1 Non-allelic gene conversion enables rapid evolutionary change at multiple

2 regulatory sites encoded by transposable elements

3

4 Christopher E. Ellison & Doris Bachtrog

5 Department of Integrative Biology, University of California Berkeley, Berkeley, CA 94720, USA

6

7

8 Transposable elements (TEs) allow rewiring of regulatory networks, and the recent 9 amplification of the ISX-element dispersed 77 functional but suboptimal binding-sites for 10 the dosage-compensation-complex to a newly-formed X-chromosome in Drosophila. Here 11 we identify two linked refining-mutations within ISX that interact epistatically to increase 12 binding affinity to the dosage-compensation-complex. Selection has increased the frequency 13 of this derived haplotype in the population, which is fixed at 30% of ISX-insertions and 14 polymorphic among another 41%. Sharing of this haplotype indicates that high levels of 15 gene-conversion among ISX-elements allow them to "crowd-source" refining-mutations, 16 and a refining-mutation that occurs at any single ISX-element can spread in two 17 dimensions: horizontally across insertion sites by non-allelic gene-conversion, and 18 vertically through the population by natural selection. These findings describe a novel 19 route by which fully functional regulatory elements can arise rapidly from TEs and 20 implicate non-allelic gene-conversion as having an important role in accelerating the 21 evolutionary fine-tuning of regulatory networks.

22

23 Introduction

24 A substantial portion of animal genomes is composed of repetitive sequences, including gene 25 duplicates, satellite DNA, and transposable elements. Gene conversion is a major force shaping 26 the evolution of repetitive regions, and interlocus or non-allelic gene conversion between 27 sequence duplicates has been studied extensively for its role in concerted evolution (Chen et al., 28 2007; Ohta, 2010). Non-allelic gene conversion also affects selection operating in gene families. 29 Compared to single-copy genes, a family of gene duplicates presents a larger mutational target, 30 and a mutation arising in any gene copy can be spread among copies by non-allelic gene 31 conversion, thereby increasing the efficiency of both positive and purifying selection (Mano and 32 Innan, 2008). Non-allelic gene conversion homogenizes the arrays of ribosomal DNA gene

copies present in the genomes of most organisms (Eickbush and Eickbush, 2007), has generated
allelic diversity within the human leukocyte antigen gene family (Zangenberg et al., 1995), and
has allowed palindromic genes on the human Y chromosome to escape degeneration (Rozen et
al., 2003).

37

38 Transposable elements give rise to families of duplicate sequences. A propensity for some TEs to 39 carry regulatory motifs and to insert adjacent to coding sequence gives them the potential for 40 being potent modulators of gene regulatory networks (Cowley and Oakey, 2013; Feschotte, 41 2008). The regulatory elements provided by these TEs, however, may be suboptimal in function, 42 and subject to subsequent fine-tuning (Polavarapu et al., 2008). Unlike regulatory elements 43 where short binding motifs (ten basepairs on average for transcription factors; Stewart et al., 44 2013) evolve *de novo* via point mutation or microsatellite expansion, binding sites that evolve 45 from TEs are initially almost identical in sequence and are nested within a larger repeat unit 46 (hundreds or thousands of basepairs in size), and may thus be subject to non-allelic gene 47 conversion. Re-wiring of the dosage compensation network in *Drosophila miranda* was driven 48 by TE-mediated amplification of a functional but suboptimal binding motif (Ellison and 49 Bachtrog, 2013). Here we show that non-allelic gene conversion is catalyzing the rapid fine-50 tuning of these suboptimal motifs by allowing sequence variants that optimize binding affinity to 51 spread among elements.

52

53 Dosage compensation in Drosophila is mediated by a male-specific ribonucleoprotein complex 54 (the male-specific lethal or MSL complex) that binds to a GA-rich sequence motif (the MSL 55 recognition motif) at a number of chromatin entry sites on the X chromosome (Alekseyenko et 56 al., 2008; Straub et al., 2008). We previously studied the acquisition of novel chromatin entry 57 sites on newly formed X chromosomes in Drosophila miranda, a species where two independent 58 sex chromosome/autosome fusions resulted in a karyotype composed of three X chromosome 59 arms, each of a different age (Alekseyenko et al., 2013; Zhou et al., 2013). XL is homologous to 60 the X chromosome of D. melanogaster and has been a sex chromosome for at least 60 million 61 years (Richards et al., 2005); chromosome XR formed roughly 15 million years ago when an 62 autosome (Muller element D) fused to XL (Carvalho and Clark, 2005), and the neo-X/neo-Y 63 chromosome pair originated around 1.5 million years ago when the Y fused to another autosome

(Muller element C) (Bachtrog and Charlesworth, 2002). Dosage compensation evolved on both
XR and the neo-X shortly after their emergence, through acquisition of novel chromatin entry
sites and co-option of the MSL regulatory network (Bone and Kuroda, 1996; Marin et al., 1996).
Interestingly, we discovered that the acquisition of dosage compensation on both XR and the
neo-X chromosome was in part mediated by the independent domestication of helitron
transposable elements that contained MSL recognition motifs, which we have termed ISXR and
ISX, respectively (Ellison and Bachtrog, 2013).

