

RESEARCH ARTICLE

Open Access

Genomic analysis of the four ecologically distinct cactus host populations of *Drosophila mojavensis*



Carson W. Allan^{1,2} and Luciano M. Matzkin^{1,2,3,4*} 

Abstract

Background: Relationships between an organism and its environment can be fundamental in the understanding how populations change over time and species arise. Local ecological conditions can shape variation at multiple levels, among these are the evolutionary history and trajectories of coding genes. This study examines the rate of molecular evolution at protein-coding genes throughout the genome in response to host adaptation in the cactophilic *Drosophila mojavensis*. These insects are intimately associated with cactus necroses, developing as larvae and feeding as adults in these necrotic tissues. *Drosophila mojavensis* is composed of four isolated populations across the deserts of western North America and each population has adapted to utilize different cacti that are chemically, nutritionally, and structurally distinct.

Results: High coverage Illumina sequencing was performed on three previously unsequenced populations of *D. mojavensis*. Genomes were assembled using the previously sequenced genome of *D. mojavensis* from Santa Catalina Island (USA) as a template. Protein coding genes were aligned across all four populations and rates of protein evolution were determined for all loci using a several approaches.

Conclusions: Loci that exhibited elevated rates of molecular evolution tend to be shorter, have fewer exons, low expression, be transcriptionally responsive to cactus host use and have fixed expression differences across the four cactus host populations. Fast evolving genes were involved with metabolism, detoxification, chemosensory reception, reproduction and behavior. Results of this study give insight into the process and the genomic consequences of local ecological adaptation.

Keywords: Genome evolution, Adaptation, *Drosophila*, Ecological genomics, Genome sequencing, Genome assembly, *Drosophila mojavensis*

Background

Increasing availability of whole-genome sequencing data provides new insights into the complex relationship between an organism and its environment. By examining changes in the genetic code both at the level of individual genes and at the whole-genome level it is possible to gain a better understanding of how local ecological conditions can shape the pattern of variation within and between ecologically distinct populations [1, 2]. A comprehensive integrative

approach combining genomic, phenotypic data has been identified as the gold standard in understanding the adaptation process [3, 4]. Yet, an examination of the genomic divergence of ecologically distinct populations can yield valuable insight into the adaptation process especially when the genomic data is placed in an ecological context [5]. This later approach can identify genomic regions and loci that exhibit a pattern of variation and evolution suggesting their role in local ecological adaptation. Furthermore, a consequence of the fixation of ecologically-relevant variants has been implicated in the evolution of barriers to gene flow and potentially the origins of reproductively isolated populations, i.e. species [6, 7].

While it has long been accepted that natural selection is a primary driver of change within species as a response to

* Correspondence: lmatzkin@email.arizona.edu

¹Department of Biological Sciences, University of Alabama in Huntsville, 301 Sparkman Drive, Huntsville, AL 35899, USA

²Department of Entomology, University of Arizona, 1140 E. South Campus Drive, Tucson, AZ 85721, USA

Full list of author information is available at the end of the article



environmental pressures, understanding the mechanism of how this selection leads to speciation is unclear [8, 9]. More recently the idea of ecological speciation, where various mechanisms work to prevent gene flow between populations causing reproductive isolation and eventually speciation, has more directly shown how selection to local ecological conditions may affect the process of speciation [6, 7]. Reproductive isolation interrupts gene flow between populations and may potentially lead to the formation of new species [10]. When different populations of a species inhabits and/or utilizes distinct resources this opens many possibilities for local differentiation that can lead to obstacles of gene flow as these populations are likely to have differing environmental pressures [6, 7]. For example, in the leaf beetle *Neochlamisus bebbianae*, different populations have distinct host preferences and larvae perform significantly worse when growing on alternative host species [8]. Host preferences and performance in this system facilitates the genetic and genomic isolation observed between the host populations, as each prefers a different microenvironment and likely does not interact and hybridize with members of the other population [11, 12].

Comparative genomic studies in mammals have shown clear evidence of positive selection both between humans, mice, and chimpanzees as well as between human populations [13–16]. Genes involved in the immune system, gamete development, sensory perception, metabolism, cell motility, and genes involved with cancer were those found to have signatures of positive selection. While in *Drosophila*, a genome level analysis of 12 species provided insight into the evolution of an ecological, morphological, physiological and behaviorally diverse genus [17]. Findings were relatively consistent with previously studies in other taxa with genes involving defense, chemosensory perception, and metabolism shown to be under positive selection [6, 13, 16, 18]. Since the *Drosophila* 12 genome project [17], several population genomics studies in *D. melanogaster* have examined variation within a single population, between clinal populations and between ancestral (African) and cosmopolitan populations to assess the consequence of population subdivision, evolution of quantitative trait variation and the adaptation to local ecological conditions [19–24]. These genome level analysis have been extended to other *D. melanogaster* species group flies with distinct life history and ecological strategies such as the *Morinda citrifolia* specialist *D. sechellia* [25] and the invasive agricultural pest *D. sukuzii* [26].

Studying the sequence level constraints as well as functional categories and networks associated with genes under positive selection is paramount to understanding the process of evolutionary change. However, it is crucial to place patterns of variation and divergence in an ecological context to have a more complete view how selection shapes variation within and between populations. In this study we explore the link between ecology and patterns of genome-

wide sequence variation in *D. mojavensis*, a fly endemic to the southwestern United States and northwestern Mexico that has become a model for the understanding of the genetics of adaptation [27]. This species of *Drosophila* is a cactophile in that both larval and adult stages reside and feed on necrotic cactus tissues [28]. *Drosophila mojavensis* has four distinct host populations that are geographically separated (Fig. 1). Individuals from all four populations can interbreed with each other and produce viable fertile offspring, and no postzygotic incompatibilities appear to exist, although some evidence indicate low levels of prezygotic isolation between some of the populations [29]. In addition to geographic separation each population utilizes a distinct necrotic cactus host species. The four populations and their respective main cactus host are: Santa Catalina Island living on prickly pear cactus (*Opuntia littoralis*), Mojave Desert living on barrel cactus (*Ferocactus cylindraceus*), Baja California living on agria cactus (*Stenocereus gummosus*), and Sonoran Desert living on organpipe cactus (*S. thurberi*). *Drosophila mojavensis* diverged from its sister species *D. arizonae*, a cactus generalist, approximately half a million years ago [30–33] with the divergence between *D. mojavensis* populations being more recent (230,000 to 270,000 years ago) [34]. Differing host species provide different local environments for each *D. mojavensis* populations. The necrotic cactus environment in which these flies reside is composed not only of plant tissues, but a number of bacteria and yeast species which are necessary for providing nutrition as well as playing a role in metabolizing cactus-derived compounds

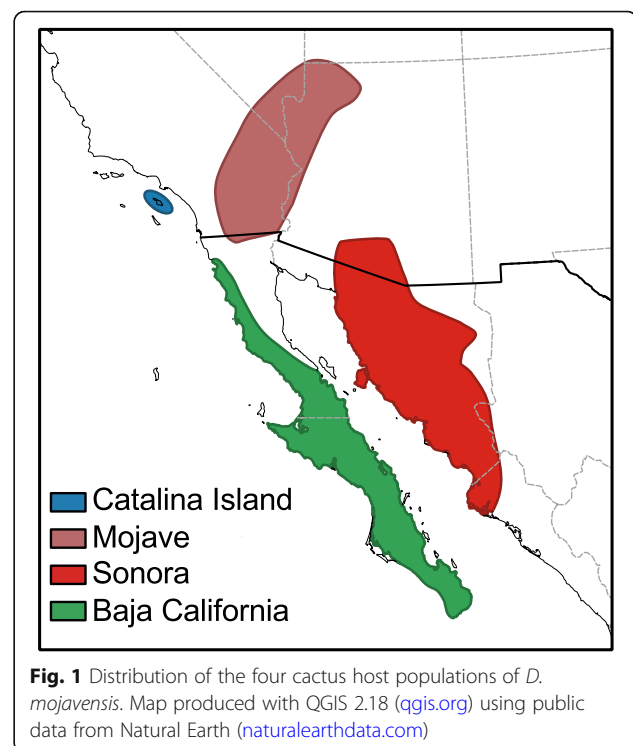


Fig. 1 Distribution of the four cactus host populations of *D. mojavensis*. Map produced with QGIS 2.18 (qgis.org) using public data from Natural Earth (naturalearthdata.com)

[35–38]. In addition to the nutritional differences of the necroses that exist between the distinct cactus species, the rots contain a number of compounds that have toxic properties which can affect the performance and viability of non-resident cactophilic *Drosophila* [39–41]. This selective pressure has resulted in the fixation of variants that facilitate the survival of *D. mojavensis* and other cactophilic *Drosophila* species to their local necrotic cactus environment [28, 42].

Population genetics on individual candidate host adaptation genes in *D. mojavensis* has shown evidence for positive selection in loci involved with xenobiotic metabolism [32]. In addition, transcriptome-wide differences have been observed in *D. mojavensis* in response to host shifts [43, 44] as well as indicating fixed expression differences between the host populations [45]. Among the loci that are differentially expressed or constitutively fixed between populations many are involved in detoxification, metabolism, chemosensory perception and behavior, supporting the role of the local necrotic cactus conditions in shaping transcriptional variation [43–45]. Taking into consideration the breadth of ecological information of *D. mojavensis* this study highlights how selection pressures caused by local ecological environments differentially shape patterns of genomic variation across the host populations and provides further insight into how selection acts on organisms and its genome level consequences.

