *Science* **277**, 834–836 (1997).
- Yalcin, B. *et al.* Genetic dissection of a behavioral quantitative trait locus shows that *Rgs2* modulates anxiety in mice. *Nat. Genet.* **36**, 1197–1202 (2004).

**Supplementary Information** is available in the online version of the paper.

**Acknowledgements** We thank C. Llorens for advice on retroelement classification and B. Dickson, N. Frankel, T. Shirangi, N. Spruston and S. Sternson for advice on the manuscript.

**Author Contributions** Y.D. and D.L.S. designed the experiments. A.B. performed QTL analysis, with help from T.M. Y.D. and K.D.L. assayed wing beat frequency. D.L.S. performed TagMap analysis. Y.D. performed all other experiments and analysis. Y.D. and D.L.S. wrote the manuscript.

**Author Information** Reprints and permissions information is available at [www.nature.com/reprints](http://www.nature.com/reprints). The authors declare no competing financial interests. Readers are welcome to comment on the online version of the paper. Correspondence and requests for materials should be addressed to D.L.S. ([sternd@janelia.hhmi.org](mailto:sternd@janelia.hhmi.org)).

**Reviewer Information** *Nature* thanks B. Prud'homme, T. Turner and the other anonymous reviewer(s) for their contribution to the peer review of this work.

## METHODS

The experiments were not randomized. The investigators were not blinded to allocation during experiments and outcome assessment, but all statistical analysis was performed automatically by our analysis pipeline.

**Strains and behavioural assays.** Fly strains used are summarized in Supplementary Table 1. For all behavioural assays, the flies were reared in standard laboratory conditions. Courtship song was recorded as described previously<sup>6</sup>. For each comparison reported in this study, the flies were recorded simultaneously, and if applicable, collected from the same vials. Song parameters were estimated as the mode of song events across 30-min recordings. Wing beat frequency was measured on tethered male flies given a bar fixation task during flight, by optically tracking wing movements with Wingbeat Analyzer (JFI Electronics Laboratory, University of Chicago)<sup>29</sup>.

**Behaviour data analysis and statistics.** Song data was segmented<sup>6</sup> and analysed (<http://www.github.com/dstern/BatchSongAnalysis>) without human intervention. For each test, our target sample size was  $n = 12$  per genotype, with individuals selected haphazardly, for power of 0.8 to detect 1 s.d. difference between treatments at  $P < 0.05$ . Outliers were systematically excluded in our song analysis pipeline using the Grubbs test with  $\alpha = 0.05$  (<http://www.mathworks.com/matlabcentral/fileexchange/3961-deleteoutliers>).  $P$  values for ANOVAs were estimated with 10,000 permutations (<http://www.mathworks.com/matlabcentral/fileexchange/44307-randanova1>). All critical experiments were replicated at least once with a similar sample size.

**QTL mapping.** QTL mapping employed 210 *sim5* and 180 *mau29* backcross progeny, which were processed into a single multiplexed shotgun genotyping (MSG) library<sup>15</sup>. Parental genomes were generated by updating the *D. simulans* r2.0.1 genome ([http://www.flybase.org/static\\_pages/feature/previous/articles/2015\\_02/Dsim\\_r2.01.html](http://www.flybase.org/static_pages/feature/previous/articles/2015_02/Dsim_r2.01.html)) with HiSeq reads from each strain (SRA accession: SRP076910). Genotypes were estimated with MSG software (<http://www.github.com/janeliaSciComp/msg>). Posterior probabilities of ancestry were thinned using *pull\_thin* ([http://www.github.com/dstern/pull\\_thin](http://www.github.com/dstern/pull_thin)) and imported into R-QTL<sup>30</sup> using *read\_cross\_msg* ([http://www.github.com/dstern/read\\_cross\\_msg](http://www.github.com/dstern/read_cross_msg)). Genome scans with a single QTL model were performed using Haley–Knott regression<sup>30</sup> and  $P$  values were estimated with 1,000 permutations.

**Fine-scale introgression mapping and high-resolution recombination mapping.** Genetic mapping was performed in three phases. First, to develop a visible marker linked to the QTL, we screened a collection of *mauW* strains that had been transformed with a *piggyBac* transposable element carrying 3XP3::EYFP, which

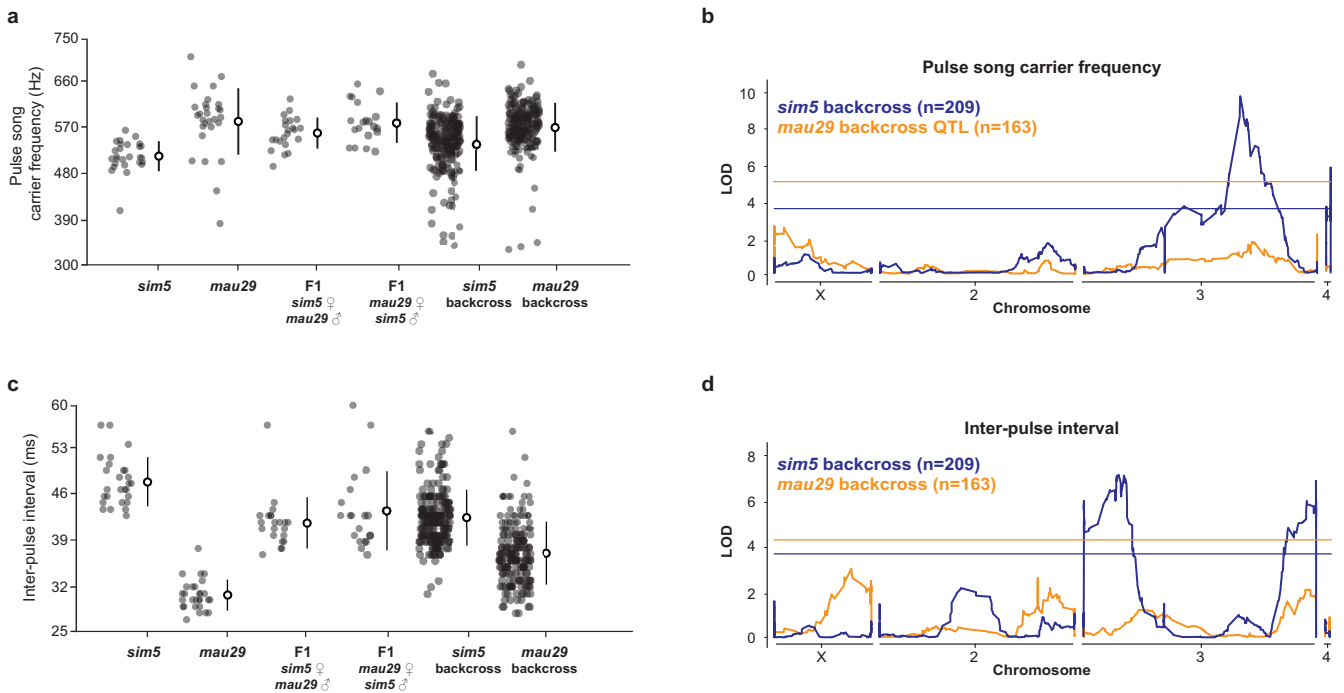
drives yellow fluorescent protein expression in the eyes. The strain *mauW446*, which carried a *piggyBac* insertion within 3 Mb of the QTL peak, was backcrossed to *sim5* for five generations. The genetic background of the resulting introgression lines was checked using whole genome sequencing. Second, the introgression line A2, which produces *D. mauritiana*-like sine song and does not contain any *D. mauritiana* DNA outside of the target region, was backcrossed to *sim5* for three to five generations for further introgression. Third, we inserted two 3XP3::DsRed markers into the introgression line A2.3: one on the left side of the interval on the EYFP-marked introgression chromosome (44.95 Mb), and the other on the right side on the *sim5* chromosome (45.12 Mb). Female flies carrying both DsRed markers were crossed to *sim5* males and male progeny without a DsRed marker were identified as recombinants. All these lines were maintained through the male lineage, which does not experience recombination. Recombination breakpoints were mapped using molecular markers (Supplementary Table 2) and sequencing.