72 ISX is highly enriched on the neo-X chromosome of D. miranda and is derived from the 73 abundant ISY element. Compared to ISY, ISX contains a 10 basepair deletion that creates a MSL 74 recognition motif, thereby allowing it to act as a chromatin entry site (Ellison and Bachtrog, 75 2013). Our previous study showed that while amplification of ISX about 1 million years ago 76 provided dozens of functional chromatin entry sites on the neo-X chromosome of D. miranda, 77 the motif dispersed by ISX is distinct from the canonical motif that is enriched within chromatin 78 entry sites on XL and XR, and shows significantly lower affinity to the MSL complex compared 79 to motifs on XL and XR (Ellison and Bachtrog, 2013). For these reasons, we postulated that the 80 ISX binding motif is suboptimal, and predicted that refining mutations should accumulate within 81 each MSL recognition motif until the neo-X chromosome becomes fully dosage compensated 82 (Ellison and Bachtrog, 2013).

83

84 **Results**

85 Variation at MSL recognition motifs among ISX insertions in D. miranda strain MSH22 86 To identify potential refining mutations that optimize MSL-binding at chromatin entry sites 87 derived from the ISX element, we characterized sequence variation within the MSL recognition 88 motifs and flanking sequence regions for all 77 insertions of the ISX element on the neo-X 89 chromosome in the sequenced reference strain MSH22 (Figure 1A). Because we have 90 previously demonstrated that ISX contains a functional MSL recognition motif but the closely 91 related ISY element does not (Ellison and Bachtrog, 2013), we sought to identify sequence 92 variants that were present in multiple ISX elements but rare or absent in ISY elements from the 93 same chromosome.

94

95 Using these criteria, we identified a sequence haplotype adjacent to the MSL recognition motif

- 96 that is common among MSH22 ISX insertions and rare among ISY elements: 57% of ISX
- 97 elements carry this haplotype versus 0.7% of neo-X ISY insertions, an asymmetry significantly
- 98 different from that expected by chance (Fisher's Exact Test; P < 2.2e-16). The haplotype consists
- 99 of two mutations (G \rightarrow T and A \rightarrow T), separated by two basepairs, which are in perfect linkage
- 100 disequilibrium among ISX but not ISY elements (Figure 1 & Figure 1-figure supplement).
- 101 Because ISX is descended from ISY and the TT alleles are rare among ISY elements, they are
- 102 likely to be derived. We hereafter refer to these mutations as the TT haplotype.
- 103

104 The TT haplotype increases MSL complex binding affinity

To determine if the TT haplotype affects binding affinity of the MSL complex, we used published ChIP-seq data of MSL3 (a component of the MSL complex) from *D. miranda* strain MSH22 (Alekseyenko et al., 2013). We compared *in vivo* MSL complex binding levels for the 44 MSH22 ISX insertions carrying the TT haplotype to the 33 insertions with the ancestral GA haplotype. The insertions with the TT alleles had significantly higher levels of MSL complex binding compared to those with the GA alleles (Wilcoxon test P=0.01; **Figure 2A**).

111

112 We previously demonstrated that insertion of an ISX element in the *D. melanogaster* genome 113 results in recruitment of the MSL complex to an ectopic autosomal location (Ellison and 114 Bachtrog, 2013). We used this same system to dissect the relationship between the TT alleles and 115 MSL complex binding affinity. Starting with a cloned ISX element (Ellison and Bachtrog, 2013), we used site-directed mutagenesis to create variants of this element that differ only with respect 116 117 to the TT haplotype. Each of the four possible haplotypes (GA, GT, TA, and TT) was engineered 118 and inserted onto D. melanogaster chromosome 2L at cytosite 38F1 using recombinase mediated 119 cassette exchange (RMCE) (Bateman et al., 2006). We then measured the effect of each of the 120 derived variants by quantifying allele-specific binding levels of the MSL complex in F1 hybrids 121 between the ancestral haplotype (GA) and each of the derived haplotypes (GT, TA, and TT). 122