Results

Number of cleaned reads and the number assembled to the Catalina Island reference genome are shown in Table 1. All three populations had approximately 88% of paired-end reads successfully assembled. Mate pair reads had lower rates of mapping ranging from 27% to 63%. Upon subsequent inspection of our reads, we determined that some of

Table 1 Number of cleaned reads and assembled reads for each population

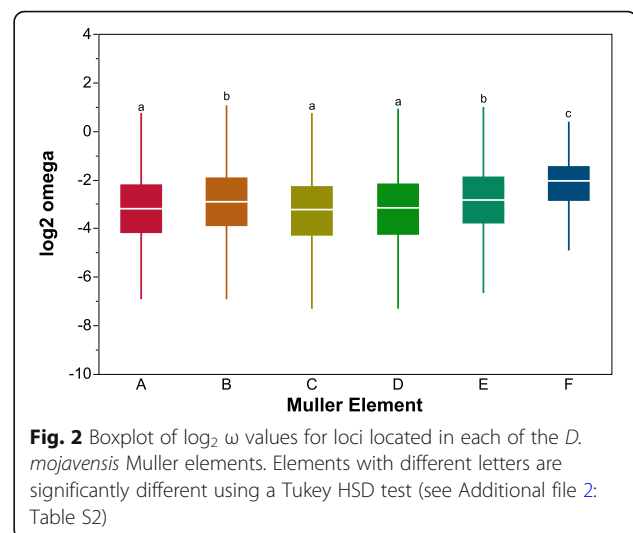
Population	Reads Mapped	Total Reads	Proportion Mapped
Baja California			
ME	12,052,662	44,912,130	0.27
PE	88,976,029	100,263,663	0.89
Total	101,028,691	145,175,793	0.70
Mojave			
ME	26,638,794	52,910,406	0.50
PE	73,196,313	83,000,942	0.88
Total	99,835,107	135,911,348	0.73
Sonora			
ME	39,962,094	63,240,890	0.63
PE	93,857,309	105,723,406	0.89
Total	133,819,403	168,964,296	0.79

ME Mate pair end reads, PE Paired end reads

our mate pair data (e.g. Baja California, see Table 1) contained high amounts of non-*Drosophila* contamination, but the mapping stringency we utilized would have prevented these contaminating sequences from mapping to the reference genome. Of the 14,680 loci annotated in the reference genome the vast majority were also present in our template-based assemblies of the other three populations. Of these annotations, a common set of 12,695 were initially processed that did not lack any premature stop codons. From this common set of loci we filtered out those that among the four populations exhibited either less than five total, zero nonsynonymous, or zero synonymous substitutions. The purpose for this filtering was to reduce the number of estimates of divergence (K_a/K_s) of low statistical confidence due to just a few mutations. This yielded a working set of 9087 loci for which all subsequent analyses were performed. The list of all loci examined, summary data, test statistics, and *D. melanogaster* ortholog information can be found in Additional file 1: Table S1.

Characteristics and patterns of divergence of *D. mojavensis* loci

Estimates of ω (K_a/K_s) were calculated using both KaKs Calculator [46] and codeml in PAML [47]. Given that the ω values were highly correlated ($r^2 = 0.88$, $P < 0.001$; see Additional file 2: Figure S1) all subsequent analyses were performed using the values obtained from codeml. The distribution of \log_2 transformed ω are shown in Additional file 2: Figure S2. Overall a total of 190 loci exhibited ω values greater than one. When examined per chromosome (Muller Element), we observed that the dot chromosome (Muller F) had the greatest mean ω , followed by the chromosomes for which segregate chromosomal inversions (Muller B and E) and than those chromosomes that lack inversions (Muller A, C and D) (Fig. 2, Additional file 2: Table S2).



To describe the characteristics of loci whose evolutionary trajectory could have been shaped by the adaptation of *D. mojavensis* populations to their respective ecological conditions we examined loci with ω values in the top 10% of the distribution, hereafter referred to as TOP10 loci. Furthermore, using codeml we performed a series of gene-wide tests of positive selection for each individual locus. Via a maximum likelihood rate test (model 7 vs. model 8) we identified 912 loci that exhibited a pattern of adaptive protein evolution. We used a smaller set of 244 loci, following an FDR correction, for all subsequent analyses, hereafter referred to as PAML-FDR loci. The set of TOP10, PAML significant loci and those with an FDR correction (PAML-FDR) can be found in Additional file 1: Table S1. The distribution of both the PAML-FDR and TOP10 loci was uniform across the *D. mojavensis* chromosomes (Additional file 2: Figures S3 and S4), with the exception that significantly fewer PAML-FDR genes were present in Muller E (Fisher's Exact test, $P = 0.02$).

Significant differences in ω values were observed across loci of differing protein coding lengths (Fig. 3). Loci smaller than 1 Kb exhibit significantly higher rate of molecular evolution, followed by those 1–2 Kb and then by gene categories of longer lengths (Additional file 2: Table S3). A similar pattern of ω values was observed for the TOP10 loci, where a significant excess of the smaller gene group (<1 Kb) was composed of TOP10 loci, and a significantly fewer were observed in the greater than 4 Kb bin (Additional file 2: Figure S5). Although the overall ω was greater in shorter loci, the proportion of these loci who exhibited a significant pattern of positive selection was significantly less (Additional file 2: Figure S6). Similarly to what was observed for gene length, genome-wide, loci with fewer exons tended to have greater levels of ω , with the highest observed from loci having two exons, then those with either only one or three exons, followed by those having four to six exons and lastly

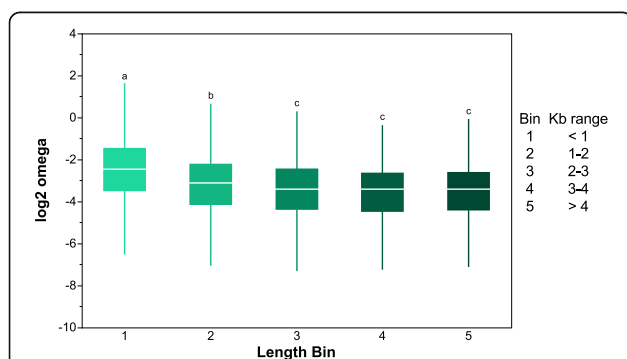


Fig. 3 Boxplot of $\log_2 \omega$ values for loci in five different coding length bins. Bins with different letters are significantly different using a Tukey HSD test (see Additional file 2: Table S3)

those with seven or more (Additional file 2: Figure S7, Table S4). TOP10 loci were overrepresented in the one and two exon categories and underrepresented in the more than seven exon category, whereas the PAML-FDR loci were uniformly distributed across all exon number categories (Additional file 2: Figures S8 and S9).

Relationship between expression and rate of molecular evolution

To assess the relationship between expression level and rate of molecular evolution we integrated our results with previous collected RNAseq data from *D. mojavensis* [48]. When examined genome-wide, genes with male-biased expression had significantly greater ω values than female-biased (Tukey HSD, $P < 0.001$) and unbiased (Tukey HSD, $P < 0.001$) expressed genes, and female-biased genes had the lowest rate (Tukey HSD, $P < 0.001$) of molecular evolution of all three expression categories (Additional file 2: Figure S10, Table S5). Among the TOP10 loci, there was a significant representation of them in the male-biased group of genes and a significant underrepresentation in the female-biased genes (Fig. 4). No significant over- or underrepresentation was observed among the PAML-FDR genes with respect to the sex biased expression categories (Additional file 2: Figure S11). Expression data was also used to assess the relationship between overall expression level and rate of molecular evolution. After removing both the female- and male-biased genes, we observed that of the 5101 remaining loci those in the lowest expression category showed the greatest ω values (Additional file 2: Figure S12, Table S6). Similarly, the TOP10 loci were overrepresented among the low expression category of loci and no differences were observed among the expression

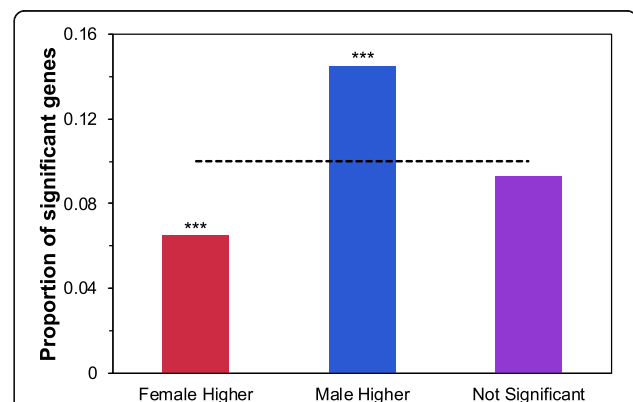


Fig. 4 Proportion of TOP10 loci that show female-bias, male-bias or unbiased gene expression. Dashed line indicates the genome wide proportion of TOP10 loci (0.10). Gene expression data is from [48]. Asterisk indicate significance via Fisher's Exact test (* $P < 0.05$, ** $P < 0.01$, *** $P < 0.001$)

categories of the PAML-FDR loci (Additional file 2: Figures S13 and S14).