**CRISPR/Cas9 genome editing.** For all CRISPR/Cas9-mediated HDR, guide RNAs (gRNAs), donor DNA, *in vitro* transcribed *D. melanogaster* codon optimized Cas9 mRNA, and a Dicer-substrate short interfering RNA (DsiRNA) targeted against *lig4* (Sequence 1: rUrCrCrUrGrCrArGrCrUrGrArUrGrCrUrUrGrCrUrGr rUrGrUrCrGrU; Sequence 2: rGrArCrArCrArGrCrArArGrCrArUrCrArGrCr rUrGrCrArGG A, synthesized by IDT) to inhibit non-homologous end joining<sup>18</sup> were co-injected into the embryos. All the injections were performed by Rainbow Transgenic Flies using the following concentrations: 0.2  $\mu\text{g}/\mu\text{l}$  gRNA source, 0.5  $\mu\text{g}/\mu\text{l}$  donor DNA, 0.1  $\mu\text{g}/\mu\text{l}$  Cas9 mRNA, and 0.1  $\mu\text{g}/\mu\text{l}$  *lig4* DsiRNA. The germline transmission rates are provided in Supplementary Table 3.

**RT-qPCR.** Total RNA was extracted from 5–7 day old males and converted to cDNA template after DNase I treatment. Real-time PCR was performed using *ACTB* as an internal control. The primer sequences are provided in Supplementary Table 4.

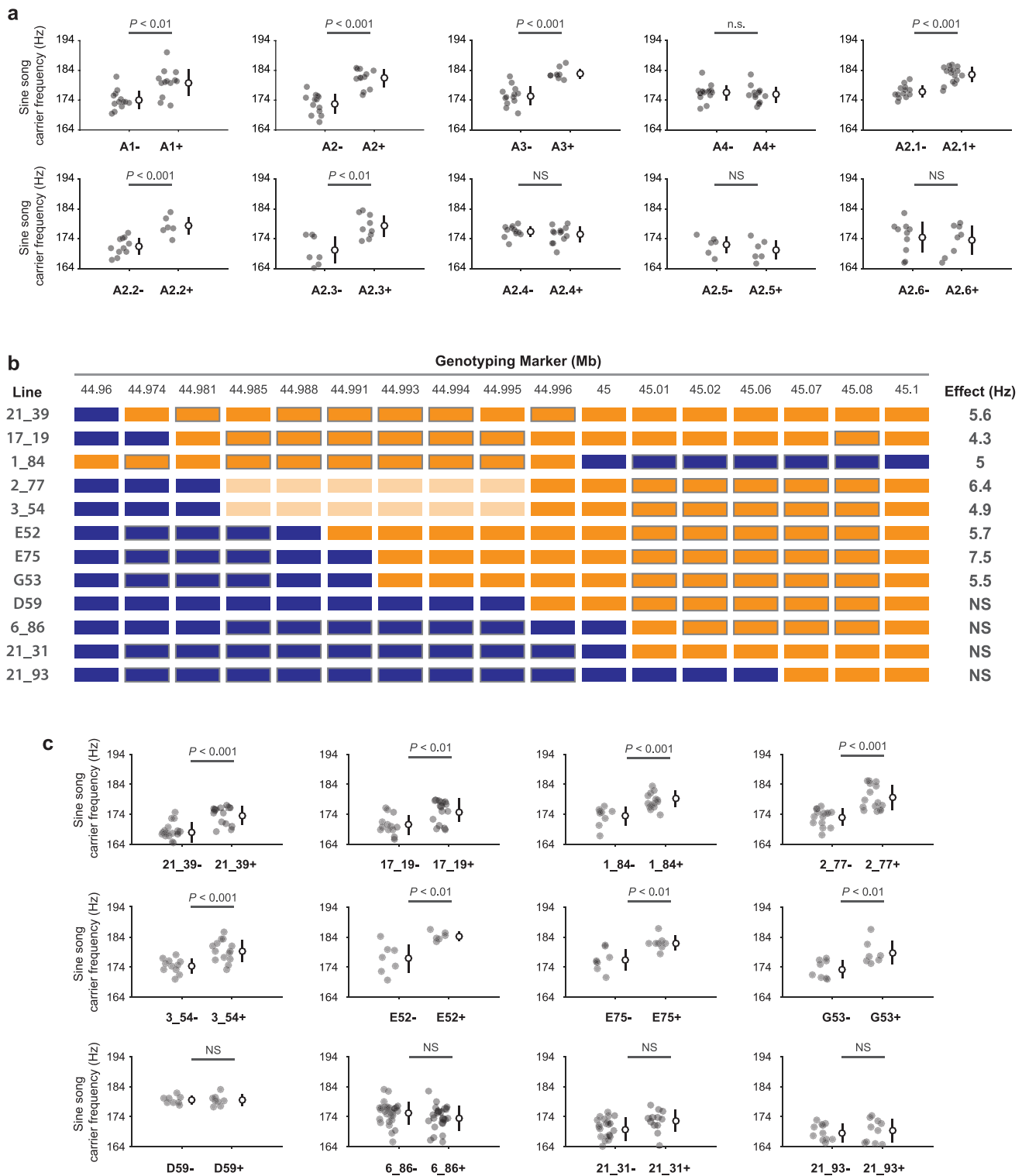
**Shellder analysis.** The name comes from the Pokémon character Shellder, who can attach to the character Slowpoke and cause Slowpoke to evolve into a new form. The sequence annotation, phylogenetic analysis, and TaqMap identification of *Shellder* (GenBank KX196449) insertion sites are described in Supplementary Methods.