123 Interestingly, each T allele, when assayed separately, has a negative effect on MSL binding

- 124 levels compared to the ancestral G or A allele (Figure 2B). However, when combined, the TT
- 125 haplotype results in significantly increased levels of MSL complex binding, relative to the

126 ancestral GA haplotype (Wilcoxon Test P = 0.0289; Figure 2B). These results suggest that there 127 is sign epistasis between the two alleles and that the high frequency TT haplotype represents a

refining/fine-tuning adaptation, since recruitment of MSL complex to the adjacent MSL

120 Terming the tuning adaptation, since recratation of the D complex to the adjacent

129 recognition motif is increased.

130

131 Non-allelic gene conversion is spreading the TT haplotype among ISX insertions

132 It is unlikely that the TT haplotype arose multiple times by parallel mutation, and there are two 133 possibilities that could explain its prevalence among MSH22 ISX insertions. First, this double 134 mutation may have occurred early during the process of ISX amplification, thus giving rise to 135 two lineages of ISX: one that carries the ancestral GA haplotype, and the other with the TT 136 haplotype. The TT-harboring elements in MSH22 would then all be descendants from the latter 137 ISX lineage. Alternatively, this mutation may have occurred only after the GA-containing ISX 138 element was fixed in the population at all 77 neo-X insertion sites, at which point it was spread 139 among independent ISX elements via non-allelic gene conversion.

140

We can distinguish between these possibilities by examining patterns of sequence polymorphism for each ISX insertion across multiple strains of *D. miranda*. A canonical signature of non-allelic gene conversion is the presence of shared polymorphisms across sequence duplicates (Arguello et al., 2006; Mansai and Innan, 2010). If gene conversion is spreading the TT haplotype among ISX insertions, we expect it to be polymorphic among individuals at several ISX insertion sites, whereas we do not expect the TT haplotype to be polymorphic at individual ISX insertions under the alternative scenario.

148

149 To genotype multiple wild-derived individuals at each of the MSH22 ISX insertions, we used 150 paired-end Illumina genomic resequencing data from 23 inbred lines of *D. miranda*, including 151 MSH22. We aligned all reads to the MSH22 reference genome and identified mate-pairs where 152 one mate was anchored in unique sequence flanking an ISX insertion. We then assembled these 153 reads to generate a contig spanning the 5' flank of the ISX element insertion, which contains the 154 MSL recognition motif, for each inbred line. Using this approach we generated population data 155 for 69 insertions out of the 77 total ISX insertions present in the MSH22 reference genome 156 assembly. Uneven sequence coverage between insertions and individuals meant that not all

insertions could be assembled for each individual. However, the majority of individuals are
represented in the majority of datasets: each insertion dataset contained ~20 lines on average (see
Dataset S1 in Dryad: Ellison & Bachtrog, 2015). Almost all ISX insertions are fixed among
strains (68 of 69) and insertion sites are identical between lines, suggesting that independent
parallel insertions are unlikely to be present within our dataset. We performed PCR and Sanger
sequencing on a subset of these regions and estimate the base-calling error rate of our Illumina
contigs to be ~0.1%.

164

165 Consistent with non-allelic gene conversion spreading the TT haplotype, we observe a strong 166 signal of allele-sharing within the sequence region flanking the MSL recognition motif among 167 ISX insertions (Figure 3). On average, 68.9% of polymorphisms observed within a given 168 insertion are shared among other insertions (though most polymorphisms are shared only 169 between a few elements). The TT haplotype is especially striking in this regard as it is 170 polymorphic in 41% of insertions (Figures 3, 4, & Figures 3-figure supplement). If population 171 subdivision contributes to this excess of allele sharing, we would expect individuals to cluster by 172 allele state at the TT locus, across all polymorphic ISX insertions. Instead, we find that different 173 individuals contribute to the TT polymorphism at each of these ISX insertions (Figure 4-figure 174 **supplement 1**), suggesting that abundant non-allelic gene conversion is the most likely 175 explanation for this observation. Interestingly, the population frequency of the TT haplotype is 176 similar among insertions that are near each other on the chromosome (permutation test P=0.018; 177 Figure 4). This is consistent with higher gene conversion rates between more closely linked ISX 178 elements generating correlated population frequencies among adjacent elements (Sasaki et al., 179 2010).

