# Estimating Divergence Dates and Substitution Rates in the *Drosophila* Phylogeny

Darren J. Obbard,*,[1] John Maclennan,[2] Kang-Wook Kim,[3] Andrew Rambaut,[1] Patrick M. O'Grady,[4] and Francis M. Jiggins[5]

[1]Institute of Evolutionary Biology, and Centre for Infection Immunity and Evolution, University of Edinburgh, Edinburgh, United Kingdom

[2]Department of Earth Sciences, University of Cambridge, Cambridge, United Kingdom

[3]Department of Animal and Plant Sciences, University of Sheffield, Sheffield, United Kingdom

[4]Department of Environmental Science, Policy and Management, University of California, Berkeley

[5]Department of Genetics, University of Cambridge, Cambridge, United Kingdom

*Corresponding author: E-mail: darren.obbard@ed.ac.uk.

Associate editor: Jody Hey

Research article

## Abstract

An absolute timescale for evolution is essential if we are to associate evolutionary phenomena, such as adaptation or speciation, with potential causes, such as geological activity or climatic change. Timescales in most phylogenetic studies use geologically dated fossils or phylogeographic events as calibration points, but more recently, it has also become possible to use experimentally derived estimates of the mutation rate as a proxy for substitution rates. The large radiation of drosophilid taxa endemic to the Hawaiian islands has provided multiple calibration points for the *Drosophila* phylogeny, thanks to the "conveyor belt" process by which this archipelago forms and is colonized by species. However, published date estimates for key nodes in the *Drosophila* phylogeny vary widely, and many are based on simplistic models of colonization and coalescence or on estimates of island age that are not current. In this study, we use new sequence data from seven species of Hawaiian *Drosophila* to examine a range of explicit coalescent models and estimate substitution rates. We use these rates, along with a published experimentally determined mutation rate, to date key events in drosophilid evolution. Surprisingly, our estimate for the date for the most recent common ancestor of the genus *Drosophila* based on mutation rate (25–40 Ma) is closer to being compatible with independent fossil-derived dates (20–50 Ma) than are most of the Hawaiian-calibration models and also has smaller uncertainty. We find that Hawaiian-calibrated dates are extremely sensitive to model choice and give rise to point estimates that range between 26 and 192 Ma, depending on the details of the model. Potential problems with the Hawaiian calibration may arise from systematic variation in the molecular clock due to the long generation time of Hawaiian *Drosophila* compared with other *Drosophila* and/or uncertainty in linking island formation dates with colonization dates. As either source of error will bias estimates of divergence time, we suggest mutation rate estimates be used until better models are available.

Key words: *Drosophila*, phylogeny, relaxed clock, Hawai'i, substitution rate, mutation rate.

## Introduction

> "It is to be stressed that all of these estimates are, at best, guesses." (Ashburner et al. 1984)

Generating reliable substitution rate estimates to place an absolute timescale on phylogenetic trees has been an area of interest for ∼50 years (Zuckerkandl and Pauling 1965). Substitution rates are typically estimated either by associating nodes of a phylogenetic tree with dated fossils or phylogeographic events (discussed in Drummond et al. 2006; Ho and Phillips 2009) or by using mutation rate estimates derived from pedigrees or laboratory studies (e.g., Cutter 2008). Assuming neutral evolution, these different approaches should yield similar estimates (reviewed in Bromham and Penny 2003). However, this is often not the case, and long-term substitution rates are sometimes much lower than short-term mutation rates (Ho et al. 2011). This may

be caused by purifying selection, which reduces substitution rates relative to mutation rates, or by sites with high mutation rates becoming saturated. Regardless of the causes, it has led to the view that the calibration points used to estimate substitution rates should ideally be of a similar age to the events being dated (Ho and Larson 2006).

Although the phylogenetic relationships within the Drosophilidae have been extensively studied (e.g., Throckmorton 1975; Ashburner et al. 1984; Grimaldi 1990; DeSalle and Grimaldi 1991; Pelandakis and Solignac 1993; Remsen and DeSalle 1998; Remsen and O'Grady 2002; Schawaroch 2002; Kopp 2006; Da Lage et al. 2007; Gao et al. 2007; O'Grady and DeSalle 2008; Robe et al. 2010; van der Linde et al. 2010; Cao et al. 2011; Gao et al. 2011; O'Grady et al. 2011; Kellermann et al. 2012; Oliveira et al. 2012; Yang et al. 2012), the absolute timescale of *Drosophila* evolution remains surprisingly uncertain. For example, published

**Open Access**

estimates for the most recent common ancestor (MRCA) of *Drosophila melanogaster* and *D. simulans* have ranged from as much as 9.4 Ma (Easteal and Oakeshott 1985) to as little as 1.2 Ma (Cutter 2008). Calibration points for Drosophilidae come from two sources, amber fossils from Baltic and Dominican deposits, and biogeographic dates. Drosophilid fossils are relatively rare, and the majority are from specimens that can be difficult to date (Grimaldi 1987). However, fossil-based estimates suggest that the family Drosophilidae originated at least 30–50 Ma (Throckmorton 1975; Grimaldi 1987) and that the genus *Drosophila* originated at least 20 Ma (Grimaldi 1987, 1988). However, more often, timescales for *Drosophila* have been based on Hawaiian phylogeography (e.g., Carson 1976; Easteal and Oakeshott 1985; Thomas and Hunt 1993; Russo et al. 1995; Tamura et al. 2004). This is possible because the Hawaiian archipelago consists of a line of volcanic islands of decreasing age: each time an island is formed, it is colonized by flies from the neighboring island, and the resulting population then diverges to form a new species. Although there are important caveats (Heads 2011), this "conveyor belt speciation" means species formation can be related to the geological age of the islands (for an introduction, see Fleischer et al. 1998; Price and Clague 2002). This has provided some widely used estimates, including two of the most often quoted dates for the divergence of *D. melanogaster* and *D. simulans* (2.3 Ma from Russo et al. 1995 and 5.4 Ma from Tamura et al. 2004).

There are two potentially important sources of error in the previous substitution rates estimated using Hawaiian *Drosophila*. First, previous studies have not modeled the process of island colonization and gene coalescence. This is important, as polymorphisms in the ancestral population will cause the divergence times between the copies of a gene in two species to exceed the species divergence times, and gene trees may not reflect species trees (Nei 1971; Peterson and Masel 2009; Charlesworth 2010). Second, recent geological and geochemical advances have improved understanding of the growth of Hawaiian volcanoes (Moore and Clague 1992; DePaolo and Stolper 1996; Sharp and Renne 2005; Sherrod et al. 2007), and the emergence of individual islands is sometimes older than previously assumed (e.g., Fleischer et al. 1998 and later).