29. Gotz, K. G. Course-control, metabolism and wing interference during ultralong tethered flight in *Drosophila melanogaster*. *J. Exp. Biol.* **128**, 35–46 (1987).
30. Broman, K. W., Wu, H., Sen, S. & Churchill, G. A. R/qtl: QTL mapping in experimental crosses. *Bioinformatics* **19**, 889–890 (2003).



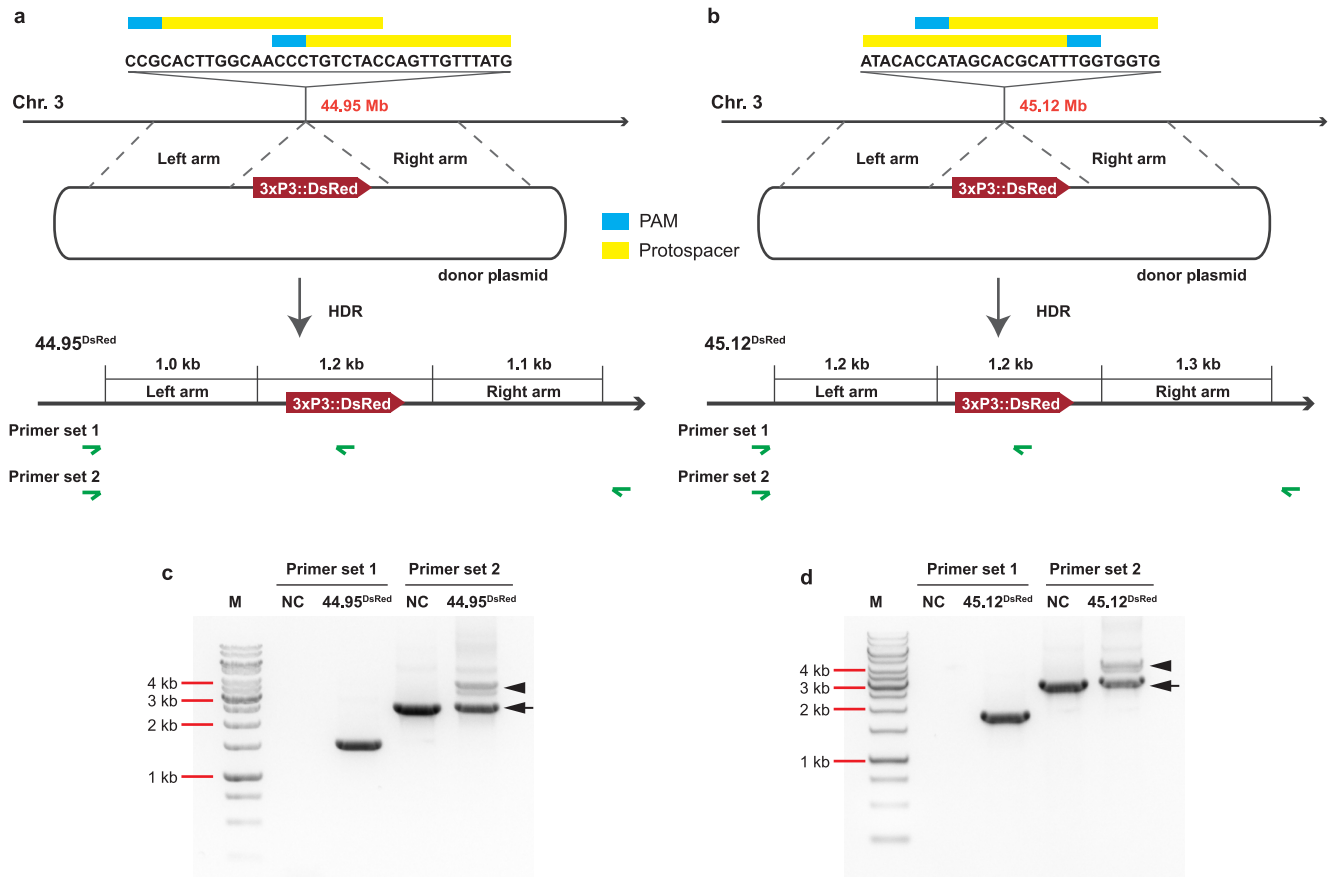
**Extended Data Figure 1 | QTL analysis of pulse song carrier frequency and inter-pulse interval.** **a**, Pulse song carrier frequency (mean  $\pm$  s.d.) in parental strains, F1 hybrids, and backcross males. **b**, QTL map of pulse song frequency. LOD, logarithm of the odds. Horizontal lines mark  $P = 0.01$ . A single QTL on chromosome 3 at 34,919,124 bp was identified in the *sim5* backcross (blue). No significant QTL were detected in the

*mau29* backcross (orange). **c**, Inter-pulse interval (mean  $\pm$  s.d.) in parental strains, F1 hybrids, and backcross males. **d**, QTL map of inter-pulse interval. Two QTLs on chromosome 3 at 7,902,342 bp and 52,144,317 bp were identified in the *sim5* backcross (blue). No significant QTL were detected in the *mau29* backcross (orange).