180

181 Selection is driving the spread of the TT haplotype through the population

182To test if selection has acted to increase the frequency of the TT haplotype in the population, we183examined patterns of polymorphisms at GA- and TT-containing ISX elements. The TT haplotype184harbors significantly less linked variation than the ancestral GA haplotype, across insertion sites

- and individuals (haplotype diversity = 0.53 vs. 0.81; resampling *P*<0.001; Figure 5A). In
- addition, ISX insertions where TT is fixed have significantly lower nucleotide diversity
- 187 compared to the insertions where GA is fixed (one-sided Wilcoxon test P=0.035; Figure

5B). Finally, the frequency spectrum at the TT haplotype also shows an excess of high frequency

derived alleles, compared to the frequency spectrum at the GA haplotype (resampling *P*=0.027;

Figure 5C). All of these patterns are expected if natural selection acting on the TT haplotype is

- 191 driving its spread through the population.
- 192

193 Discussion

194 Recent work in a variety of eukaryotes suggests that transposable elements may be major drivers 195 of regulatory evolution (Cowley and Oakey, 2013; Feschotte, 2008). Their high transposition 196 rate and ability to supply ready-to use regulatory elements across the genome implies that they 197 may rapidly wire new genes into regulatory networks (Feschotte, 2008). We recently showed 198 that domesticated TEs contribute to rewiring of the dosage compensation network in *D. miranda*, 199 but appear to supply only suboptimal binding sites for the MSL complex (Ellison and Bachtrog, 200 2013). Here, we identify a derived haplotype with two mutations that interact epistatically to 201 increase binding affinity for the MSL complex. We show that these fine-tuning mutations spread 202 among independent ISX insertions by non-allelic gene conversion, and through the population by 203 natural selection (Figure 6). Relative to regulatory elements that evolve in isolation, a family of 204 regulatory motifs dispersed by TEs presents a larger mutational target, and a mutation arising in 205 any element contained within a larger repeat unit (the TE) can spread among copies by non-206 allelic gene conversion. Consequently, the rate of evolutionary fine-tuning at such regulatory 207 elements can be greatly accelerated by increasing their effective population size (Mano and 208 Innan, 2008). Thus, transposable elements can "crowd-source" beneficial mutations to rapidly 209 fine-tune regulatory networks.

210

211 Our transgenic experiments show that each individual T allele actually decreases the binding 212 affinity for the MSL complex relative to the ancestral GA haplotype. Thus, TA or GT haplotypes 213 should be selected against in the population if present on a functional ISX element. Consistent 214 with the deleterious effect of individual T alleles, the TA and GT haplotypes are present on some 215 ISY elements but completely absent from ISX, i.e., we find the two T mutations to be in perfect 216 linkage disequilibrium among ISX elements but not ISY (Figure 1B). While most ISY elements 217 carry the ancestral GA haplotype, a small fraction (0.7% of neo-X ISY insertions) instead carry 218 the derived TT haplotype. It is therefore possible that the TT haplotype was introduced onto the

219 ISX background by non-allelic gene conversion from ISY. Under this scenario, the large family 220 of ISY elements in the *D. miranda* genome could be acting as a reservoir of natural variation, 221 where complex mutations can accumulate in the absence of epistasis. Non-allelic gene 222 conversion could then transfer these haplotypes to related repetitive elements (such as ISX). 223 While many of these haplotypes are likely to be neutral or deleterious, some may be beneficial, 224 as in the case of the TT haplotype. Such a scenario avoids the waiting time for a double 225 mutation, as well as the fitness valley that would have to be traversed if the two mutations were 226 to occur sequentially on the ISX background.

227

228 To conclude, our findings suggest that TE-dispersed binding motifs follow an evolutionary

trajectory that is fundamentally different from those that arise by other means. The

230 complementary roles of TEs in dispersing regulatory motifs, and gene conversion in spreading

subsequent refining mutations, combine to allow for the rapid rewiring and fine-tuning of gene

regulatory networks. This process adds a new layer of complexity onto how TEs influence

regulatory innovation, as well as a new context in which gene conversion affects genome