An independent time calibration for the *Drosophila* phylogeny has recently become available through precise laboratory estimates of the mutation rate in *D. melanogaster* (e.g., Haag-Liautard et al. 2007). By equating mutation rate with synonymous substitution rate and by assuming all species share the same rate per unit time (based on 10 generations per year for *D. melanogaster*), Cutter (2008) estimated MRCA dates for members of the *melanogaster* group using genome-wide data, including an estimate of ~1.2 Ma for the common ancestor of *D. melanogaster* and *D. simulans* (0.6–1.9 Ma, 80% confidence interval [CI]). This seems remarkably close to the Hawaiian-calibrated date of 2.3 Ma provided by Russo et al. (1995), but it differs markedly from the 5.4 Ma of Tamura et al. (2004) (3.2–7.5 Ma, 95% CI). Furthermore, as many synonymous sites evolve under purifying selection in *Drosophila* (e.g., Halligan and Keightley

2006), calibrations based on mutation rate are likely to underestimate the true age.

In this study, we generate data for seven species of Hawaiian *Drosophila* and use explicit coalescent models to associate sequence divergence with the formation of Hawaiian Islands. We then use an experimental estimate for the mutation rate, derived by sequencing three *D. melanogaster* lines (Keightley et al. 2009), to provide an alternative calibration. A Bayesian phylogenetic approach (Drummond et al. 2006) allows us to naturally account for sources of uncertainty, such as errors in mutation rate estimates and the dates of island formation. Furthermore, by applying relaxed-clock models, we hope to provide a more realistic view of uncertainty due to variable rates of evolution across the phylogeny. These new divergence date estimates can be used to date key events in *Drosophila* and provide alternative hypotheses for the evolutionary history of this important model clade.

## Materials and Methods

### Speciation and Colonization Models for Hawai'i

Assuming colonization follows island formation rapidly, we can associate speciation with datable volcanic activity (Fleischer et al. 1998; Price and Clague 2002). However, in reality, there will be a delay before islands are colonized, and this adds uncertainty. The delay might be negligible: for example, Surtsey (formed 32 km from Iceland in 1963) was colonized by more than 50 plant species within 50 years (Fridriksson 2000). Alternatively, for endemic Hawaiian *Drosophila*, the delay before colonization could be much longer because they are adapted to cool, high-elevation rainforest (never below 500 m and usually above 1,000 m, P.M.O. personal observation and Hardy 1965). This requires soil that may take thousands of years to form and requires the volcano reach a substantial height. Their ecology also makes them susceptible to on-going eruptions, so that it is not the delay before colonization that is important, but the delay before the last (permanent) colonization. Although closed canopy forest can develop on lava flows within 400 years (Atkinson 1970), ash-fall data suggest that vegetation is regularly extirpated, and the newly formed landscape does not support Hawaiian *Drosophila* (Carson et al. 1990). The youngest volcano supporting *Drosophila* (Kilauea) is still active, but it is regularly recolonized from other volcanoes on the same island (Carson et al. 1990), and if it were not contiguous with a larger island, regular recolonization from a distant island would reset the speciation clock. If the only issue is soil formation and community assembly, then a colonization delay of ~10,000 years is negligibly small compared with the timescale of island formations (see later). Alternatively, if speciation requires the cessation of all volcanic activity above 1,000 m, then the delay may approach the hundreds of thousands of years that will impact date estimates.

A second major source of uncertainty is the difference between the time at which two lineages cease interbreeding (which we assume is the colonization date) and the coalescence time of sequences sampled from those two species.
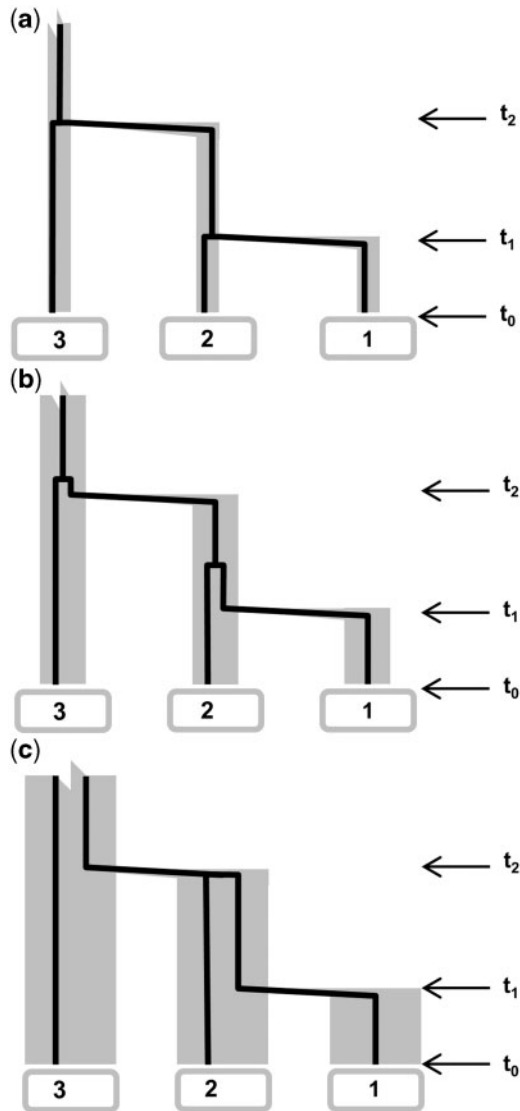
**Fig. 1.** Models linking speciation with coalescence. Panels (a–c) depict the coalescence of lineages (black lines) sampled at time $t_0$ from three island-endemic species (labeled 1, 2, and 3). In this idealized model, speciation events can be associated with the colonization of island 2 from island 3 at time $t_2$ and the colonization of island 1 from island 2 at time $t_1$. Associating coalescence with speciation is less easy. In panel (a), population sizes are extremely small (depicted by the width of the gray background), so that coalescence happens extremely rapidly and the coalescence date of sequences sampled from islands 1 and 2 can be approximated to the date at which island 1 was colonized. In panel (c) within-island population sizes are effectively infinite, but the bottleneck at colonization is very tight. Thus, coalescence does not occur on the island until the colonization bottleneck is reached, when coalescence becomes very rapid. Under this model, the coalescence date of island lineages 1 and 2 can be approximated by the formation of island 2 at time $t_2$. In panel (b), population sizes are intermediate, and coalescence time is determined by the effective population size on the ancestral island.