**Extended Data Figure 2 | Details of fine-scale mapping.** **a**, Sine song frequency of the introgression lines shown in Fig. 2a. **b**, Genotype and phenotype data of informative recombinant lines. The genotyping markers are listed on the top, according to their physical locations on chromosome 3. For each recombinant line, the genotyping results are represented using coloured bars: blue, homozygote for *sim5*; orange, heterozygote for *sim5* and *mauW*; and light orange, unknown. For boxed bars, the genotypes were assigned assuming no additional recombination events in this region.

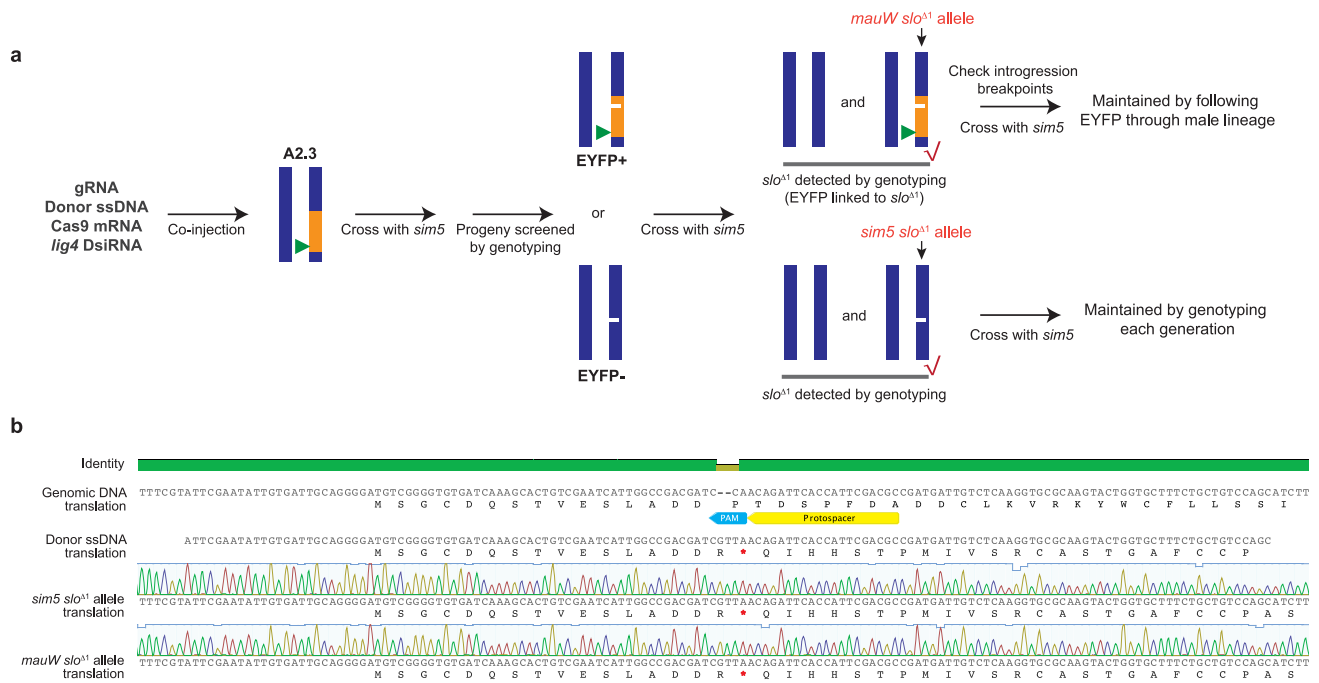
The sine song frequency effect is summarized on the right. Line 1\_84 was recovered by selecting flies with stronger DsRed and without EYFP. (This strategy was abandoned later due to the challenge of distinguishing two copies of DsRed from one. Please refer to Fig. 2b.) **c**, Sine song frequency phenotypes of the recombinant lines in panel **b**. Comparison of sine song frequency (mean  $\pm$  s.d.) was made between heterozygous introgression males (+) and their pure *sim5* sibling brothers (-) using one-way ANOVA. NS, non-significant.



**Extended Data Figure 3 | Targeted insertion of DsRed markers via CRISPR/Cas9-mediated HDR.** **a, b,** Schematics of the targeted insertions of the markers 44.95<sup>DsRed</sup> (**a**) and 45.12<sup>DsRed</sup> (**b**). PAM, protospacer adjacent motif. **c, d,** PCR validation of 44.95<sup>DsRed</sup> (**c**) and 45.12<sup>DsRed</sup> (**d**) using the primer sets indicated in the panels above. A heterozygous