- evolution.
- 235

236 References

- Alekseyenko, A.A., Ellison, C.E., Gorchakov, A.A., Zhou, Q., Kaiser, V.B., Toda, N., Walton, Z., Peng, S., Park,
 P.J., Bachtrog, D., *et al.* (2013). Conservation and de novo acquisition of dosage compensation on newly
 evolved sex chromosomes in Drosophila. Genes Dev 27, 853-858.
- Alekseyenko, A.A., Peng, S., Larschan, E., Gorchakov, A.A., Lee, O.K., Kharchenko, P., McGrath, S.D., Wang,
 C.I., Mardis, E.R., Park, P.J., *et al.* (2008). A sequence motif within chromatin entry sites directs MSL
 establishment on the Drosophila X chromosome. Cell *134*, 599-609.
- Arguello, J.R., Chen, Y., Yang, S., Wang, W., and Long, M. (2006). Origination of an X-linked testes chimeric gene
 by illegitimate recombination in Drosophila. PLoS Genet 2, 745-754.
- Bachtrog, D., and Charlesworth, B. (2002). Reduced adaptation of a non-recombining neo-Y chromosome. Nature
 416, 323-326.
- Bateman, J.R., Lee, A.M., and Wu, C.T. (2006). Site-specific transformation of Drosophila via phiC31 integrasemediated cassette exchange. Genetics *173*, 769-777.
- 249 Bentley, D.R., Balasubramanian, S., Swerdlow, H.P., Smith, G.P., Milton, J., Brown, C.G., Hall, K.P., Evers, D.J.,
- Barnes, C.L., Bignell, H.R., *et al.* (2008). Accurate whole human genome sequencing using reversible
 terminator chemistry. Nature (London) 456, 53-59.
- Bone, J.R., and Kuroda, M.I. (1996). Dosage compensation regulatory proteins and the evolution of sex

- chromosomes in Drosophila. Genetics *144*, 705-713.
- Bradley, R.K., Roberts, A., Smoot, M., Juvekar, S., Do, J., Dewey, C., Holmes, I., and Pachter, L. (2009). Fast
 statistical alignment. PLoS Comput Biol *5*, e1000392.
- Carvalho, A.B., and Clark, A.G. (2005). Y chromosome of D. pseudoobscura is not homologous to the ancestral
 Drosophila Y. Science *307*, 108-110.
- Chen, J.-M., Cooper, D.N., Chuzhanova, N., Ferec, C., and Patrinos, G.P. (2007). Gene conversion: mechanisms,
 evolution and human disease. Nature Reviews Genetics *8*, 762-775.
- Cowley, M., and Oakey, R.J. (2013). Transposable elements re-wire and fine-tune the transcriptome. PLoS Genet 9, e1003234.
- Eickbush, T.H., and Eickbush, D.G. (2007). Finely orchestrated movements: Evolution of the ribosomal RNA genes.
 Genetics *175*, 477-485.
- Ellison, C.E., and Bachtrog, D. (2013). Dosage Compensation via Transposable Element Mediated Rewiring of a
 Regulatory Network. Science (Washington D C) *342*, 846-850.
- Ellison, C.E. and Bachtrog, D. (2015) Data from: Non-allelic gene conversion enables rapid evolutionary change at
 multiple regulatory sites encoded by transposable elements. Dryad Digital Repository.
- 268 doi:10.5061/dryad.dg483
- Feschotte, C. (2008). Transposable elements and the evolution of regulatory networks. Nature Reviews Genetics 9,
 397-405.
- Hall, G.S., and Little, D.P. (2007). Relative quantitation of virus population size in mixed genotype infections using
 sequencing chromatograms. Journal of Virological Methods *146*, 22-28.
- 273 Langmead, B., and Salzberg, S.L. (2012). Fast gapped-read alignment with Bowtie 2. Nat Methods 9, 357-359.
- Mano, S., and Innan, H. (2008). The evolutionary rate of duplicated genes under concerted evolution. Genetics *180*,
 493-505.
- Mansai, S.P., and Innan, H. (2010). The Power of the Methods for Detecting Interlocus Gene Conversion. Genetics
 184, 517-U292.
- Marin, I., Franke, A., Bashaw, G.J., and Baker, B.S. (1996). The dosage compensation system of Drosophila is coopted by new evolved X chromosomes. Nature *383*, 160-163.
- 280 Ohta, T. (2010). Gene conversion and evolution of gene families: an overview. Genes (Basel) 1, 349-356.
- Peng, Y., Leung, H.C.M., Yiu, S.M., and Chin, F.Y.L. (2012). IDBA-UD: a de novo assembler for single-cell and
 metagenomic sequencing data with highly uneven depth. Bioinformatics (Oxford) 28, 1420-1428.
- Polavarapu, N., Marino-Ramirez, L., Landsman, D., McDonald, J.F., and Jordan, I.K. (2008). Evolutionary rates and
 patterns for human transcription factor binding sites derived from repetitive DNA. BMC Genomics *9*, Article
 No.: 226.
- Richards, S., Liu, Y., Bettencourt, B.R., Hradecky, P., Letovsky, S., Nielsen, R., Thornton, K., Hubisz, M.J., Chen,
 R., Meisel, R.P., *et al.* (2005). Comparative genome sequencing of Drosophila pseudoobscura: chromosomal,
 gene, and cis-element evolution. Genome Res *15*, 1-18.
- Rozen, S., Skaletsky, H., Marszalek, J.D., Minx, P.J., Cordum, H.S., Waterston, R.H., Wilson, R.K., and Page, D.C.