Sequences sampled from sister species will be derived from different alleles in the ancestral population, and these alleles may have shared a common ancestor long before migration onto the new island occurred (Nei 1971, see illustration in fig. 1). This effect will be greatest when the effective

population size is large, and the true coalescent time will lie between colonization of the donor island and migration onto the new island (fig. 1). Under a simple panmictic constant population model, the expected waiting time to coalescence for two alleles is $2 N_e$ generations, which for *Drosophila* is likely to be on the order of hundreds of thousands of years.

Coalescence in the ancestral population does not appear to have been explicitly modeled in previous studies of Hawaiian *Drosophila*, and some estimates have implicitly assumed the model illustrated in figure 1c, thus providing an upper limit on the age of speciation events. This model (fig. 1c) reflects a tight bottleneck at speciation causing rapid coalescence where local population sizes are otherwise extremely large (larger than continental populations of *Drosophila*) and thereby links between-species coalescence with the colonization of the ancestral island. For example, Bonacum et al. (2005) provided an upper limit for the MRCA of *D. hemipeza* and other species in the planitibia subgroup by associating the split with the formation of O'ahu (2.6–3 Ma) and the split between *D. differens*/*D. planitibia* and *D. silvestris*/*D. heteroneura* with the formation of Moloka'i (2 Ma). As *D. hemipeza* is endemic to O'ahu and *D. differens* endemic to Moloka'i (fig. 2), this model dates the divergence of the gene sequences as being over a million years before the species split. Although the widely cited studies of Russo et al. (1995) and Tamura et al. (2004) each used only used a single calibration date and chose to extrapolate over a much longer time frame, both also provided an upper-limit estimate by dating the MRCA of *D. picticornis* and the planitibia group by the formation of Kaua'i (then estimated at 5.1 Ma), which is over a million years before the earliest possible speciation date (the formation of O'ahu).

In this study, we take two alternative approaches. In both cases, we envision that the speciation event is directly associated with island formation, e.g., a species present on O'ahu colonized Moloka'i when Moloka'i was formed, giving rise to *D. differens* and leaving the population on Oahu that became *D. hemipeza*. This dates the separation of these species to the formation of Moloka'i. If population sizes are small and coalescence within populations correspondingly rapid, then coalescence may be close to the date at which Moloka'i formed (fig. 1a). If effective population sizes are large, coalescence will occur at some earlier date on O'ahu (fig. 1b). We model both scenarios: in model A (A1 and A2), we constrain sequence divergence times for all loci to the date for the newly formed island, whereas in model C (C1 and C2), we allow independent coalescence for each locus, assuming a constant effective population size on the ancestral island. Note that although model C does not account for the possibility of a bottleneck at colonization, this should not affect our results as the effective population size of Hawaiian *Drosophila* (see Results) means that coalescence is likely to occur more recently than the colonization of the ancestral island. Nevertheless, results from an "upper limit" model of colonization (model B1 and B2), which associates speciation with the age of the donor island (fig. 1c), are also provided in supplementary table S1, Supplementary Material online, for comparison with previous work.
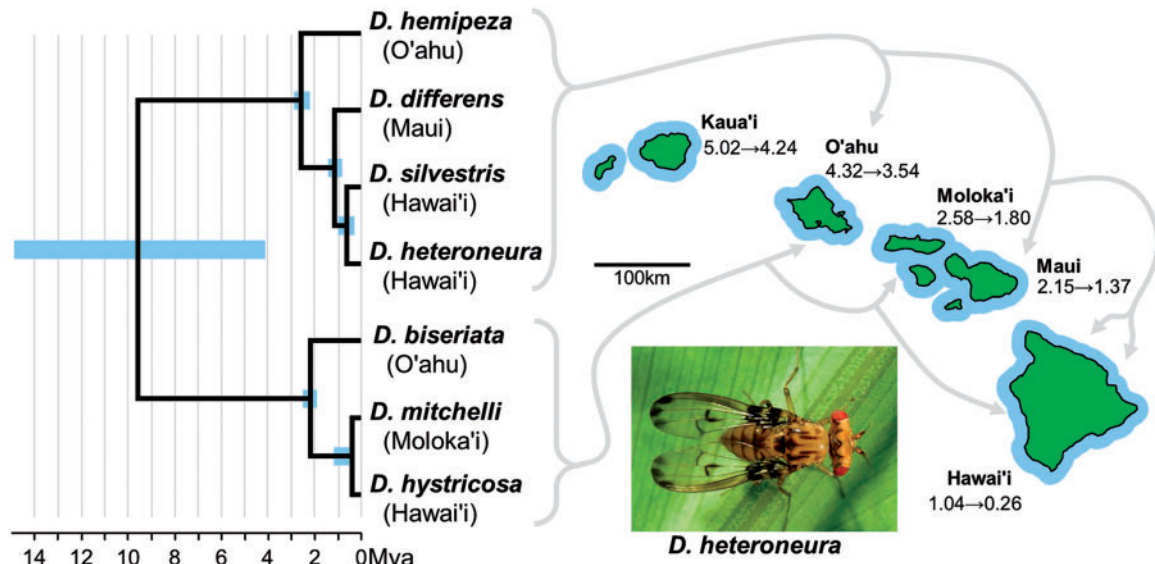
**Fig. 2.** *Drosophila* speciation on the Hawaiian Islands. As new Hawaiian islands are formed to the east, species from the nearest extant island are able to colonize the new island and become reproductively isolated (gray arrows). This "conveyer belt" speciation process has allowed the Hawaiian members of the Drosophilidae to radiate rapidly, forming a large and specioise group that display extreme morphological and behavioral diversity. This diversity includes the striking but now highly endangered "picture-wing" group (such as *D. heteroneura*, inset) that have been a major focus of *Drosophila* evolutionary ecology. The phylogeny (left) illustrates the inferred topology and speciation times of the seven species sampled for this study, with dates derived from model A1 (see main text), which sets speciation dates to the surface emergence of the first volcano of each island. Island dates are given as a span from the time of inferred surface emergence to shield completion for the oldest volcano on the island.

## Geological Dates for the Hawaiian Islands

Each island is composed of a mosaic of overlapping volcanoes, for example, Hawai'i itself comprises material from Kilauea and Mauna Loa, Mauna Kea, Hualalai, and Kohala. Volcanic growth takes place in stages (Stearns 1946; Clague and Dalrymple 1987; DePaolo and Stolper 1996): the eruption of basalts build a seamount on the ocean floor (Moore et al. 1982), followed by rapid growth of the volcano over several hundred thousand years (Moore and Clague 1992) until it breaks the sea surface and forms a shield volcano, eventually reaching 3 km above sea level (Peterson and Moore 1987). Volcanic material varies according to depth: deep-water eruptions form bulbous pillow lavas; shallow-water eruptions tend to form fragmental hyaloclastite; and subaerial flows are mostly massive sheet flows. After shield building, subaerial eruptions of alkali basalt take place and form a thin carapace over the shield volcanoes (MacDonald and Katsura 1964): this is known as the postshield stage and marks the end of significant growth of the volcano.