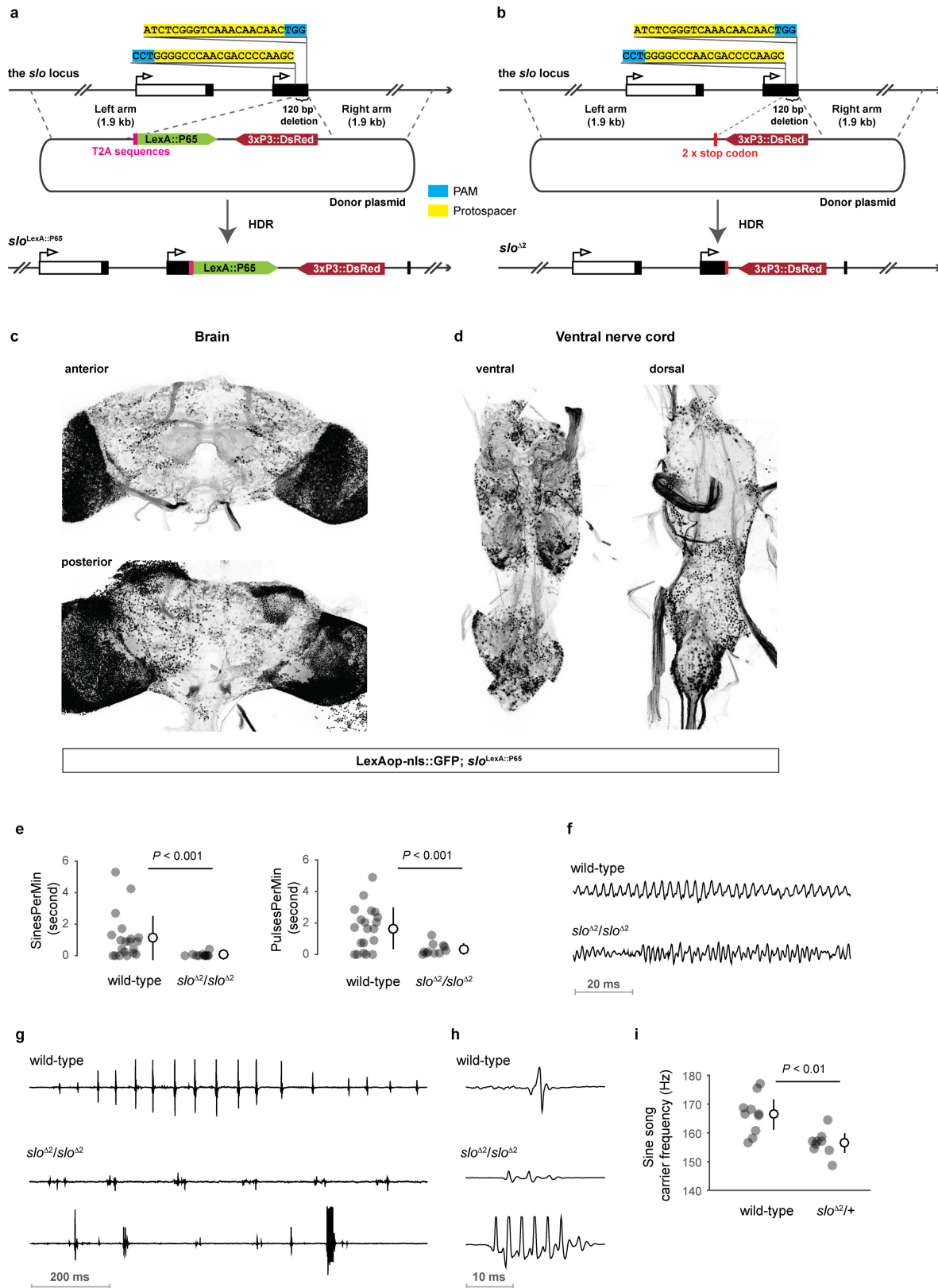
fly was used for PCR and the negative control (NC) used a fly of the introgression line A2.3. The lower band (arrow) seen with primer set 2 therefore represents the wild-type allele. The PCR fragment containing the DsRed insertion is indicated by the white arrowhead. M, GeneRuler 1 kb DNA ladder.





**Extended Data Figure 4 | Generation and validation of the *slo*-null allele *slo<sup>Δ1</sup>*.** **a**, Schematic of the generation of the *sim5* and *mauW slo<sup>Δ1</sup>* alleles in the genetic background of the introgression line A2.3. Blue and orange bars represent *sim5* and *mauW* DNA, respectively. Green triangle denotes the EYFP marker. gRNA, guide RNA; ssDNA, single-stranded

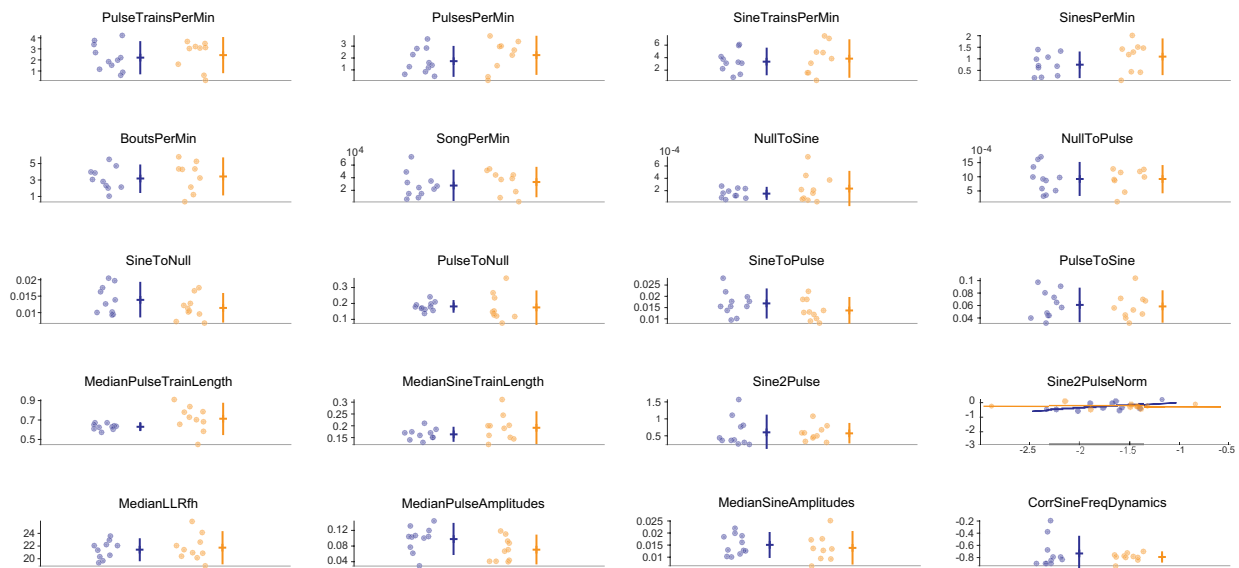
DNA; DsiRNA, Dicer-substrate short interfering RNA. **b**, Sequence verification of *slo<sup>Δ1</sup>*. Successful integration of ssDNA was confirmed by sequencing the cloned PCR products amplified from a heterozygous fly. Red asterisks indicate the introduced stop codon. PAM, protospacer adjacent motif.