| 290 | (2003). Abundant gene conversion between arms of palindromes in human and ape Y chromosomes. Nature |
|------------|---|
| 291 | (London) 423, 873-876. |
| 292 | Sasaki, M., Lange, J., and Keeney, S. (2010). Genome destabilization by homologous recombination in the germ |
| 293 | line. Nature Reviews Molecular Cell Biology 11, 182-195. |
| 294 | Stewart, A.J., Hannenhalli, S., and Plotkin, J.B. (2013). Why Transcription Factor Binding Sites Are Ten |
| 295 | Nucleotides Long. Genetics 193. |
| 296 207 | Straub, T., Grimaud, C., Gilfillan, G.D., Mitterweger, A., and Becker, P.B. (2008). The Chromosomal High-Affinity |
| 297 | Binding Sites for the Drosophila Dosage Compensation Complex. PLoS Genet 4, Article No.: e1000302. |
| 298 299 | sequence alignment editor and analysis workbench. Bioinformatics (Oxford) 25, 1189-1191. |
| 300 | Zangenberg, G., Huang, MM., Arnheim, N., and Erlich, H. (1995). New HLA-DPB1 alleles generated by |
| 301 | interallelic gene conversion detected by analysis of sperm. Nat Genet 10, 407-414. |
| 302 | Zhou, Q., Ellison, C.E., Kaiser, V.B., Alekseyenko, A.A., Gorchakov, A.A., and Bachtrog, D. (2013). The |
| 303 | Epigenome of Evolving Drosophila Neo-Sex Chromosomes: Dosage Compensation and Heterochromatin |
| 304 | Formation. PLoS Biology 11, Article No.: e1001711. |
| 305 | |
| 306 | |
| 307 | Acknowledgments: This work was funded by NIH grants (R01GM076007 and R01GM093182) to D.B. and a NIH |
| 308 | postdoctoral fellowship to C.E.E. All DNA-sequencing reads generated in this study are deposited at the National |
| 309 | Center for Biotechnology Information Short Reads Archive (<u>www.ncbi.nlm.nih.gov/sra</u>) under the BioProject |
| 310 | ID PRJNA270105. We thank Molly Przeworski, Jeffrey Fawcett, Isabel Gordo and Monty Slatkin for comments on |
| 311 | the manuscript and Daniel Weissman for helpful discussions. |
| 312 | |
| 313 | Competing interests : The authors have no competing financial and non-financial interests. |
| 314 | |
| 315 | |
| 316 | Materials and Methods |
| 317 | |
| 318 | Resequencing of D. miranda wild lines |
| 319 | Isofemale lines were established from individuals collected in Northern California and inbred for |
| 320 | several generations. DNA was extracted from 1-8 females per line using the Qiagen PureGene |
| 321 | kit and fragmented by nebulization. Paired-end Illumina libraries were constructed using |
| 322 | standard protocols (Bentley et al., 2008) and sequenced on an Illumina Genome Analyzer II |
| 323 | machine. |
| 324 | |

325 ISX assembly and variant identification

326 Resequencing data were mapped to version 2.2 of the *D. miranda* MSH22 reference assembly 327 (GenBank: AJMI00000000.2) using bowtie2 (Langmead and Salzberg, 2012). ISX locations 328 were identified in (Ellison and Bachtrog, 2013). Paired-end read alignments were evaluated 329 within 2 kilobase windows flanking each ISX insertion and reads with mapping quality of 20 or 330 greater were extracted along with their mate. The extracted mate pairs were then assembled 331 using IDBA-UD, for each line separately (Peng et al., 2012). Contigs were aligned using FSA 332 (Bradley et al., 2009) and visualized with Jalview (Waterhouse et al., 2009). A custom Perl script 333 (available at https://github.com/chris-ellison/MSAvariants) was used to identify sequence 334 variants within the alignments. We also PCR amplified eight of the ISX insertions where the TT 335 haplotype was polymorphic. We confirmed that this polymorphism was present at each of these 336 insertions and estimated the base-calling accuracy of the assemblies by sequencing the PCR 337 products using Sanger technology.