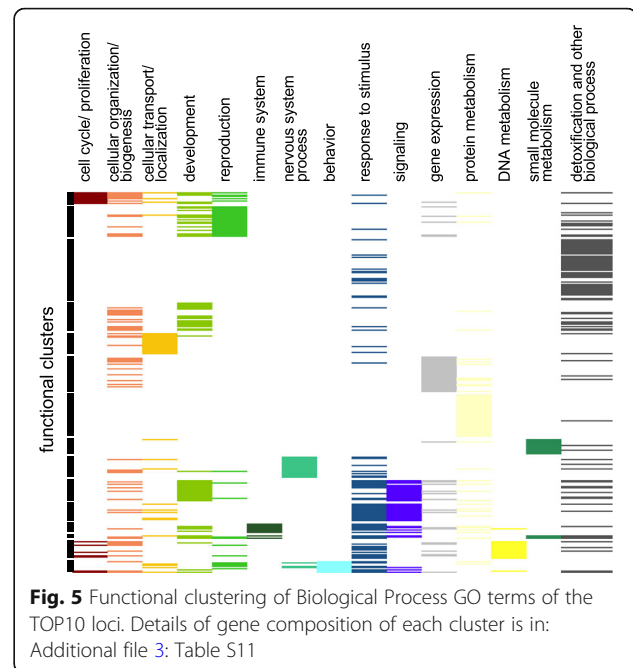
We also integrated our genomic data with two prior ecological transcriptional studies. We compare rates of molecular evolution of loci that are differentially expressed in response to cactus host utilization [44] as well as those loci who exhibit fixed significant expression differences between the four host populations in the absence of cactus compounds (i.e. constitutive differences) [45]. To remove the potential confounding effect of those loci that show a pattern of positive selection, we removed those loci from the subsequent expression analysis. For both datasets, loci that are either differentially expressed in response to necrotic cactus ($P < 0.001$ post FDR correction) or those that show constitutive differences between the populations ($P < 0.001$ post FDR correction) have a significantly greater value of ω (ANOVA, $P < 0.001$, for both comparisons) (Additional file 2: Figure S15, Table S7).

Functional gene groups analysis

Of our 9087 genes in our filtered dataset, approximately 14% (1238) genes did not have orthologous calls back to loci in the *D. melanogaster* reference genome (Additional file 2: Figure S16). Of the remaining set of genes with *D. melanogaster* orthologs, less than half of the genes (3649) had at least one gene ontology (GO) term. The percentage of loci without *D. melanogaster* orthologous in the TOP10 and PAML-FDR genes was greater (40 and 23%, respectively). Overall only 336 and 144 loci had at least one GO term for the TOP10 and PAML-FDR datasets, respectively. Clustering of biological process and molecular function GO terms within the TOP10 and PAML-FDR dataset illustrated some distinct functional groups. Figure 5 illustrates the biological process functional clusters for TOP10 genes, in which clusters associated with reproduction/development, detoxification and response to stimuli, and behavior are present. A network analysis of the same set of loci indicates similar functional networks as well as those associated with defense and chromatin regulation and remodeling (Fig. 6). Functional and network clustering for molecular function GO terms, KEGG and the PAML-FDR dataset can be found in Additional file 2: Figures S17-S20, Additional file 3: Table S11. Among molecular functions, in the TOP10 dataset, serine endopeptidase activity appeared to be overrepresented (Additional file 2: Table S8).

Discussion

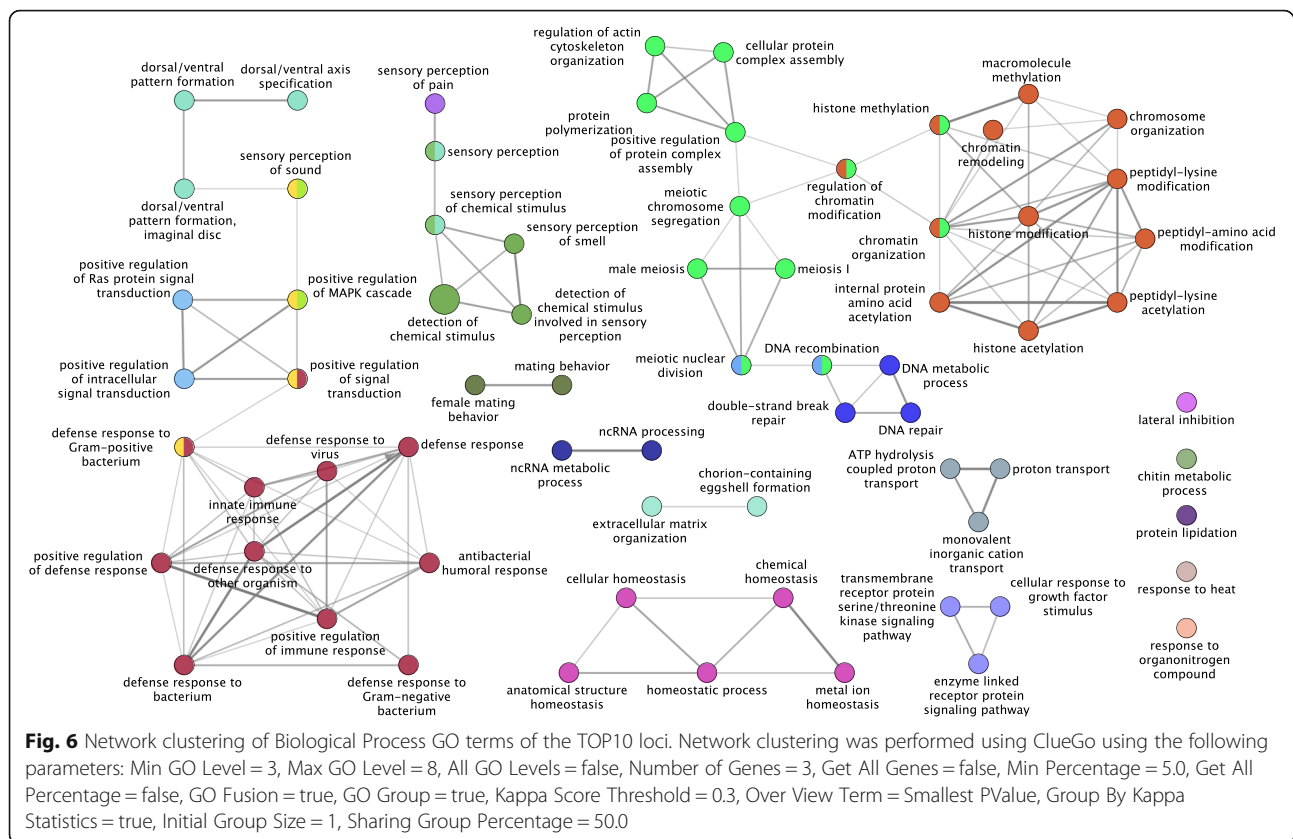
In this study we sequenced, assembled and analyzed the genomes of each of the four cactus host populations of *D. mojavensis* for the purpose of assessing the genomic consequences of the adaptation to local ecological



conditions. Overall, we were able to analyze the sequence, pattern of divergence and structure of 9087 genes. And although the four genomes examined diverged relatively recently [30–34], for several loci, sufficient number of substitutions occurred for us to begin to assess the changes associated with cactus host adaptation.

Unlike what is present in *D. melanogaster*, *D. mojavensis* chromosomes are all acrocentric and its karyotype is composed of six Muller elements [49]. In *D. melanogaster* element A is the X chromosome and elements B/C and D/E form large metacentric chromosomes (2L/2R and 3L/3R, respectively), while the F element or dot chromosome is reduced in sized and highly heterochromatic [50, 51]. In *D. mojavensis* we observed the highest rate of molecular evolution in the small F element, followed by elements B and E, and then the remaining autosomal elements and the X chromosome (Fig. 2).

Selection on the X chromosome has been examined in a number of studies with somewhat variable results [52]. Analysis of several melanogaster group species has shown significant elevated ω values for genes on the X chromosome [17]. From population genetics theory it is generally predicted that the X chromosome would show elevated rates of evolution due to its reduced population size and level of recombination [52]. A subsequent genomic analysis of the X chromosome across more distant *Drosophila* species (*D. melanogaster*, *D. pseudoobscura*, *D. miranda* and *D. yakuba*) failed to find evidence of increased protein evolution on the X chromosome [53]. It is difficult to make any conclusions about the lack of a pattern of accelerated X chromosome evolution found



here, it may be possible that there has not been enough divergence time between these populations for factors such as effective population size to have a measurable effect. The greatest ω values were present in the dot chromosome which in *D. mojavensis* is heterochromatic and has a highly reduced level of recombination [54], which would make it highly susceptible to sweeps and hence higher rates of molecular evolution.

Within *D. mojavensis* there are polymorphic inversions in Muller elements B and E [55], both exhibited overall higher chromosomal-wide levels of ω (Fig. 3). Lower levels of recombination and higher divergence rates have been known to occur around the inversion breakpoint regions in *Drosophila* [56]. One possible explanation for the elevated rates of molecular evolution in these chromosomes is the distinct karyotypes of the sequenced lines (Additional file 2: Table S9). One consequence of a template-based assembly as performed in this study, is that chromosomal structural differences can be largely wiped away. A more detailed analysis of the consequence of chromosomal inversion on the evolutionary trajectories of associated loci will be performed in future analyses of de novo assemblies of *D. mojavensis* genomes from all host populations [57] as well as from sibling species (*D. arizonae* and *D. navojoa*) (unpublished data, Matzkin). Furthermore, these new chromosome-level genome

assemblies of *D. mojavensis* and related species will allow us to determine the fraction of loci with high ω that are de novo and unique to the *D. mojavensis* lineage.

Genes across the genome as well as those with evidence of positive selection or in the top 10% of ω values were assessed for a number of characteristics. Genome-wide loci exhibiting greater ω values tended to be shorter, have fewer exons (3 or less), have low expression, be differentially expressed in response to cactus host use and have fixed expression differences across the four cactus host populations of *D. mojavensis* (Fig. 3; Additional file 2: Figures S7, S12, S15). Overall this pattern of divergence was similar when examining the TOP10 or PAML-FDR loci. Previous genomic analyses in *D. melanogaster* and related species have observed similar characteristics of loci with elevated ω values. This indicates that although the phylogenetic scale of the present study is limited (within *D. mojavensis*) the forces shaping genome evolution between diverged species can also be observed between recently isolated populations within species.