**Extended Data Figure 5 | Expression pattern and song phenotype of *slo*.**

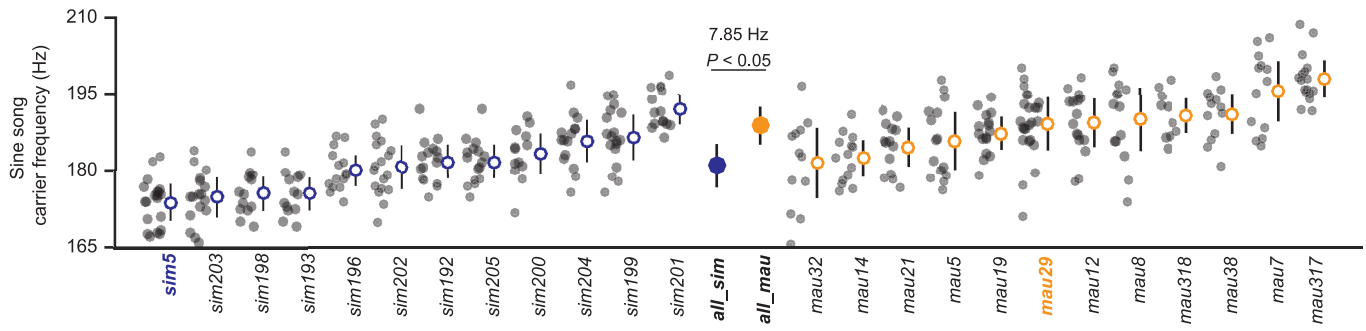
**a**, Schematic of the generation of the *slo*<sup>LexA::P65</sup> allele. T2A self-cleavage peptide sequences, LexA::P65, a 120 bp deletion, and a 3XP3::DsRed insertion were introduced into the first common exon shared by all *slo* transcripts in *D. melanogaster*. **b**, Schematic of the generation of the *slo*<sup>Δ2</sup> allele. Two stop codons, a 120 bp deletion, and a 3XP3::DsRed insertion were introduced into the same locus indicated in panel **a**. The same allele was generated in each of the following *D. simulans* strains: *sim5*, *sim202*, *sim203*, *sim205*. PAM, protospacer adjacent motif.

**c, d**, Brain (**c**) and ventral nerve cord (**d**) of *D. melanogaster* LexAop-nls::GFP, *slo*<sup>LexA::P65</sup> males showed widespread expression of *slo*. **e–h**, Courtship song phenotype of *sim5 slo*<sup>Δ2</sup> males. The *slo*-null males produced very little song (**e**). SinesPerMin and PulsesPerMin measure the average amount of sine song and pulse song produced per minute, respectively. The sine and pulse events of the *slo*-null males were severely disrupted: sine song waveform (**f**); pulse train (**g**); pulse song waveform (**h**). **i**, *slo*<sup>Δ2</sup> heterozygotes sang sine song at 9.9 Hz lower frequency than wild-type. Data represent mean ± s.d. *P* values by one-way ANOVA.



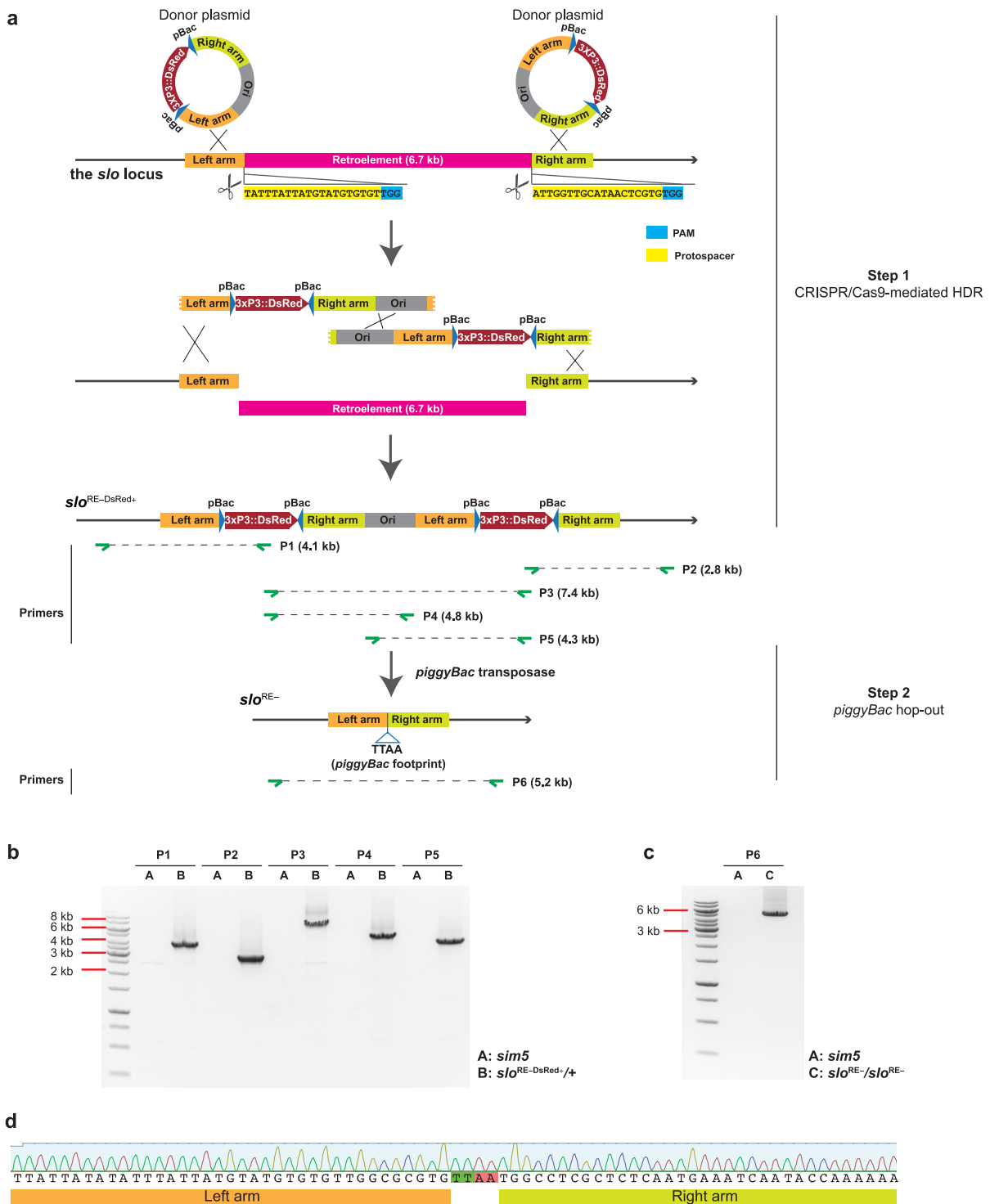
**Extended Data Figure 6 | Additional song phenotypes of *slo* hemizygotes.** Blue indicates the genotype *mauW<sup>-</sup>/sim5<sup>+</sup>* and orange indicates *mauW<sup>+</sup>/sim5<sup>-</sup>*. None of these song phenotypes are significantly different by one-way ANOVA ( $P > 0.05$ ). Data represent mean  $\pm$  s.d. The definitions of each song phenotype are as follows: PulseTrainsPerMin, average number of pulse trains per minute; PulsesPerMin, average pulse song duration in seconds per minute; SineTrainsPerMin, average number of sine trains per minute; SinesPerMin, average sine song duration in seconds per minute; BoutsPerMin, average number of song bouts per minute; SongPerMin, average song duration per minute; NullToSine, transition probability from no song to sine song; NullToPulse, transition probability from no song to pulse song; SineToNull, transition probability from sine song to no song; PulseToNull,