Genetic studies have tended to use K–Ar geological dates, rather than the more accurate Ar–Ar estimates now available (e.g., Sharp et al. 1996; Sharp and Renne 2005; Sherrod et al. 2007), and some have assumed that island age could be established directly from the oldest surface exposed rocks. However, the load of the volcano results in subsidence and burial of older land surfaces, so that volcanic rocks formed when the island first emerged above sea level may now have subsided far below the surface of the ocean (Moore and Clague 1992; DePaolo and Stolper 1996). Therefore, such studies may have underestimated the age of the islands, which will tend to lead to substitution rates being overestimated—the

opposite bias to the assumptions about the colonization and coalescent process.

We use a growth model for typical Hawaiian volcanoes (DePaolo and Stolper 1996), which predicts that the combined duration of the preshield and shield stages is close to 1 My and that a typical volcano would breach the sea surface ~0.22 My after the beginning of its growth or ~0.78 My before the eruption of the postshield phase lava (table 1). In the first set of phylogenetic models (A1 and C1), we combine this growth model with the available K–Ar or Ar–Ar dates (Sherrod et al. 2007 and references therein) to provide an improved estimate of emergence of the oldest volcano on each island. As the environment may be inhospitable to *Drosophila* during the shield building phase, we also use this growth model to estimate shield-completion dates for the oldest volcano on each island and use these as colonization dates in a second series of genetic models (model A2 and C2).

Uncertainty associated with extrapolating the Mauna Kea growth model of DePaolo and Stolper (1996) to other Hawaiian volcanoes is likely to be the most important source of uncertainty in estimating island emergence dates and is unfortunately difficult to constrain. The error is likely to be >0.11 My, because the emergence of Kau'ai estimated from the volcano growth model is 0.11 My younger than the hard minimum age provided by surface-exposed rocks, indicating that substantial errors must be associated with the growth models. Although DePaolo and Stolper (1996) estimate an average lifetime of 1 My for each volcano, others estimate growth durations between 0.65 and 1.5 My (Guillou et al. 2000). If this range is interpreted as a 2σ Gaussian error,

**Table 1.** Volcano Age Estimates.

| Island | Volcano | Oldest surface rocks | | Shield completion | | Shield emergence | |
|---|---|---|---|---|---|---|---|
| | | $t_m$ | $\sigma_m$ | $t_{ps}$ | $\sigma_{ps}$ | $t_e$ | $\sigma_e$ |
| Kau'ai | Waimea Canyon basalt | 5.13 | 0.06 | 4.24 | 0.09 | 5.02 | 0.19 |
| O'ahu | Wai'anae | 3.93 | 0.08 | 3.54 | 0.04 | 4.32 | 0.17 |
| Moloka'i | West Moloka'i | 1.99 | 0.08 | 1.8 | 0.04 | 2.58 | 0.17 |
| Maui | West Maui | 1.83 | 0.07 | 1.37 | 0.08 | 2.15 | 0.19 |
| Hawai'i | Kohala | 0.46 | 0.03 | 0.26 | 0.01 | 1.04 | 0.17 |

NOTE.—All ages are given in units of My before present. $t_m$ is the hard minimum age provided by the age of the oldest exposed rocks on the surface, with uncertainty $\sigma_m$ derived from the geochemical dating techniques. $t_{ps}$ is the end of the shield-building phase of the volcano with uncertainty $\sigma_{ps}$. $t_e$ is the age of volcano emergence from the growth model, with uncertainty $\sigma_e$.

then the $1\sigma$ error on the age of emergence estimates is 0.17 My. This uncertainty is broadly consistent with the discrepancy between recent dating of samples from the Hawai'i Scientific Drilling Project (HSDP-2) core through Mauna Kea and the prediction of the DePaolo and Stolper (1996) model, where the measured ages are ~0.15 My older than those predicted near the base of the core (Sharp and Renne 2005). This extrapolation error can be combined with the uncertainty of the K–Ar or Ar–Ar dates to provide an estimate of the uncertainty on the age of island emergence (table 1). Using a Bayesian approach allows us to include such uncertainty by informing the prior distribution for each constrained node date. In models A1 and C1 (island emergence), we incorporated a combined estimate of uncertainty associated with extrapolating the growth model and the analytical errors associated with the sample dates (table 1: $\sigma_e$). In models A2 and C2, we incorporated only the uncertainty associated with technical measurement error, as this date can be estimated directly from the oldest surface-exposed subaerial eruptions of alkali basalt (table 1: $\sigma_{ps}$).

Changes in sea level, due to changes in the volume of oceanic water and to subsidence of the islands, are also potentially important because they partly determine which separate islands were once contiguous: Moloka'i and Maui may once have been connected, although Maui has probably always been separate from Hawai'i, and O'ahu is unlikely to have been contiguous with Kaua'i (discussed in Fleischer et al. 1998). In this study, we choose not to account for this additional level of complexity, as its impact is likely to be small compared with other sources of uncertainty (see earlier) and it is harder to model (Price and Clague 2002).

### Sequencing from Hawaiian *Drosophila*

To improve the Hawaiian calibration of the *Drosophila* phylogeny, we sequenced 22 loci (317–577 bp) giving a total of ~10 kbp in each of seven species (Genbank JQ413030–JQ413183). Species were selected from two clades, which are both thought to have independently speciated with island formation: members of the planitibia subgroup (*D. hemipeza, D. differens, D. silvestris,* and *D. heteroneura,* provided by K. Kaneshiro in 2002 from laboratory stocks W40B14, Z79Z1, U26B9, and W48B6, respectively; DNA was extracted on that date) and members of the mitchelli subgroup (*D. biseriata, D. hystricosa,* and *D. mitchelli,* collected and

identified by P.M.O.). These species are all restricted to high-elevation rainforest: *D. hemipeza* above 600 m; *D. differens* 900–1,300 m; *D. silvestris* 1,100–1,700 m; *D. heteroneura* 900–1,800 m; *D. biseriata* above 550 m; *D. hystricosa* above 460 m; and *D. mitchelli* 900–1,200 m (data from collections reported in Hardy 1965). Under our colonization model, four of the six nodes on the phylogeny are datable (fig. 2). This is because both *D. silvestris* and *D. heteroneura* occur on Hawai'i and likely speciated since its formation and because the common ancestor of the planitibia and mitchelli subgroups likely predates the origin of the extant islands.