338

339 Transgenesis

We used the QuikChange Lightning site-directed mutagenesis kit from Agilent Technologies and
the ISX element cloned in (Ellison and Bachtrog, 2013) to engineer four ISX variants that
differed only with respect to the TT haplotype: ISX-GA, ISX-GT, ISX-TA, and ISX-TT. Each
construct was injected by BestGene Inc. (Chino Hills) into *D. melanogaster* embryos carrying a
RMCE landing site at cytosite 38F1 on chromosome 2L (Bloomington Drosophila Stock Center

- 345 strain #27388). Transformants were verified by PCR and Sanger sequencing.
- 346

347 Quantification of allele-specific binding levels of the MSL complex

- 348 Male third instar larvae (~250 mg) were collected from F1 hybrids between ISX-GA and each of
- the other three engineered lines: ISX-GT, ISX-TA, and ISX-TT. Chromatin immunoprecipitation
- 350 was performed for four biological replicates of each of these lines using the MSL2 d-300
- 351 primary antibody from Santa Cruz Biotechnology Inc. and the protocol described in
- 352 (Alekseyenko et al., 2013). Primers flanking the ISX MRE region were used to generate
- 353 heterozygous amplicons from the MSL2 IP and input control. Sanger chromatograms were used
- in conjunction with polySNP software (Hall and Little, 2007) to calculate relative abundance of
- 355 ISX alleles within the IP and input control amplicons. Abundance of the 'T' alleles in the IP

- amplicons relative to the ancestral G/A alleles was calculated and normalized by the same valuesfrom the input control.
- 358

359 **Permutation and resampling tests**

360 To determine if the TT frequency among neighboring ISX elements was correlated, we clustered

361 elements within 100kb of each other and calculated the standard deviation in TT allele frequency

362 within clusters. We then compared these values to 1000 permutations where TT allele frequency

- 363 was randomly shuffled between ISX locations.
- 364 The haplotype diversity and allele frequency spectrum resampling tests were performed by
- drawing, without replacement, two groups of size 617 and 674, respectively, from the pool of

366 1,291 ISX sequences. The intergroup difference in haplotype diversity, as well as the number of

derived variants with frequency of 0.75 or greater, was calculated for each of 1000 replicates and

368 compared to the difference between the TT and GA groups.

369

370

- 371 Figure Legends
- 372

Figure 1. TE-derived MSL recognition element (MRE) motifs from the neo-X chromosome of *Drosophila miranda*

375 (A) The MSL recognition motif (MRE) plus 20 basepairs of flanking sequence were extracted from all 77 ISX 376 transposable elements located on the neo-X chromosome in the MSH22 reference genome assembly. The multiple 377 sequence alignment of these 77 sequence regions (arranged from top-to-bottom in the order in which they are found 378 on the chromosome) shows that there is sequence variation among elements both within and adjacent to the 21 379 basepair MRE motif. Each variant has been classified as ancestral or derived based on its frequency in the ISX 380 progenitor element, ISY. The derived allele frequency for each variant in this region is shown for ISX as well as 139 381 ISY elements from the neo-X chromosome (see Figure 1-figure supplement for ISY alignment). Red arrows point 382 to the derived TT haplotype that is common among ISX elements but rare in ISY. (B) Barplot showing the 383 frequencies of all haplotypes at the GA/TT locus, for ISY and ISX elements separately. Two haplotypes are present 384 within ISX elements (GA and TT) and the two alleles within each haplotype are in perfect linkage disequilibrium. In 385 contrast, the majority of ISY elements harbor the GA haplotype, but these two alleles are not in perfect linkage 386 disequilibrium among ISY elements. Rather, five additional allelic combinations are present at low frequencies in 387 this location among ISY, but not ISX elements.

388

389 Figure 1 - Supplement. Alignment of ISY elements from the D. miranda MSH22 genome assembly

390 139 ISY elements from the MSH22 neo-X chromosome were identified and 200 basepairs from their 5' flanks were391 aligned. The black arrows point to the sites where the derived 'T' alleles are common among ISX elements. In

- 392 contrast, only a single ISY element from the neo-X chromosome harbors the TT haplotype.
- 393

394 Figure 2. The TT haplotype increases MSL binding affinity

(A) MSL3 ChIP-seq data from *D. miranda* strain MSH22 shows that the ISX insertions carrying the TT haplotype
 recruit significantly higher levels of MSL complex compared to those with the GA haplotype (Wilcoxon test

397 *P*=0.01). (B) Engineered ISX elements that differ only with respect to the TT haplotype bind different levels of MSL

398 complex. There is an epistatic interaction between the two 'T' alleles such that separately, they decrease MSL