The first comparative genomic study within the *D. melanogaster* group species [58] observed an association between coding length and ω , which they partially attributed to a positive correlation between K_s and protein length. Longer genes have more of these

mutations and this may explain in part why genes with high ω values are likely to be shorter. In this study we did not observe such correlation, in fact the relationship is negative ($P < 0.001$), but explains very little of the variation in K_s ($r^2 = 0.004$) (Additional file 2: Figure S21). Therefore, it is difficult to infer the effect of the association between K_s and protein length, and the lack of positive correlation might be a function of the close relationship between the genomes studied here. The negative association between intron number and rate of molecular evolution has been previously suggested to be due to the presence of exonic splice site enhancers which help in the correct removal of introns from the transcription sequence. As mutations in these regions are more likely to be conserved changes here could cause an intron to not be removed or part of an exon to be removed instead [59]. The link between intron presence and ω values may also help explain why TOP10 genes tend to be shorter as long genes are more likely to have introns [60]. The correlation between gene length and rate of molecular evolution could also be explained as a result of the increased level of interactions between sites of larger exons [61]. In this study a negative correlation between ω and exon length ($r^2 = 0.08$, $P < 0.001$) was observed (Additional file 2: Figure S22). These interactions between residues of a protein, commonly refer to as Hill-Robertson interference [62], have a tendency to buffer against the accumulation of amino acid substitutions and can explain a significant portion of the pattern of molecular evolution in genomes [63].

Highly expressed genes tend to have a higher level of constraint as indicated by the tendency of having lower rates of molecular evolution. This has been previously explained as being a result of selection against mutations that alter transcriptional and translational efficiency as well as selection for the maintenance of correct folding (translational robustness) [58, 64–68]. Given our coarse transcription data we were not able to tease apart which of the above-mentioned forces might more strongly shape the rate of molecular evolution in these genomes. Nonetheless we observed a clear negative relationship across the four *D. mojavensis* genomes between transcriptional level and ω . In addition to overall expression, both tissue and sex-bias expression have been known shape the evolutionary trajectories of genes [63, 69–71]. Male, or more specifically testes expressed genes have been associated with elevated rates of molecular evolution in *Drosophila* and across many taxa [72]. Many of these loci are believed to be under strong sexual selection, which would explain their accelerated rate of molecular evolution. As predicted we observed an overall higher rate of molecular evolution in male-biased genes. Even female-biased loci exhibited a

significant greater ω than unbiased genes. Previous behavioral and molecular studies in *D. mojavensis* have shown that this species experiences strong and recurrent bouts of sexual selection [29, 73–79].

Loci indicating a pattern of positive selection and those with elevated ω appear to be associated with a wide range of metabolic processes. These changes are likely a result of the distinct nutritional and xenobiotic environment the different *D. mojavensis* populations experience. The chemical composition of the cacti and the species of yeast found in each rot varies [35–42] and thus the populations have likely needed to optimize the recognition, avoidance and processing of these necrosis-specific compounds through changes in metabolism, physiology and behavior.

One aspect of metabolism that has likely been shaped by cactus host adaptation is the detoxification of cactus compounds, as the distinct cactus hosts have different chemical compositions. Expression studies have shown that genes involved in detoxification are enriched when flies develop in an alternative necrotic cactus species. Fitness costs of living on the alternative cactus have also been shown to be quite high with those flies having low viability (<40%) [44, 80, 81]. Out of all GO terms examined in this study, the only ones that were consistently overrepresented were those associated with serine-type endopeptidase activity. These type of proteins perform a number of function within organisms, among them is their targeting of organophosphorus toxins [82]. These compounds are often used in pesticides and are found to inhibit serine hydrolase function in both insects and vertebrates [82]. While the apparent positive selection on these genes could be due to a response to pesticides they might experience in the field, but more likely they may be evolving in response to the effects of the toxic or nutritional compounds found in cactus rots.

Cactophilic *Drosophila* have been shown to deploy a number of enzymatic strategies to ameliorate the deleterious consequences of ingesting cactus necrosis-derived compounds. Many of the previously identified proteins playing a role in detoxification in cactophiles (Glutathione S-transferases, Cytochrome P450s, Esterases and UDP-glycosyltransferase) have been associated with detoxification in a broad number of taxa [83–87]. In fact, in recent comparative genomic analysis of the cactophilic *D. buzzatii* [88] and *D. aldrichi* [89], a number of metabolic genes, including those associated with detoxification were shown to be under positive selection. In the present genomic analysis of the *D. mojavensis* genome we observed that the largest functional cluster (Fig. 5) was composed of several genes belonging to known detoxification protein families, such as Cytochrome P450 and Glutathione S-transferases (Gst). Furthermore, previous transcriptional studies have indicated

that these same categories of detoxification loci are differentially expressed when *D. mojavensis* are utilizing necrotic cactus tissues [43, 44]. A population genetics analysis of *GstD1* has indicated a pattern of adaptive amino acid evolution at this locus in the Sonora and Baja California populations [32]. The location of the fixed residue fixed in the lineages leading to these two populations indicated potential functional consequences and a recent kinetic analysis of these proteins have support this prediction (Matzkin, unpublished data).

The diversity of bacterial species found on each necrotic cactus provides, directly or indirectly, nutritional resources for the fly populations, but also are composed of potentially distinct pathogenic organisms [90, 91]. A number of genes with elevated rates of molecular evolution in this study are linked to a range of processes involved with the immune response. As each population is faced with a different composition of threats, the evolutionary arms race between flies and their pathogens creates further divergence between the populations as they face different pathogenic landscapes. Studies in other species, such as *D. simulans*, have found that genes with immune related functions were found to have higher rates of positive selection than the genome average [92]. Exposure to bacterial pathogens in *D. mojavensis* could occur while utilizing the necrotic cactus substrate, but as has been previously suggested [93], via sexual transmission.

A number of the TOP10 loci in this study perform functions associated with sensory perception and behavior (Fig. 6). *Drosophila mojavensis* larvae actively seek out patches of preferred yeast species [94] and across the four host populations there are distinct larval foraging strategies [95]. More specifically genes involved in chemosensory behavior were observed to have elevated ω values in these genomes. Across *Drosophilids*, there have been a number of studies indicating the links between the evolution of chemosensory genes and host specialization [96–98]. In *D. sechellia*, a specialist species, was found to be losing olfactory receptor genes at a faster rate than its sibling generalist species *D. simulans* [99]. In *D. mojavensis* each cactus species rot contains different compounds and thus have distinct set of volatiles emanating from the necroses [40, 41]. These chemical differences have shaped the feeding and oviposition behavior of flies as has been shown by the exposure of adults to cactus volatiles [100–102]. Recent analysis of populations differentiation in odorant and gustatory receptors have shown that unlike what might be initially predicted a number of the changes in these receptors suggests that effects at the level of signal transduction in addition to odorant recognition [103]. Further functional analysis is needed to better understand the evolution and

functional changes of chemosensory pathways associated with the adaptation to necrotic cacti.

In addition to their role in xenobiotic metabolism, serine proteases have been shown to be involved in the network of proteins associated with reproductive interactions in several taxa. In *D. melanogaster* accessory gland proteins (ACP), such as sex peptide, are found to perform a wide range of functions ranging from stimulating ovulation and reducing a female's remating rate to helping to defend against infections [104–106]. Knockouts of serine proteases have been shown to interfere with the behavioral and physiological effects of the male-derived sex peptide [106]. In *D. mojavensis* and its sister species *D. arizonae* a large number of proteases are expressed in female reproductive tracts and several have been shown to be under strong positive selection [76, 107–109]. In addition to ACPs being transferred via the ejaculate, gene transcripts have been found to be deposited by males into females during copulation [75]. Some of these male-derived transcripts could alter the female's transcriptional response, while other may potentially be translated within females. Furthermore, the loci of several of these male-transferred transcripts show a pattern of strong and continuous positive selection, likely as the result of persistent sexual selection [74]. While there seems to be no postzygotic effects of sexual isolation within the *D. mojavensis* populations there is some evidence of prezygotic isolation, where certain populations prefers to mate with members of its own population [29]. The pattern of positive selection and/or elevated rate of molecular evolution for proteases and reproductive loci in the present study may highlight the continuing genomic consequence of sexual selection in this species.

Conclusions

Local ecological adaptation can shape the pattern variation at multiple levels (life history, behavior and physiological), and the imprint of this multifaceted selection can be observed at the genomic level. In this first ever genome-wide analysis of the pattern of molecular evolution across the four ecologically distinct populations of *D. mojavensis*, we have begun to describe the genomic consequences of the adaptation of these cactophilic *Drosophila* to their respective environments. Given that across the four populations are known differences in cactus host use, which encompass differences in both toxic and nutritional compounds, but as well as necrotic host density, temperature, exposure to desiccation and likely pathogens and predators, it was expected that a number of functional classes of loci might be under selection. Among genes with elevated rates of change are those involved in detoxification, metabolism, chemosensory perception, immunity, behavior and reproduction. We observed general patterns of variation across the genomes indicating that loci with elevated rates of molecular evolution tended to be shorter, with fewer exons and have low overall

expression. Furthermore, fast evolving loci also were more likely to be differentially expressed in response to cactus host use and have fixed inter-population expression differences, indicating that both transcriptional and coding sequence changes have been involved in the local ecological adaptation of *D. mojavensis*.