Loci were either selected because they have previously been advocated for phylogenetic inference in *Drosophila* or from the set of protein-coding loci that have 1:1 orthologs in the 12 completed drosophilid genomes (Clark et al. 2007), conditional on the fact that they include an exon longer than 700 bp, and do not have unusually high codon usage bias. Polymerase chain reaction (PCR) primers were designed to match the *D. grimshawi* sequence in regions that are highly conserved across the 12 genomes (see supplementary table S2, Supplementary Material online, for primers). Following PCR, unincorporated primers and deoxynucleotide triphosphates (dNTPs) were removed using exonuclease I and shrimp alkaline phosphatase, and the products were then sequenced in both directions using BigDye™ v3.1 (Applied Biosystems) and an ABI capillary sequencer (GenePool facility, University of Edinburgh). Sequence chromatograms were inspected by eye to confirm the validity of all variants within and between species and assembled using SeqMan (DNAstar Inc., Madison, WI).

### Sequences from Other Species

Divergence times for the 12 genomes were estimated using 50 protein coding loci (~63 kbp per species), all of which were 1:1 orthologs (Clark et al. 2007). We limited our choice to loci for which <10% of sites in the data matrix were missing and which had low overall codon-usage bias (the frequency of optimal codons, as identified by Vicario et al. 2007). Sequences were downloaded from ftp://flybase.net/genomes/12_species_analysis/clark_eisen/alignments/ and had previously been aligned and masked for regions of poor alignment (supplementary table S3, Supplementary Material online).

To estimate divergence times for the nine species of the *D. melanogaster* subgroup, we selected 36 protein coding loci (~38 kbp per species; supplementary table S4, Supplementary Material online). For six loci (*Adh*, *Amyrel*, *Est-6*, *per*, *rep4*, and *Ry*), sequences derived from multiple strains or isofemale lines were publically available in more than three species, and using these data, we ran an additional coalescent model (*BEAST; see Heled and Drummond 2010 and later). These sequences were derived from Genbank accessions (supplementary table S4, Supplementary Material online), from the *D. simulans* genome project (Begun et al. 2007), or from the Malawi accessions of the *D. melanogaster* population genomics project release 1.0 (http://www.dpgp .org/) (Langley et al. 2012).

## Models of Speciation and Sequence Evolution

All phylogenetic inference was performed with BEAST (Drummond and Rambaut 2007), and throughout we assume a relaxed-clock model of evolution in which rate variation between branches is modeled by a log-normal distribution (Drummond et al. 2006). Sequence evolution was modeled using a Hasegawa-Kishino-Yano (HKY) model (Hasegawa et al. 1985) with coding sequence partitioned into 1st + 2nd versus 3rd positions and rate variation between sites described by a four-category discrete gamma distribution. This model generally outperforms models which do not account for variation between codon positions (Shapiro et al. 2006). Base frequencies were estimated from the data, and base frequencies and between-site rate variation were unlinked between the codon-position partitions. All sequence-evolution parameters were given their default BEAST priors, except for tree shape (which was assumed to follow a birth–death speciation process), and parameters relating to the relaxed clock (see later). Unless otherwise specified, all loci in each analysis were assumed to share the same clock, substitution, and tree models. The models were each run at least twice, and stationarity was confirmed visually for all parameters by comparison between replicates of Markov chain Monte Carlo (MCMC) traces for all parameters. MCMC chain lengths varied between $10^8$ and $2 \times 10^9$ steps depending on the length needed to reach a sufficient effective sample size. Unless otherwise specified, the effective sample size for the posterior of each parameter was >500. Tree topologies were, with one exception, unconstrained; however, after preliminary runs confirmed that the topology linking 11 of the 12 genomes was very highly supported, the topology of these 11 species was constrained for the 250-gene data set to reduce computation times (the position of *D. willistoni* was unconstrained).

We first used Hawaiian island dates to infer a third-position substitution rate using the 22 newly sequenced loci from seven species of Hawaiian *Drosophila*. We did not limit the data to 4-fold degenerate codons, as it is not essential to use neutral sites in this analysis. We fitted four alternative models depending on the link between volcanic date and colonization/speciation date and the link between colonization/speciation date and sequence coalescence date

(summarized in table 2). These models were volcano emergence models (A1 and C1); shield completion models (A2 and C2); sequence divergence associated with formation of the newly colonized island (A1 and A2); and sequence coalescence on the donor island determined by effective population size (C1 and C2). In models A1 and A2, the four datable nodes of the gene tree model were constrained by the use of fully informative normal prior distributions based on volcano dates and the associated uncertainty in their estimates. For models C1 and C2, we used the *BEAST model (Heled and Drummond 2010) and constrained the four datable nodes of the species tree model using the volcano-derived prior distributions and fixed $N_e$ at $10^6$ for each species. This is equivalent to lower long-term estimates for *D. melanogaster* but likely an upper limit for Hawaiian *Drosophila*.

To estimate divergence times of the 12 published genomes using the new Hawaiian phylogeographic calibration, we used the posterior distributions of the rate estimates for third-position sites as fully informative priors on the mean of the lognormal distribution for third-position substitution rates in the 12-genome data set (e.g., Raghwani et al. 2012). Although substitution at these rapidly evolving sites will be prone to saturation, the assumption of a constant relationship between the rates at different codon positions provides relative date information. This permits date estimates from divergences that would otherwise be saturated if only third positions were used. For models A1 and A2, we found the posterior distribution of the third-position substitution rate could be closely approximated by an offset gamma distribution, whereas for models C1 and C2, a normal distribution provided a good approximation (supplementary table S5, Supplementary Material online, for summary statistics and parameters for approximating distributions). We similarly estimated sequence divergence times across the *melanogaster* subgroup using 36 loci and the Hawai'i-derived priors on substitution rate to calibrate the timescale.

We then estimated sequence divergence times across the 12 *Drosophila* genomes using the experimentally inferred neutral mutation rate of $3.46 \times 10^{-9}$ bp$^{-1}$ gen$^{-1}$ from *D. melanogaster* (Keightley et al. 2009), converted to a time-based estimate of 0.0346 bp$^{-1}$ My$^{-1}$ by assuming 10 generations per year (e.g., Cutter 2008). This was combined with the error associated with the per-generation mutation rate estimate (Keightley et al. 2009) to provide a fully informative normally distributed prior for the mean of the lognormal distribution that describes between-branch variation in substitution rates (mean 0.0346 substitutions bp$^{-1}$ My$^{-1}$ and standard deviation 0.00281). The variance of the log-normal distribution that describes rate variation between branches was estimated from the data using the prior: $\Gamma(1,0.1)$.