transition probability from pulse song to no song; SineToPulse, transition probability from sine song to pulse song within song bouts; PulseToSine, transition probability from pulse song to sine song within song bouts; MedianPulseTrainLength, median length of pulse trains; MedianSineTrainLength, median length of sine trains; Sine2Pulse, ratio of amount of sine song to amount of pulse song; Sine2PulseNorm, Sine2Pulse normalized by the amount of song; MedianLLRfh, median score of the log likelihood ratio of fit of pulse to pulse model versus white noise (measure of pulse shape); MedianPulseAmplitudes, median amplitude of pulses; MedianSineAmplitudes, median amplitude of sines; CorrSineFreqDynamics, slope of sine song carrier frequency within song bouts.



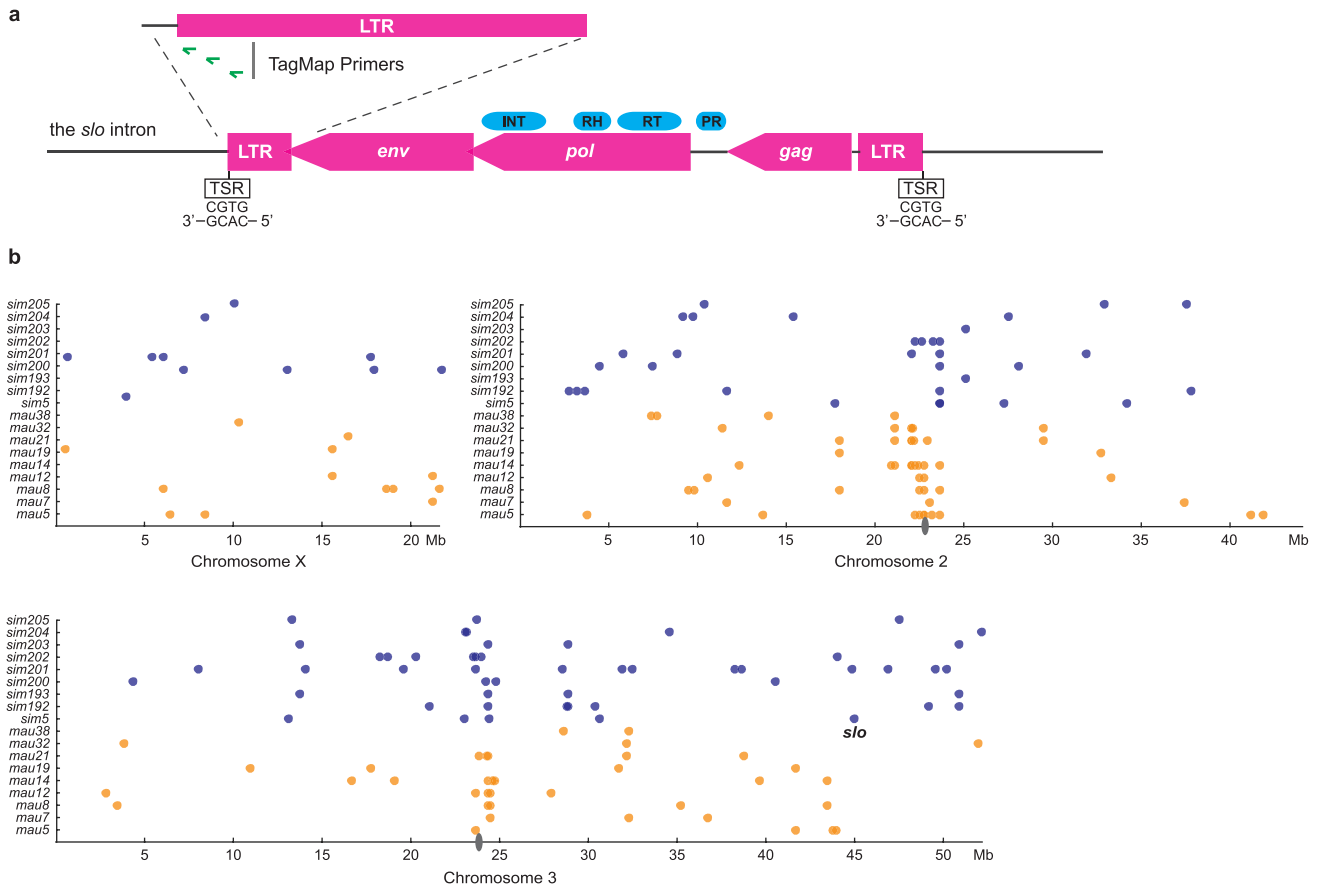
**Extended Data Figure 7 | Sine song frequency phenotypes of 12 *D. simulans* wild-type isolates (blue) and 12 *D. mauritiana* wild-type isolates (orange).** Open circles with lines represent mean  $\pm$  s.d. Closed circles with lines represent mean  $\pm$  s.d. of the means of all *D. simulans*

(*all\_sim*) and *D. mauritiana* strains (*all\_mau*).  $P$  value by one-way ANOVA. All strains were recorded simultaneously through multiple recording sessions.