Methods

Drosophila mojavensis lines and sample preparation

Fly lines MJBC 155 collected in La Paz, Baja California in February 2001, MJ 122 collected in Guaymas, Sonora in 1998, and MJANZA 402–8 collected in ANZA-Borrego Park, California in April 2002 were used as the source lines for the sequencing of three *D. mojavensis* populations. These lines were highly inbred to reduce the heterozygosity of their DNA. Summary of the karyotype of each of the lines sequenced as well as the Catalina Island template genome stock (15081–1352.00) can be found in Additional file 2: Table S9. The flies were grown for two generations in banana molasses media [95] supplemented with ampicillin (125 µg/ml) and tetracycline (12.5 µg/ml), to prevent the isolation of bacterial DNA in addition to the flies'. DNA was extracted from homogenized whole male flies using a combination of phenol/chloroform DNA extraction and Qiagen DNeasy spin-columns to achieve the required amount of DNA material. RNase A was used to reduce RNA contamination. Gel electrophoresis was run on each sample to check the quality of the extraction. Any samples with RNA contamination were run through a Qiagen QIAquick PCR Purification Kit spin column to filter contaminants. Extracted DNA was sent to the HudsonAlpha Institute for Biotechnology Genomic Services Lab (Huntsville, Alabama) for sequencing. One hundred base pair paired-end and mate pair sequencing was done on an Illumina HiSeq 2000 with one lane for each.

Genome assembly

Paired-end and mate pair Illumina reads were filtered and trimmed using step one of the A5 Pipeline [110]. This step uses SGA [111] and TagDust [112] with the quality scores from the Illumina FASTQ files to reduce the number of low quality reads. A5 was run on the Dense Memory Cluster of the Alabama Super Computer Center with four processing cores and 64 gigabytes of memory allocated for each run. With the reads cleaned they were assembled to the template genome. The reference genome of the Catalina Island population of *D. mojavensis* was assembled as part of the *Drosophila* 12 Genomes Consortium [17]. Version 1.04 of the reference genome was retrieved from FlyBase version FB2015_02 [113]. From the reference sequence, genome scaffolds [114] containing the protein-coding genes previously mapped to a chromosome, were extracted for use as a

template for the assembly; these scaffolds are detailed in Additional file 2: Table S10. The reference templates as well as the Illumina reads were imported into Geneious 8.1. Assembly was done separately for paired-end and mate pair data. Using Geneious 8.1 and its Map to Reference feature the cleaned reads were assembled to each of the template scaffolds. BAM files were exported for each paired-end and mate pair assembly. SAMtools [115] was used to merge BAM files to create an assembly with both types of reads. This merged BAM file was imported into Geneious 8.1 where consensus sequences were determined for each scaffold using majority calling to limit the number of ambiguities. GTF files for each scaffold used were retrieved from FlyBase version FB2015_02 [113]. These annotations were transferred to each of the new genomes by aligning each assembled genome scaffold to the reference genome scaffold using Mauve Genome Alignment [116] with default settings except for selecting assume collinear genomes. After alignment, annotations were transferred from the reference to the new assembly. The resulting scaffolds were exported in GenBank format. Using the EMBOSS program, extractfeat [117], CDS sequences were extracted from the assembled scaffolds. Sequence files for each gene were concatenated and then aligned using the default settings of the aligner Muscle 3.8.31 [118]. Only the longest transcript for each gene was used as some genes have multiple splice variants.

Molecular evolution analysis

To generate substitution counts for filtering, the software KaKs Calculator 1.2 [46] was used. Files of aligned genes were converted to AXT format using the Perl script parseFastaIntoAXT.pl including in the package. After conversion each gene was run through the software using the NG method [119]. The output files for each loci were concatenated and then imported into JMP 10 for filtering.

Values for ω were calculated using codeml part of the PAML 4.9 package [47]. Aligned genes were converted to PHYLIP format using BioPerl [120]. As PAML requires a phylogenetic tree to be provided for its calculations a neighbor joining tree was constructed in MEGA 5 [121]. This was done by concatenating all exons from each population and then aligning them using Mauve Genome Alignment [116]. The alignment was converted to MEG format using MEGA and a neighbor joining tree was built using the default settings. The tree was exported in newick format for use by PAML. Genes were removed from analysis if they were not divisible by three, these genes were manually screened and if alignment errors appeared to be the cause, these were manually corrected. Screening was done for stop codons within the sequences by translating the DNA sequence

to protein sequence with Transeq, part of the EMBOSS package [117] and any genes with internal stop codons were removed.

Using the BioPython PAML module [122], control files were built for each gene alignment with default values taken except codon frequency was set to F3x4. Site-class models 0, 7, and 8 were used to calculate the ω values [123–125]. Model 0 is a single ratio based omega value for the entire gene. Model 7 is a null model with 10 classes, which does not allow for positive selection while model 8 adds an additional class that allows for positive selection. Both the ω values and log likelihood values were extracted from each output file and the data was organized in Microsoft Excel. If model 8 significantly better fits the data this is evidence of positive selection [47]. Significance values were found by taking the difference between the log likelihood values of the two outputs and multiplying them by two. This value was then compared a chi-square distribution to find *P* values for each gene. Genes with less than five total substitutions as determined by KaKs Calculator [46] were filtered out and not considered. This was done to help deal with the low power of these methods when there are very few changes between the populations. Genes with few changes are more likely to cause the software to either return an undefined result or to reach the maximum ω the software allows. In addition, genes with either no nonsynonymous or no synonymous changes were also removed. This yielded a total of 9087 genes that were used in the analysis. Histograms of a \log_2 transformation of the ω values were produced using JMP 10. A comparison between the \log_2 transformations of the NG Ka/Ks and the omega value from model 0 of codeml was generated with JMP 10.

The length of each gene's coding sequence was extracted from the PHYLIP sequence headers. This was to determine if genes with longer length have significantly different omega values. Genes were binned based on length and an ANOVA with post-hoc Tukey test using JMP 10 was used to compare length bins for significance. Intron data was extracted from the reference genome annotation using Geneious 8.1. Based on this, genes were binned based on the number of exons. ANOVA with post-hoc Tukey test in JMP 10 compared the bin sets for significant difference in omega. To determine if there was a significant difference in omega between genes present on each Muller element ANOVA with post-hoc Tukey test was used in JMP 10 to compare omega value distribution on each element.

Expression analysis

Previous transcriptional studies provided differential expression data for cactus host shifts [44] and between

populations [45]. Loci that were found to be significant with codeml model 7 and 8 were removed from this analysis. The model 0 omega for loci with a FDR significance greater than 0.001 for third-instar larva from the *D. mojavensis* Sonora population that were raised on agria cactus rot was compared to non-significant loci using ANOVA in JMP 10. Comparison of model 0 omega between FDR significant loci and non-significant loci was also done for differential expression between third-instar larva of the four host populations with ANOVA in JMP 10.

To explore the relationship between omega and gene expression level RNAseq data from [48] was retrieved for whole male and female *D. mojavensis* flies as aligned BAM files. Differential expression was calculated by using edgeR [126] to look for genes with significantly higher male or female expression. Box plots of omega model 0 for genes with significant male or female expressed genes as well as genes without sex based expression were compared using ANOVA with post-hoc Tukey test in JMP 10. Average adjusted (+0.25) \log_2 RPKM of non-sex biased genes was plotted against \log_2 omega model 0 and linear regression was performed on the data with JMP 10.

Gene ontology terms analysis

Network graphs were generated using Cytoscape 3.2.1 [127] with the add-on app ClueGO 2.2.5 [128]. GO term and KEGG pathway data used was from the June 2016 release. A custom *D. melanogaster* reference set was used for analysis based on *D. melanogaster* genes with a *D. mojavensis* ortholog that was present in the unfiltered dataset as retrieved from FlyBase version FB2017_06 [112]. Both the TOP10 and PAML-FDR genes were run on, biological processes, molecular function and KEGG terms. Data for GO term summary tables was retrieved from FlyBase version FB2017_06 *D. melanogaster* release 6.19 [113]. For each *D. mojavensis* gene with a *D. melanogaster* ortholog, GO term summaries were phrased from the FlyBase GO Summary Ribbons for molecular function and biological process. Clustering done with JMP 10 using the Ward method and 15 groups allowed.

Supplementary information

Supplementary information accompanies this paper at <https://doi.org/10.1186/s12864-019-6097-z>.

Additional file 1. This file is **Table S1**, which contains the names, descriptive characteristics, test statistics and *D. melanogaster* ortholog information for all the *D. mojavensis* loci examined in the study.

Additional file 2. This file includes the supplementary **Tables S2-S10** and **Figures S1-S22**.

Additional file 3. This file contains supplementary **Table S11** which is a multi-worksheet Excel document containing the gene ontology analysis

for both Biological Process and Molecular Function for the PAML-FDR and TOP10 loci.

Abbreviations

2 L: Left arm of 2nd chromosome in *D. melanogaster*; 2R: Right arm of 2nd chromosome in *D. melanogaster*; 3 L: Left arm of 3rd chromosome in *D. melanogaster*; 3R: Right arm of 3rd chromosome in *D. melanogaster*; ACP: Accessory gland protein; ANOVA: Analysis of Variance; BAM: Binary Alignment Map; CDS: Coding sequence; EMBOSS: European Molecular Biology Open Software Suite; FDR: False Discovery Rate; GO: Gene Ontology; Gst: Glutathione S-transferase; Ka: Number of nonsynonymous substitution per nonsynonymous site; kb: Kilobase; KEGG: Kyoto Encyclopedia of Genes and Genomes; Ks: Number of synonymous substitution per synonymous site; MEGA: Molecular Evolutionary Genetics Analysis software; PAML: Phylogenetic Analysis of Maximum Likelihood program; PAML-FDR: PAML significant loci post-FDR correction; PHYLIP: Phylogeny Inference Package; RPKM: Reads Per Kilobase per Million mapped reads; TOP10: Loci with ω values in the top 10% of the distribution

Acknowledgements

The authors greatly acknowledge the work of Laurel Brandsmeier in this project.