Any attempt to equate the mutation rate with the substitution rate requires completely neutral sites, which may not exist in *Drosophila* (e.g., Halligan and Keightley 2006). However, 4-fold degenerate sites in genes that do not have high levels of codon usage bias are likely to be a close approximation (Halligan and Keightley 2006) and have the added advantage of being relatively easy to align compared with other alternatives (e.g., short introns). We therefore limited

**Table 2.** Estimated Third-Codon Position Substitution Rates in Hawaiian *Drosophila*.

| Model | Model assumptions | | Root (Ma)[a] | Substitutions/bp/million years[b] | | |
|---|---|---|---|---|---|---|
| | Colonization date | $N_e$ | | Mean[c] | Lower 95% bound | Upper 95% bound |
| A1 | Island emergence | Very small | 9.6 (4.1–15) | 0.008 | 0.005 | 0.014 |
| C1 | Island emergence | $10^6$ | 10 (12–15) | 0.006 | 0.005 | 0.007 |
| A2 | Shield completion | Very small | 4.6 (2.4–7.1) | 0.019 | 0.011 | 0.029 |
| C2 | Shield completion | $10^6$ | 7.8 (6.1–9.8) | 0.010 | 0.008 | 0.012 |

[a]The inferred root date of the Hawaiian *planitibia* and *mitchelli* subgroups, with 95% highest posterior density intervals.
[b]Full distributions and model approximations are given in supplementary table S5, Supplementary Material online.
[c]Estimate of the mean of the lognormal distribution for the rate of third-codon positions (substitutions per base per million years).

the data set to 4-fold codons, and applied the prior distribution to third positions only (i.e., 4-fold degenerate sites).

We took an identical approach to estimate species divergence times in the *D. melanogaster* subgroup using the 36 locus data set. Because the importance of lineage sorting within the *D. melanogaster* subgroup is well established (e.g., Pollard et al. 2006) and because the shorter timescale of the *D. melanogaster* subgroup relative to the 12-species data set makes the impact of incomplete lineage sorting proportionately larger, we also used a *BEAST model (Heled and Drummond 2010) to allow sequence coalescence within a species-tree framework. In the absence of multiple sequences per species, we were unable to estimate effective population sizes and, therefore, chose to fix it at $10^6$ for each species, approximately equivalent to long-term effective population size estimated from *D. melanogaster* (Andolfatto and Przeworski 2000). To assess the impact of assuming this constant and universal population size, we ran an additional analysis using the six (of 36) loci for which multiple sequence accessions were available in more than three species (supplementary table S4, Supplementary Material online, for loci and species). This allowed us to estimate the effective population size from the data but did not allow for bottlenecks at speciation.

To assess the major sources of uncertainty, we ran additional analyses of the 12-genome data set using the experimental mutation rate. We explored the relative roles of uncertainty in the mutation rate, the overall quantity of sequence data, and the strictness of the clock by varying each in turn. In turn, we constrained the uncertainty in the mutation rate to 1/10 of its true value, analyzed a data set of 250 loci in place of 50 loci, and enforced a near-strict molecular clock (branch-to-branch standard deviation $10^{-6}$ in place of the estimated ∼0.3). We also explored the consequences of prior choice on the latter parameter, by selecting alternative priors for the standard deviation of the uncorrelated lognormal distribution that differed by roughly two orders of magnitude in their mean: prior $\Gamma(1,10)$ versus $\Gamma(1,0.1)$ for all other analyses.

## Results

### Tree Topology and Rate Inference from the Hawaiian Islands

We first assessed whether the phylogeny of the Hawaiian *Drosophila* matched the order of island formation, as this is a prerequisite for using the island ages to infer a rate of

sequence evolution. Although no locus individually supported the expected topology to the exclusion of other topologies (i.e., each was <95% of the posterior tree set), the concatenated data set strongly supported the expected topology (∼100% of the posterior tree set: supplementary fig. S1, Supplementary Material online) and only three loci individually provided substantial support for an alternative topology (supplementary fig. S1, Supplementary Material online).

Assuming that colonization occurs shortly after island emergence and that the population size is small, we estimate that the mean third-codon position substitution rate for Hawaiian *Drosophila* is 0.008 bp$^{-1}$ My$^{-1}$ 95% highest posterior density interval (0.005–0.014) (table 2, model A1). This estimate is fairly robust to assumptions about the effective population size: even if these Hawaiian species have effective population sizes typical of continental species of *Drosophila* ($N_e = 10^6$), then the mean substitution rate only decreases to 0.006 (0.005–0.007) bp$^{-1}$ My$^{-1}$ (table 2, model C1). This also allows us to estimate when the planitibia and mitchelli subgroups diverged: model A1 leads to an estimate of 9.6 (4–15) Ma, and model C1 gives estimates of 10 (12–15) Ma (fig. 2, table 2).

To assess which of these models is likely to be most realistic, we can estimate the population-scaled mutation rate ($\theta = 4N_e\mu$) of Hawaiian *Drosophila* from the genetic diversity of synonymous sites ($\pi_s$). As we only had a single (diploid) specimen of each species, we have calculated the mean of $\pi_s$ across all the loci we sequenced. The members of the mitchelli subgroup were all wild-caught specimens, and they had a mean $\pi_s$ of 0.0034 (*D. biseriata*, $\pi_s = 0.0030$; *D. hystricosa*, $\pi_s = 0.0017$ and *D. mitchelli*, $\pi_s = 0.0056$). The planitibia subgroup species examined herein were all derived from laboratory stocks, and we found a mixture of homozygous and heterozygous loci. If we assume the completely homozygous loci are the result of laboratory inbreeding, we can gain a rough estimate of $\pi_s$ from the remaining genes. This gave a mean of $\pi_s = 0.0040$, which is very similar to the mitchelli subgroup (*D. hemipeza*, $\pi_s = 0.0009$; *D. differens*, $\pi_s = 0.0048$; *D. silvestris*, $\pi_s = 0.0013$; and *D. heteroneura*, $\pi_s = 0.0091$). These values are ∼4- to 5-fold lower than those found in African populations of *D. melanogaster*, suggesting a long-term effective population size that is proportionately smaller and thus closer to model A1 than model C1.