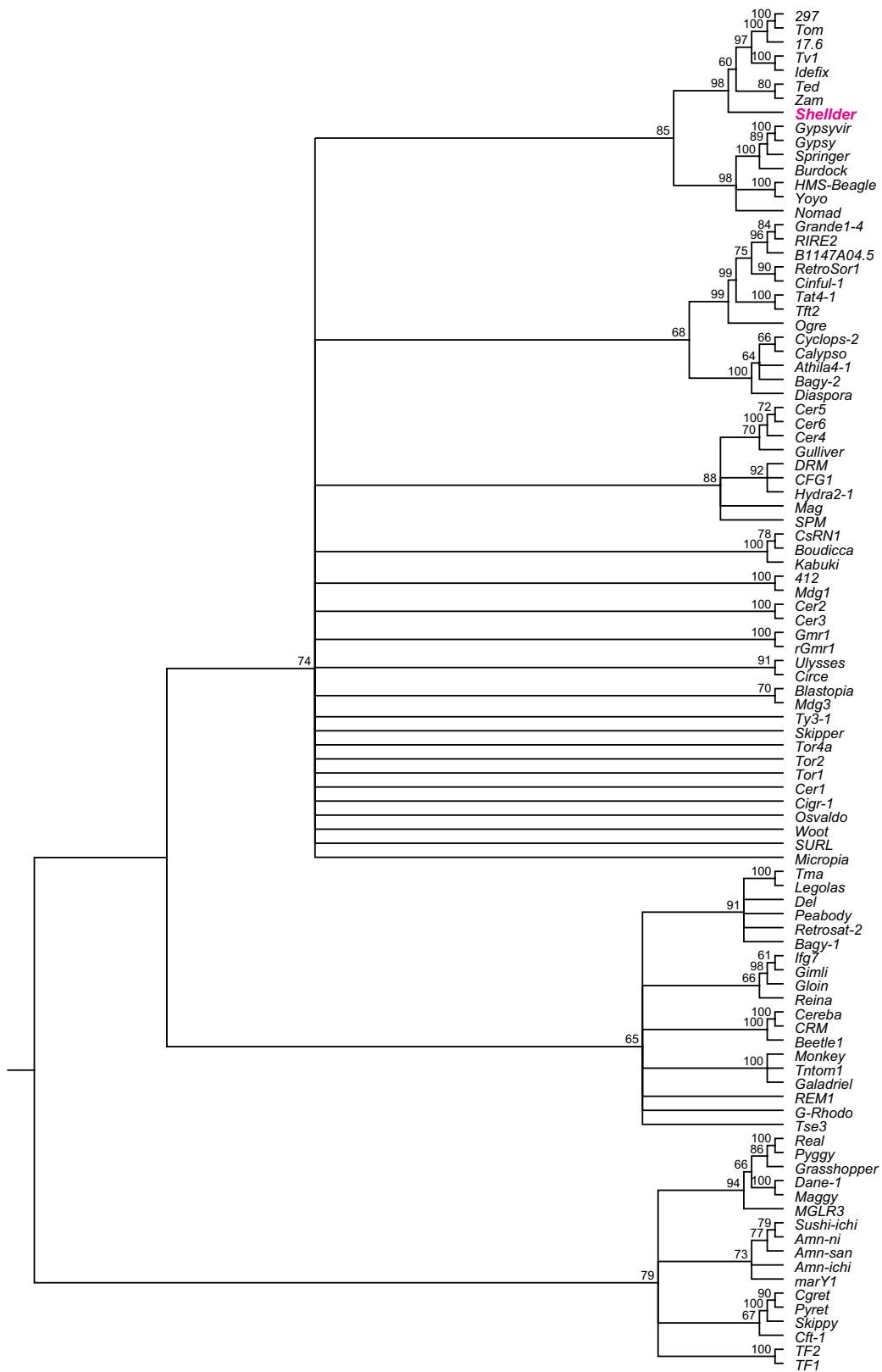
**Extended Data Figure 8 | Targeted deletion of the retroelement insertion at the *slo* intron in *sim5*.** **a**, Putative schematic of the targeted deletion of the retroelement insertion using a two-step strategy. First, the retroelement was replaced by a 3XP3::DsRed marker cassette flanked by *PiggyBac* transposon ends (labelled as pBac), using CRISPR/Cas9-mediated HDR. This generated the  $slo^{RE-DsRed+}$  allele. Second, the 3XP3::DsRed marker was hopped out of the genome using *piggyBac* transposase to generate the  $slo^{RE-}$  allele. In the first step, two independent

HDR events appear to have occurred between the donor plasmids and the *slo* locus. All three  $slo^{RE-DsRed+}$  alleles we generated support this rare recombination type (data not shown), which possibly reflects the difficulty of repairing a large deletion with a single donor plasmid. PAM, protospacer adjacent motif. **b**, **c**, PCR validation of the  $slo^{RE-DsRed+}$  allele (**b**) and the  $slo^{RE-}$  allele (**c**) using the primer sets indicated in panel **a**. M, GeneRuler 1 kb DNA ladder. **d**, Sequence verification of the  $slo^{RE-}$  allele. The *piggyBac* footprint 'TTAA' is highlighted.



**Extended Data Figure 9 | Identification of putative *Shellder* copies in *D. simulans* and *D. mauritiana* populations.** **a**, Schematic of the *Shellder* insertion at the *slo* locus in *sim5*. *Shellder* contains three open reading frames (ORFs) resembling the *gag*, *pol*, and *env* genes of a retrovirus, flanked by 458 bp long terminal repeats (LTRs). Putative core protein domains are indicated: PR, protease; RT, reverse transcriptase; RH, RNase H; INT, integrase. In other retroviruses, the *pol* ORF often includes a 5' protease domain. *Shellder* contains a conserved protease domain 5' of the predicted *pol* start codon. It is possible that the *Shellder* *pol* ORF uses a

non-ATG start codon. TSR, target site repeat. P1, P2, and P3 represent the three LTR-specific primers used for TagMap. **b**, Putative *Shellder* copies in *D. simulans* (blue) and *D. mauritiana* (orange) wild-type strains identified by TagMap. The *slo* locus insertion in *sim5* is indicated and is unique amongst these samples. *Shellder* insertions are enriched near centromeres (grey ovals), but can also be found in the euchromatic regions. Precise mapping locations are provided in Supplementary Table 5. This is probably an incomplete survey of *Shellder* copies, because TagMap may be biased towards detecting young and intact copies of transposable elements.



Extended Data Figure 10 | Phylogenetic position of *Shellder* in the Ty3/Gypsy family based on reverse transcriptase protein sequences. Bootstrap values (%) are indicated on the branches only when they exceed 60. Branch lengths are not drawn proportional to genetic distance.