Authors' contributions

CWA performed the assembly and analysis of the genomic data and was involved in the writing of the manuscript. LMM conceived of and designed the study, was involved in the analysis and the writing of the manuscript. All authors read and approved the final manuscript.

Funding

This work was supported by a Junior Faculty Distinguished Research award from the University of Alabama in Huntsville and partly supported by grants from the National Science Foundation (DEB-1219387 and IOS-1557697) to LMM.

Availability of data and materials

The sequence reads datasets supporting the conclusions of this article are available in the NCBI Sequence Read Archive (SRA) under the study accession number SRP190536 and the alignment files for each loci can be accessed via OSF (<https://osf.io/2759a>). Additionally, datasets supporting the conclusions of this article are included within the article's supplementary information files.

Ethics approval and consent to participate

Not applicable

Consent for publication

Not applicable

Competing interests

The authors declare that they have no competing interests.

Author details

¹Department of Biological Sciences, University of Alabama in Huntsville, 301 Sparkman Drive, Huntsville, AL 35899, USA. ²Department of Entomology, University of Arizona, 1140 E. South Campus Drive, Tucson, AZ 85721, USA. ³BIO5 Institute, University of Arizona, 1657 East Helen Street, Tucson, AZ 85721, USA. ⁴Department of Ecology and Evolutionary Biology, University of Arizona, 1041 E. Lowell St., Tucson, AZ 85721, USA.

Received: 11 April 2019 Accepted: 11 September 2019

Published online: 12 October 2019

References

- Feder ME, Mitchell-Olds T. Evolutionary and ecological functional genomics. *Nat Rev Genet.* 2003;4:649–55.
- Stapley J, Reger J, Feulner PGD, Smadja C, Galindo J, Ekblom R, Bennison C, Ball AD, Beckerman AP, Slate J. Adaptation genomics: the next generation. *Trends Ecol Evol.* 2010;25(12):705–12.
- Barrett RDH, Hoekstra HE. Molecular spandrels: tests of adaptation at the genetic level. *Nat Rev Genet.* 2011;12(11):767–80.
- Storz JF, Wheat CW. Integrating evolutionary and functional approaches to infer adaptation at specific loci. *Evolution.* 2010;64(9):2489–509.
- Ungerer MC, Johnson LC, Herman MA. Ecological genomics: understanding gene and genome function in the natural environment. *Heredity.* 2008;100(2):178–83.
- Nosil P. *Ecological Speciation.* Oxford: Oxford University Press; 2012.
- Rundle HD, Nosil P. Ecological speciation. *Ecol Lett.* 2005;8(3):336–52.
- Funk DJ. Isolating a role for natural selection in speciation: host adaptation and sexual isolation in *Neochlamisus bebbianae* leaf beetles. *Evolution.* 1998;52(6):1744–59.
- Wu CI, Ting CT. Genes and speciation. *Nat Rev Genet.* 2004;5(2):114–22.
- Feder JL, Opp SB, Wlazlo B, Reynolds K, Go W, Spisak S. Host Fidelity is an effective premating barrier between sympatric races of the apple maggot Fly. *Proc Natl Acad Sci.* 1994;91(17):7990–4.
- Funk DJ, Egan SP, Nosil P. Isolation by adaptation in *Neochlamisus* leaf beetles: host-related selection promotes neutral genomic divergence. *Mol Ecol.* 2011;20(22):4671–82.
- Egan SP, Janson EM, Brown CG, Funk DJ. Postmating isolation and genetically variable host use in ecologically divergent host forms of *Neochlamisus bebbianae* leaf beetles. *J Evol Biol.* 2011;24(10):2217–29.
- Nielsen R, Bustamante C, Clark AG, Glanowski S, Sackton TB, Hubisz MJ, Fedel-Alon A, Tanenbaum DM, Civello D, White TJ, et al. A scan for positively selected genes in the genomes of humans and chimpanzees. *PLoS Biol.* 2005;3(6):976–85.
- Clark AG, Glanowski S, Nielsen R, Thomas PD, Kejarawal A, Todd MA, Tanenbaum DM, Civello D, Lu F, Murphy B, et al. Inferring nonneutral evolution from human-chimp-mouse orthologous gene trios. *Science.* 2003;302(5652):1960–3.
- Bustamante CD, Fedel-Alon A, Williamson S, Nielsen R, Hubisz MT, Glanowski S, Tanenbaum DM, White TJ, Sninsky JJ, Hernandez RD, et al. Natural selection on protein-coding genes in the human genome. *Nature.* 2005;437(7062):1153–7.
- Kosiol C, Vinar T, da Fonseca RR, Hubisz MJ, Bustamante CD, Nielsen R, Siepel A. Patterns of positive selection in six mammalian genomes. *PLoS Genet.* 2008;4(8):e1000144.
- Consortium DG. Evolution of genes and genomes on the *Drosophila* phylogeny. *Nature.* 2007;450(7167):203–18.
- Yang Z. The power of phylogenetic comparison in revealing protein function. *Proc Natl Acad Sci.* 2005;102(9):3179–80.
- Pool JE, Corbett-Detig RB, Sugino RP, Stevens KA, Cardeno CM, Crepeau MW, Duchon P, Emerson JJ, Saelao P, Begun DJ, et al. Population genomics of sub-saharan *Drosophila melanogaster*: African diversity and non-African admixture. *PLoS Genet.* 2012;8(12):e1003080.
- Langley CH, Stevens K, Cardeno C, Lee YC, Schrider DR, Pool JE, Langley SA, Suarez C, Corbett-Detig RB, Kolaczowski B, et al. Genomic variation in natural populations of *Drosophila melanogaster*. *Genetics.* 2012;192(2):533–98.
- Bergland AO, Tobler R, Gonzalez J, Schmidt P, Petrov D. Secondary contact and local adaptation contribute to genome-wide patterns of clinal variation in *Drosophila melanogaster*. *Mol Ecol.* 2016;25(5):1157–74.
- Campo D, Lehmann K, Fjeldsted C, Souaiaia T, Kao J, Nuzhdin SV. Whole-genome sequencing of two north American *Drosophila melanogaster* populations reveals genetic differentiation and positive selection. *Mol Ecol.* 2013;22(20):5084–97.
- Grenier JK, Arguello JR, Moreira MC, Gottipati S, Mohammed J, Hackett SR, Boughton R, Greenberg AJ, Clark AG. Global diversity lines-a five-continent reference panel of sequenced *Drosophila melanogaster* strains. *G3-Genes Genom Genet.* 2015;5(4):593–603.
- Pool JE. The mosaic ancestry of the *Drosophila* genetic reference panel and the *D. melanogaster* reference genome reveals a network of epistatic fitness interactions. *Mol Biol Evol.* 2015;32(12):3236–51.
- Shiao MS, Chang JM, Fan WL, Lu MY, Notredame C, Fang S, Kondo R, Li WH. Expression divergence of chemosensory genes between *Drosophila sechellia* and its sibling species and its implications for host shift. *Genome Biol Evol.* 2015;7(10):2843–58.
- Chiu JC, Jiang XT, Zhao L, Hamm CA, Cridland JM, Saelao P, Hamby KA, Lee EK, Kwok RS, Zhang GJ, et al. Genome of *Drosophila suzukii*, the spotted wing *Drosophila*. *G3-Genes Genom Genet.* 2013;3(12):2257–71.
- Matzkin LM. Ecological genomics of host shifts in *Drosophila mojavensis*. *Adv Exp Med Biol.* 2014;781(781):233–47.