These analyses assume that new islands are colonized soon after they emerge, which is likely to have occurred given the

**Table 3.** Divergence Dates (Ma) of *Drosophila* Species Based on Data from the 12 Published Genome Sequences.

| Dated nodes | Previously published estimates[a] | | Mutation rate[b] | Hawai'i shield-completion date[c] | | Hawai'i surface emergence date[d] | |
|---|---|---|---|---|---|---|---|
| | Russo et al. (1995) | Tamura et al. (2004) | | Model A2 | Model C2 | Model A1 | Model C1 |
| *simulans* complex | | 0.9 (0–1.9) | 0.5 (0.3–0.7) | 0.4 (0.2–0.5) | 1.3 (0.8–1.9) | 1.5 (0.6–2.6) | 2.3 (1.4–3.2) |
| *melanogaster-simulans* | 2.3 (1–3.6) | 5.4 (3.2–7.5) | 1.4 (0.9–1.9) | 1.1 (0.7–1.4) | 3.6 (2.4–5.0) | 4.2 (1.8–7.1) | 6.3 (4.3–8.6) |
| *yakuba-erecta* | . | 10 (6–15) | 2.4 (1.6–3.3) | 1.9 (1.3–2.6) | 6.4 (4.2–8.8) | 7.5 (3.0–12) | 11 (7.4–15) |
| *melanogaster* subgroup | 6.1 (3.9–8.3) | 13 (8–17) | 3.3 (2.4–4.4) | 2.7 (2.0–3.5) | 9.1 (6.6–12) | 11 (4.7–18) | 16 (11–21) |
| *melanogaster-annanassae* | . | 44 (27–62) | 15 (11–21) | 12 (8.4–16) | 40 (27–54) | 47 (20–79) | 70 (49–93) |
| *virilis/repleta*-hawaiian | 32 (26–38) | 43 (26–60) | 13 (9.6–17) | 12 (8.8–15) | 40 (30–54) | 47 (21–79) | 70 (50–91) |
| *mel-obscura* groups | 25 (19–31) | 55 (33–76) | 24 (17–31) | 19 (14–24) | 63 (46–82) | 73 (33–123) | 109 (80–142) |
| *Drosophila-Sophophora* | 40 (33–46) | 63 (39–87) | 32 (25–40) | 26 (21–32) | 89 (67–113) | 103 (47–170) | 154 (120–193) |

NOTE.—Estimates are the posterior means with 95% highest posterior density intervals.
[a]Derived from two previously published articles. Bounds are ±2x the reported standard error.
[b]Using only 4-fold degenerate codons, the prior of the rate of the third position was constrained to the estimate provided by Keightley et al. (2009).
[c]Using all codons and a Hawaiian calibration that associates speciation dates with the estimated completion of the first shield for that island.
[d]Using all codons and a Hawaiian calibration that associates speciation dates with the estimated emergence of the first volcano above the surface.

**Table 4.** Divergence Dates for Species within the Melanogaster Subgroup.

| Dated nodes | Mutation rate | | | Hawai'i shield-completion date | | Hawai'i surface emergence date | |
|---|---|---|---|---|---|---|---|
| | Linked gene tree | *BEAST[a] coalescent ($N_e = 10^6$) | *BEAST coalescent 6 loci (estimated $N_e$) | Model A2 | Model C2 | Model A1 | Model C1 |
| *simulans* complex | 0.7 (0.5–0.9) | 0.6 (0.4–0.7) | 0.5 (0.2–0.9) | 1.0 (0.5–1.6) | 2.0 (1.5–2.5) | 2.2 (1.0–3.8) | 3.4 (2.6–4.3) |
| *melanogaster–simulans* | 1.4 (1.1–1.8) | 1.3 (1.1–1.7) | 1.3 (0.8 −1.8) | 2.2 (1.2–2.5) | 4.3 (3.3–5.6) | 4.9 (2.4–8.5) | 7.5 (5.8–9.6) |
| *yakuba–santomea* | 0.6 (0.4–0.8) | 0.6 (0.4–0.7) | 0.8 (0.4–1.2) | 0.9 (0.5–1.4) | 1.7 (1.3–5.5) | 2.0 (0.9–3.5) | 3.0 (2.2–4.0) |
| *yakuba–teissieri* | 1.4 (1.1–1.7) | 1.4 (1.1–1.7) | 1.6 (1.1–1.3) | 2.2 (1.2–3.5) | 4.3 (3.3–5.5) | 4.9 (2.3–8.4) | 7.4 (5.9–9.5) |
| *erecta–orena* | 1.4 (1.0–1.8) | 1.6 (1.2–1.9) | 1.2 (0.5–2.0) | 2.2 (1.2–3.4) | 4.2 (3.1–5.5) | 4.7 (2.2–8.3) | 7.2 (5.4–9.4) |
| *erecta–yakuba* | 2.7 (2.2–3.3) | 3.0 (2.5–3.6) | 2.9 (2.0–4.0) | 4.3 (2.5–6.8) | 8.5 (6.6–11) | 9.5 (4.6–16) | 15 (12–18) |
| *melanogaster* subgroup | 3.4 (2.7–4.0) | 3.5 (2.9–4.1) | 3.4 (2.4–4.5) | 5.5 (2.9–8.3) | 11 (8.4–13) | 12 (5.8–21) | 18 (15–23) |

[a]May be estimated poorly, effective sample size (ESS) for these parameters was 70–100, based on two combined BEAST MCMC chains, each of $2 \times 10^9$ steps.

patterns of colonization that are observed on contemporary islands (see Materials and Methods). However, if colonization is delayed until the volcanic shield-building phase is complete (table 1), then this would roughly double our estimates of the substitution rates (table 2; models A2 and C2). This in turn leads to our estimates of divergence dates being considerably more recent (tables 3 and 4).

### MRCA Dates in the Genus *Drosophila* Calibrated by Hawaiian Islands

The substitution rates at third-codon positions estimated using Hawaiian *Drosophila* allow us to estimate divergence dates of the 12 species of *Drosophila* with published genomes (table 3). In these analyses, we estimated gene coalescent dates, so these will be older than speciation dates. Using 50 genes with low codon-usage bias from the published *Drosophila* genomes and assuming the colonization-on-emergence model for Hawai'i (and that Hawaiian flies have extremely small population sizes: model A1), we find the posterior means for the MRCA of *D. melanogaster* and *D. simulans* to be 4.2 (1.8–7.1) Ma (fig. 3 and table 3; model A1). This increases to 6.3 (4.3–8.6) Ma if Hawaiian population

sizes are large (fig. 3 and table 3; model C1). The corresponding dates for the MRCA of the subgenera *Drosophila* and *Sophophora* are 103 (47–170) Ma and 154 (120–193) Ma (fig. 3 and table 3). Dates of the other main divergences are listed in table 3, and a tree calibrated with model A1 is illustrated in figure 4.