- 399 complex binding relative to the ancestral allele, but together in the TT haplotype, they increase MSL complex
- 400 binding (Wilcoxon Test *P*=0.028 for both comparisons [GT *vs* TT and TA *vs* TT]). The rectangles and error bars
- 401 show the average and standard deviation of values from four biological replicates for each condition.
- 402

403 Figure 3. ISX variation among wild lines of *D. miranda*

404 For each ISX insertion identified within the *D. miranda* MSH22 reference genome assembly (alignment shown at

- 405 left, see also Figure 1), we characterized sequence variation across *D. miranda* individuals. The TT haplotype
- 406 (magenta lines) was fixed across individuals at 30% of insertions (see example alignment, top right), polymorphic at
- 407 41% of insertions (example shown middle right), and absent at 29% of insertions (bottom right). Allele sharing

- 408 between insertions occurs at sites other than the TT haplotype, but these sites tend to be shared across fewer
- 409 insertions (see heatmap, bottom right). Figure 3-figure supplement shows the population alignment across all ISX
- 410 insertions on the neo-X.
- 411
- 412

413 Figure 3 - Supplement. Shared polymorphism across sixty-nine ISX insertions

- 414 The 5' 200 basepairs of the ISX element was assembled for an average of 20 individuals, for each of 69 ISX
- 415 insertions. Each stripe corresponds to the population data for a given insertion and nucleotides are colored as in
- 416 Figure 1. Solid lines point to columns of the alignment containing polymorphisms that are shared between multiple
- 417 ISX insertions. For these columns, the heatmap is shaded to reflect the degree of allele-sharing, which ranges from
- 418 3% of insertions to 41% of insertions. The 'T' letters under the heatmap mark the location of the TT haplotype
- shown in Figure 1.
- 420

421 Figure 4. Population frequency of TT haplotype across ISX insertions

- 422 The location of all ISX elements on the *D. miranda* neo-X chromosome, as inferred from the the MSH22 reference
- 423 genome assembly, is shown by vertical green bars. The derived TT haplotype (frequency shown in red), is
- 424 polymorphic at 27 of 69 ISX insertions, a pattern consistent with non-allelic gene conversion.
- 425

426 Figure 4 – Supplement. ISX genotype across insertions and individuals

- 427 Heatmap showing each of the 27 ISX insertions (rows) where the TT haplotype is polymorphic among individuals.
- 428 Columns show the genotype of each individual, for each of these insertions. Each individual has a mixture of TT and
- 429 GA ISX insertions, suggesting that TT polymorphism among lines is not due to population subdivision.
- 430

431 Figure 5 Selection shapes patterns of variation at the TT haplotype.

- 432 (A) Haplotype diversity across all ISX sequences. Assembled ISX contigs were combined for all insertions and
- 433 individuals. The 25 basepairs flanking each side of the TT region were extracted from a total of 1,291 sequences and
- 434 split into two groups based on whether they contained the TT or GA haplotype. Haplotype diversity was then
- 435 calculated for each group. The difference between groups is significantly larger than expected by chance
- 436 (resampling P < 0.001), with the sequences containing the TT haplotype having less haplotype diversity compared to
- those containing the GA haplotype. (B) Nucleotide diversity across all ISX sequences. We compared nucleotide
- 438 diversity for ISX insertions where all individuals carried the ancestral GA haplotype to those where the derived TT
- 439 haplotype was fixed. ISX insertions that are fixed for the TT haplotype have significantly reduced nucleotide
- 440 diversity compared to insertions fixed for the GA haplotype (one-sided Wilcoxon test P=0.035). (C) Allele-
- 441 frequency spectrum across ISX sequences. The allele frequency spectrum was calculated separately for TT and GA-
- 442 carrying ISX sequences, across all insertions and individuals. Consistent with incomplete hitchhiking under positive
- selection, the TT frequency spectrum shows an excess of high frequency derived alleles, compared to the GA
- 444 spectrum (resampling *P*=0.027)..

445

446 Figure 6. Non-allelic gene conversion spreads refining mutations among TE-derived MRE motifs

- 447 Shared polymorphism of the TT haplotype among ISX insertions suggests a model where a mutation that refines
- regulatory activity arose once at a single TE-derived regulatory element, and spread across elements via non-allelic
- gene conversion. Over evolutionary time, such a mutation spreads in two dimensions: horizontally among TE-
- 450 derived regulatory elements and vertically through the population, until it is fixed across elements and across
- 451 individuals. The TT haplotype is at the midpoint of this process. Across ISX insertions, it is fixed, absent, and
- 452 polymorphic, in approximately equal proportions.
- 453