28. Heed WB. Ecology and genetics of Sonoran desert *Drosophila*. In: Brussard PF, editor. *Ecological genetics: the interface*. New York: Springer-Verlag; 1978. p. 109–26.
29. Markow TA. Sexual isolation among populations of *Drosophila mojavensis*. *Evolution*. 1991;45:1525–9.
30. Reed LK, Nyboer M, Markow TA. Evolutionary relationships of *Drosophila mojavensis* geographic host races and their sister species *Drosophila arizonae*. *Mol Ecol*. 2007;16(5):1007–22.
31. Matzkin LM, Eanes WF. Sequence variation of alcohol dehydrogenase (*Adh*) paralogs in cactophilic *Drosophila*. *Genetics*. 2003;163:181–94.
32. Matzkin LM. The molecular basis of host adaptation in Cactophilic *Drosophila*: molecular evolution of a glutathione S-transferase gene (*GstD1*) in *Drosophila mojavensis*. *Genetics*. 2008;178(2):1073–83.
33. Matzkin LM. Population genetics and geographic variation of alcohol dehydrogenase (*Adh*) paralogs and glucose-6-phosphate dehydrogenase (*G6pd*) in *Drosophila mojavensis*. *Mol Biol Evol*. 2004;21(2):276–85.
34. Smith G, Lohse K, Etges WJ, Ritchie MG. Model-based comparisons of phylogeographic scenarios resolve the intraspecific divergence of cactophilic *Drosophila mojavensis*. *Mol Ecol*. 2012;21(13):3293–307.
35. Starmer WT. Analysis of the community structure of yeasts associated with the decaying stems of Cactus. I. *Stenocereus gummosus*. *Microb Ecol*. 1982;8(1):71–81.
36. Starmer WT. Associations and Interactions Among Yeasts, *Drosophila* and their habitats. In: Barker JSF, Starmer WT, editors. *Ecological genetics and evolution: the cactus-yeast-Drosophila model system*. New York: Academic Press; 1982. p. 159–74.
37. Fogleman JC, Starmer WT. Analysis of the community structure of yeasts associated with the decaying stems of cactus. III. *Stenocereus thurberi*. *Microb Ecol*. 1985;11(2):165–73.
38. Starmer WT, Lachance MA, Phaff HJ, Heed WB. The biogeography of yeasts associated with decaying cactus tissue in North America, the Caribbean, and Northern Venezuela. *Evol Biol*. 1990;24:253–96.
39. Fellows DF, Heed WB. Factors affecting host plant selection in desert-adapted cactophilic *Drosophila*. *Ecology*. 1972;53:850–8.
40. Kircher HW. Chemical composition of cacti and its relationship to Sonoran Desert *Drosophila*. In: Barker JSF, Starmer WT, editors. *Ecological genetics and evolution: the cactus-yeast-Drosophila model system*. New York: Academic Press; 1982. p. 143–58.
41. Fogleman JC, Abril JR. Ecological and evolutionary importance of host plant chemistry. In: Barker JSF, Starmer WT, MacIntyre RJ, editors. *Ecological and evolutionary genetics of Drosophila*. New York: Plenum Press; 1990. p. 121–43.
42. Fogleman JC, Danielson PB. Chemical interactions in the cactus-microorganism-*Drosophila* model system of the Sonoran Desert. *Am Zool*. 2001;41(4):877–89.
43. Matzkin LM, Watts TD, Bitler BG, Machado CA, Markow TA. Functional genomics of cactus host shifts in *Drosophila mojavensis*. *Mol Ecol*. 2006;15:4635–43.
44. Matzkin LM. Population transcriptomics of cactus host shifts in *Drosophila mojavensis*. *Mol Ecol*. 2012;21(10):2428–39.
45. Matzkin LM, Markow TA. Transcriptional differentiation across the four cactus host races of *Drosophila mojavensis*. In: Michalak P, editor. *Speciation: natural processes, genetics and biodiversity*. Hauppauge: Nova Science Publishers Inc.; 2013. p. 119–36.
46. Zhang Z, Li J, Zhao XQ, Wang J, Wong GK, Yu J. KaKs_Calculator: calculating Ka and Ks through model selection and model averaging. *Genomics Proteomics Bioinformatics*. 2006;4(4):259–63.
47. Yang ZH. PAML 4: Phylogenetic analysis by maximum likelihood. *Mol Biol Evol*. 2007;24(8):1586–91.
48. Graveley BR, Brooks AN, Carlson J, Duff MO, Landolin JM, Yang L, Artieri CG, van Baren MJ, Boley N, Booth BW, et al. The developmental transcriptome of *Drosophila melanogaster*. *Nature*. 2011;471(7339):473–9.
49. Wasserman M. Cytological and phylogenetic relationships in the Repleta Group of the Genus *Drosophila*. *Proc Natl Acad Sci*. 1960;46(6):842–59.
50. Riddle NC, Elgin SCR. The *Drosophila* dot chromosome: where genes flourish amidst repeats. *Genetics*. 2018;210(3):757–72.
51. Bridges CB. Salivary chromosome maps with a key to the banding of the chromosomes of *Drosophila melanogaster*. *J Hered*. 1935;26(2):60–4.
52. Singh ND, Petrov DA. Evolution of gene function on the X chromosome versus the autosomes. *Gene Protein Evolution*. 2007;3:101–18.
53. Thornton K, Bachtrog D, Andolfatto P. X chromosomes and autosomes evolve at similar rates in *Drosophila*: no evidence for faster-X protein evolution. *Genome Res*. 2006;16(4):498–504.
54. Leung W, Shaffer CD, Reed LK, Smith ST, Barshop W, Dirkes W, Dothager M, Lee P, Wong J, Xiong D, et al. *Drosophila* Muller F elements maintain a distinct set of genomic properties over 40 million years of evolution. *G3*. 2015;5(5):719–40.
55. Ruiz A, Heed WB, Wasserman M. Evolution of the Mojavensis cluster of cactophilic *Drosophila* with descriptions of two new species. *J Hered*. 1990;81:30–42.
56. Hasson E, Eanes WF. Contrasting histories of three gene regions associated with in (3L) Payne of *Drosophila melanogaster*. *Genetics*. 1996;144(4):1565–75.
57. Jaworski CC, Allan CW, Matzkin LM. Chromosome-level hybrid *de novo* genome assemblies as an attainable option for non-model organisms. *bioRxiv*. 2019. <https://doi.org/10.1101/748228>.
58. Larracuente AM, Sackton TB, Greenberg AJ, Wong A, Singh ND, Sturgill D, Zhang Y, Oliver B, Clark AG. Evolution of protein-coding genes in *Drosophila*. *Trends Genet*. 2008;24(3):114–23.
59. Blencowe BJ. Exonic splicing enhancers: mechanism of action, diversity and role in human genetic diseases. *Trends Biochem Sci*. 2000;25(3):106–10.
60. Hawkins JD. A survey on intron and exon lengths. *Nucleic Acids Res*. 1988; 16(21):9893–908.
61. Comeron JM, Guthrie TB. Intragenic Hill-Robertson interference influences selection intensity on synonymous mutations in *Drosophila*. *Mol Biol Evol*. 2005;22(12):2519–30.
62. Hill WG, Robertson A. Effect of linkage on limits to artificial selection. *Genet Res*. 1966;8(3):269–94.
63. Fraïsse C, Puixeu Sala G, Vicoso B. Pleiotropy modulates the efficacy of selection in *Drosophila melanogaster*. *Mol Biol Evol*. 2018;36(3):500–15.
64. Drummond DA, Bloom JD, Adami C, Wilke CO, Arnold FH. Why highly expressed proteins evolve slowly. *Proc Natl Acad Sci*. 2005;102(40):14338–43.
65. Wilke CO, Drummond DA. Population genetics of translational robustness. *Genetics*. 2006;173(1):473–81.
66. Pal C, Papp B, Hurst LD. Highly expressed genes in yeast evolve slowly. *Genetics*. 2001;158(2):927–31.
67. Akashi H. Gene expression and molecular evolution. *Curr Opin Genet Dev*. 2001;11(6):660–6.
68. Nuzhdin S, Wayne M, Harmon K, McIntyre L. Common pattern of evolution of gene expression level and protein sequence in *Drosophila*. *Mol Biol Evol*. 2004;21(7):1308–17.
69. Zhang Z, Hambuch TM, Parsch J. Molecular evolution of sex-biased genes in *Drosophila*. *Mol Biol Evol*. 2004;21(11):2130–9.
70. Grath S, Parsch J. Sex-biased gene expression. *Annu Rev Genet*. 2016;50:29–44.
71. Meisel RP. Towards a more nuanced understanding of the relationship between sex-biased gene expression and rates of protein-coding sequence evolution. *Mol Biol Evol*. 2011;28(6):1893–900.
72. Swanson WJ, Vacquier VD. The rapid evolution of reproductive proteins. *Nat Rev Genet*. 2002;3(2):137–44.
73. Bono JM, Markow TA. Post-zygotic isolation in cactophilic *Drosophila*: larval viability and adult life-history traits of *D. mojavensis/D. arizonae* hybrids. *J Evol Biol*. 2009;22(7):1387–95.
74. Bono JM, Matzkin LM, Hoang K, Brandsmeier L. Molecular evolution of candidate genes involved in post-mating-prezygotic reproductive isolation. *J Evol Biol*. 2015;28(2):403–14.
75. Bono JM, Matzkin LM, Kelleher ES, Markow TA. Postmating transcriptional changes in reproductive tracts of con- and heterospecifically mated *Drosophila mojavensis* females. *Proc Natl Acad Sci*. 2011;108(19):7878–83.
76. Kelleher ES, Markow TA. Reproductive tract interactions contribute to isolation in *Drosophila*. *Fly*. 2007;1(1):33–7.
77. Knowles LL, Markow TA. Sexually antagonistic coevolution of a postmating-prezygotic reproductive character in desert *Drosophila*. *Proc Natl Acad Sci*. 2001;98(15):8692–6.
78. Krebs RA, Markow TA. Courtship behavior and control of reproductive isolation in *Drosophila mojavensis*. *Evolution*. 1989;43:908–13.
79. Pitnick S, Miller GT, Schneider K, Markow TA. Ejaculate-female coevolution in *Drosophila mojavensis*. *P Roy Soc B-Biol Sci*. 2003;270(1523):1507–12.
80. Etges WJ, Heed WB. Sensitivity to larval density in populations of *Drosophila mojavensis*: influences of host plant variation on components fitness. *Oecologia*. 1987;71:375–81.
81. Etges WJ. Direction of life history evolution in *Drosophila mojavensis*. In: Barker JSF, Starmer WT, MacIntyre RJ, editors. *Ecological and evolutionary genetics of Drosophila*. New York: Plenum Press; 1990. p. 37–56.
82. Casida JE, Quistad GB. Serine hydrolase targets of organophosphorus toxicants. *Chem Biol Interact*. 2005;157:277–83.

83. Luque T, O'Reilly DR. Functional and phylogenetic analyses of a putative *Drosophila melanogaster* UDP-glycosyltransferase gene. *Insect Biochem Mol Biol*. 2002;32(12):1597–604.
84. Ranson H, Rossiter L, Ortelli F, Jensen B, Wang XL, Roth CW, Collins FH, Hemingway J. Identification of a novel class of insect glutathione S-transferases involved in resistance to DDT in the malaria vector *Anopheles gambiae*. *Biochem J*. 2001;359:295–304.
85. Ranson H, Hemingway J. Glutathione transferases. In: Gilbert LI, Iatrou K, Gill SS, editors. *Comprehensive Molecular Insect Science*, vol. 5. Amsterdam: Elsevier; 2005. p. 383–402.
86. Feyereisen R. Insect cytochrome P450. In: Gilbert LI, Iatrou K, Gill SS, editors. *Comprehensive Molecular Insect Science*, vol. 4. Amsterdam: Elsevier; 2005. p. 1–77.
87. Li XC, Schuler MA, Berenbaum MR. Molecular mechanisms of metabolic resistance to synthetic and natural xenobiotics. *Annu Rev Entomol*. 2007;52:231–53.
88. Guillen Y, Rius N, Delprat A, Williford A, Muyas F, Puig M, Casillas S, Ramia M, Egea R, Negre B, et al. Genomics of ecological adaptation in Cactophilic *Drosophila*. *Genome Biol Evol*. 2015;7(1):349–66.
89. Rane RV, Pearce SL, Li F, Coppin C, Schiffer M, Shirriffs J, Sgro CM, Griffin PC, Zhang G, Lee SF, et al. Genomic changes associated with adaptation to arid environments in cactophilic *Drosophila* species. *BMC Genomics*. 2019;20(1):52.
90. Foster JLM, Fogleman JC. Identification and ecology of bacterial communities associated with Necroses of 3 Cactus species. *Appl Environ Microb*. 1993;59(1):1–6.
91. Foster J, Fogleman J. Bacterial succession in necrotic tissue of *Agria cactus* (*Stenocereu gummosus*). *Appl Environ Microb*. 1994;60(2):619–25.
92. Schlenke T, Begun D. Natural selection drives *Drosophila* immune system evolution. *Genetics*. 2003;164(4):1471–80.
93. Markow TA. Assortative fertilization in *Drosophila*. *Proc Natl Acad Sci*. 1997;94(15):7756–60.
94. Fogleman JC, Stamer WT, Heed WB. Larval selectivity for yeast species by *Drosophila mojavensis* in natural substrates. *Proc Natl Acad Sci*. 1981;78(7):4435–9.
95. Coleman JM, Benowitz KM, Jost AG, Matzkin LM. Behavioral evolution accompanying host shifts in cactophilic *Drosophila* larvae. *Ecol Evol*. 2018;8(14):6921–31.
96. Vosshall LB, Stocker RE. Molecular architecture of smell and taste in *Drosophila*. *Annu Rev Neurosci*. 2007;30:505–33.
97. McBride CS, Arguello JR. Five *Drosophila* genomes reveal nonneutral evolution and the signature of host specialization in the chemoreceptor superfamily. *Genetics*. 2007;177(3):1395–416.
98. Arguello JR, Cardoso-Moreira M, Grenier JK, Gottipati S, Clark AG, Benton R. Extensive local adaptation within the chemosensory system following *Drosophila melanogaster*'s global expansion. *Nat Commun*. 2016;7. <https://doi.org/10.1038/ncomms11855>.
99. McBride CS. Rapid evolution of smell and taste receptor genes during host specialization in *Drosophila sechellia*. *Proc Natl Acad Sci*. 2007;104(12):4996–5001.
100. Newby BD, Etges WJ. Host preference among populations of *Drosophila mojavensis* (Diptera: Drosophilidae) that use different host cacti. *J Insect Behav*. 1998;11(5):691–712.
101. Date P, Dweck HKM, Stensmyr MC, Shann J, Hansson BS, Rollmann SM. Divergence in olfactory host plant preference in *D. mojavensis* in response to Cactus host use. *PLoS One*. 2013;8(7):e70027.
102. Date P, Crowley-Gall A, Diefendorf AF, Rollmann SM. Population differences in host plant preference and the importance of yeast and plant substrate to volatile composition. *Ecol Evol*. 2017;7(11):3815–25.
103. Diaz F, Allan CW, Matzkin LM. Positive selection at sites of chemosensory genes is associated with the recent divergence and local ecological adaptation in cactophilic *Drosophila*. *BMC Evol Biol*. 2018;18:144.
104. Wolfner MF. The gifts that keep on giving: physiological functions and evolutionary dynamics of male seminal proteins in *Drosophila*. *Heredity*. 2002;88:85–93.
105. Avila FW, Sirot LK, LaFlamme BA, Rubinstein CD, Wolfner MF. Insect seminal fluid proteins: identification and function. *Annu Rev Entomol*. 2011;56:21–40.
106. Findlay GD, Sitnik JL, Wang W, Aquadro CF, Clark NL, Wolfner MF. Evolutionary rate covariation identifies new members of a protein network required for *Drosophila melanogaster* female post-mating responses. *PLoS Genet*. 2014;10(1):e1004108.
107. Kelleher ES, Pennington JE. Protease gene duplication and proteolytic activity in *Drosophila* female reproductive tracts. *Mol Biol Evol*. 2009;26(9):2125–34.
108. Kelleher ES, Swanson WJ, Markow TA. Gene duplication and adaptive evolution of digestive proteases in *Drosophila arizonae* female reproductive tracts. *PLoS Genet*. 2007;3(8):1541–9.
109. Kelleher ES, Watts TD, LaFlamme BA, Haynes PA, Markow TA. Proteomic analysis of *Drosophila mojavensis* male accessory glands suggests novel classes of seminal fluid proteins. *Insect Biochem Mol Biol*. 2009;39(5–6):366–71.
110. Tritt A, Eisen JA, Facciotti MT, Darling AE. An integrated pipeline for de novo assembly of microbial genomes. *PLoS One*. 2012;7(9):e42304.
111. Simpson JT, Durbin R. Efficient de novo assembly of large genomes using compressed data structures. *Genome Res*. 2012;22(3):549–56.
112. Lassmann T, Hayashizaki Y, Daub CO. TagDust-A program to eliminate artifacts from next generation sequencing data. *Bioinformatics*. 2009;25(21):2839–40.
113. Gramates LS, Marygold SJ, dos Santos G, Urbano JM, Antonazzo G, Matthews BB, Rey AJ, Tabone CJ, Crosby MA, Emmert DB, et al. FlyBase at 25: looking to the future. *Nucleic Acids Res*. 2017;45(D1):D663–71.
114. Schaeffer SW, Bhutkar A, McAllister BF, Matsuda M, Matzkin LM, O'Grady PM, Rohde C, Valente VLS, Aguade M, Anderson WW, et al. Polytene chromosomal maps of 11 *Drosophila* species: the order of genomic scaffolds inferred from genetic and physical maps. *Genetics*. 2008;179(3):1601–55.
115. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R, Proc GPD. The sequence alignment/map format and SAMtools. *Bioinformatics*. 2009;25(16):2078–9.
116. Darling AC, Mau B, Blattner FR, Perna NT. Mauve: multiple alignment of conserved genomic sequence with rearrangements. *Genome Res*. 2004;14(7):1394–403.
117. Rice P, Longden I, Bleasby A. EMBOS: the European molecular biology open software suite. *Trends Genet*. 2000;16(6):276–7.
118. Edgar RC. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Res*. 2004;32(5):1792–7.
119. Nei M, Gojobori T. Simple methods for estimating the numbers of synonymous and nonsynonymous nucleotide substitutions. *Mol Biol Evol*. 1986;3(5):418–26.
120. Stajich JE, Block D, Boulez K, Brenner SE, Chervitz SA, Dagdigian C, Fuellen G, Gilbert JGR, Korf I, Lapp H, et al. The bioperl toolkit: Perl modules for the life sciences. *Genome Res*. 2002;12(10):1611–8.
121. Tamura K, Peterson D, Peterson N, Stecher G, Nei M, Kumar S. MEGA5: Molecular evolutionary genetics analysis using maximum likelihood, evolutionary distance, and maximum parsimony methods. *Mol Biol Evol*. 2011;28(10):2731–9.
122. Talevich E, Invergo BM, Cock PJA, Chapman BA. BioPhylo: a unified toolkit for processing, analyzing and visualizing phylogenetic trees in Biopython. *BMC Bioinformatics*. 2012;13:209.
123. Nielsen R, Yang ZH. Likelihood models for detecting positively selected amino acid sites and applications to the HIV-1 envelope gene. *Genetics*. 1998;148(3):929–36.
124. Goldman N, Yang ZH. Codon-based model of nucleotide substitution for protein-coding DNA-sequences. *Mol Biol Evol*. 1994;11(5):725–36.
125. Yang ZH, Nielsen R, Goldman N, Pedersen AMK. Codon-substitution models for heterogeneous selection pressure at amino acid sites. *Genetics*. 2000;155(1):431–49.
126. Robinson MD, McCarthy DJ, Smyth GK. edgeR: a Bioconductor package for differential expression analysis of digital gene expression data. *Bioinformatics*. 2010;26(1):139–40.
127. Shannon P, Markiel A, Ozier O, Baliga NS, Wang JT, Ramage D, Amin N, Schwikowski B, Ideker T. Cytoscape: A software environment for integrated models of biomolecular interaction networks. *Genome Res*. 2003;13(11):2498–504.
128. Bindea G, Mlecnik B, Hackl H, Charoentong P, Tosolini M, Kirilovsky A, Fridman WH, Pages F, Trajanoski Z, Galon J. ClueGO: a Cytoscape plug-in to decipher functionally grouped gene ontology and pathway annotation networks. *Bioinformatics*. 2009;25(8):1091–3.

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.