We also used the Hawaiian estimates of the substitution rate to infer sequence divergence times within the *melanogaster* subgroup, based on a data set of 36 publically available loci. Using the same approach as for the 12 genomes, we applied these calibrations to the 36-locus data set (table 4 and fig. 4 for a tree calibrated with Hawaiian model A1), this resulted in estimates of gene coalescent times for the MRCA of *D. melanogaster* and *D. simulans* of 4.9 (2.4–8.5) Ma (model A1) and 12 (5.8–21) Ma for the MRCA of the whole subgroup.

### MRCA Dates in the Genus *Drosophila* Calibrated by Mutation Rate

Using laboratory measurements of the mutation rate in *D. melanogaster* as an estimator of the substitution rate at 4-fold degenerate sites, we date the MRCA of the subgenera
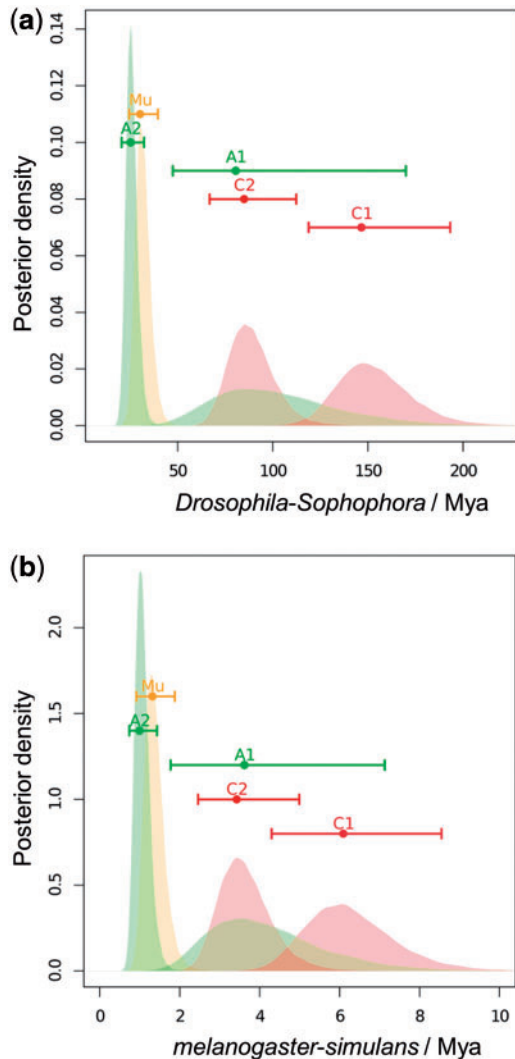
**Fig. 3.** Posterior distributions for key node dates under different models. Posterior distributions (shaded curves) and 95% highest posterior density intervals (solid lines) for the MRCA of the subgenera *Drosophila* and *Sophophora* (panel *a*), and the MRCA of *Drosophila melanogaster* and *D. simulans* (panel *b*), under five different models. Model "Mu" (yellow) uses an experimental estimate of the mutation rate, models "A1" and "A2" (green) use a Hawaiian calibration with small effective population size (i.e., instant coalescence), and models "C1" and "C2" (red) use a Hawaiian calibration with coalescence in an effective population size of $N_e = 10^6$ (see main text for details).

*Drosophila* and *Sophophora* to 32 (24–40) Ma (table 3). In the same analysis, we find the MRCA of *D. melanogaster* and *D. simulans* to be 1.4 (0.9–1.9) Ma. The mutation rate estimates are therefore considerably more recent than dates generated using the Hawaiian models A1 and C1. A summary of all node dates is presented in table 3 and illustrated in figure 4.

To obtain divergence dates within the *melanogaster* subgroup, we used the experimentally estimated mutation rate together with the data set of 36 loci (fig. 5 and supplementary table S4, Supplementary Material online). When all loci were constrained to follow the same inferred topology, we estimated the MRCA of *D. melanogaster* and *D. simulans* to be 1.4 (1.1–1.8) Ma, the MRCA of *D. erecta* and *D. yakuba* to be

2.7 (2.2–3.3) Ma, and the MRCA of the *melanogaster* subgroup to be 3.4 (2.7–4.1) Ma. These estimates are very similar to those derived from the 12-genome data set (compare figs. 4 and 5). In general, the mutation rate calibration results in dates that are considerably more recent than the Hawaiian dates (although some of the Hawaiian dates are inferred with extremely low precision; tables 3 and 4).

## Sources of Uncertainty

Variation in the substitution rate among different branches may be an important source of uncertainty when trying to estimate the divergence dates of the different groups of *Drosophila*. Because we only used small numbers of species, we expected that there would be little information to estimate relaxed-clock parameters, and we therefore chose to set a tight prior on this parameter, close to zero (i.e., consistent with a nearly constant molecular clock). Nevertheless, in the analysis of the 12-genome data set calibrated with the mutation rate, the inferred rate variation between branches was affected little by our choice of prior for this parameter: log-normal priors for the variance of rates across branches that differed in their mean by approximately two orders of magnitude resulted in posterior estimates that differed by only 20% (posterior 0.29 [0.18–0.41] vs. posterior 0.35 [0.19–0.52]). Although this does not provide a rigorous test, it does suggest that there is information to estimate the amount of between-branch rate variation and that the 12 *Drosophila* genomes do not conform to a strict molecular clock (posterior estimates of the variance for the between-branch log-normal distribution do not overlap zero).

The uncertainty associated with our estimates of divergence dates reflects the imprecision with which the *D. melanogaster* substitution/mutation rate has been estimated, the unknown rate variation among branches of the tree, and limits imposed by the finite amount of sequence data used for tree inference. Models that use more loci (n = 250 loci; MRCA of the genus *Drosophila* estimated at 29 [23–36] Ma) or less uncertainty in the mutation rate (standard deviation [SD] of the mutation rate estimate fixed at 1/10th of its true value; 33 [26–41] Ma) both still had large credibility intervals on node dates. However, if we assume there is no between-branch rate variation (i.e., approaching a strict clock), the date estimates become slightly more precise (30 Ma [25–35]; SD for the lognormal distribution fixed at $1 \times 10^{-6}$, compared with its estimated value of ca. 0.3).

Within the *melanogaster* subgroup, lineage sorting has resulted in incongruence between gene trees and the species tree for many loci (e.g., Pollard et al. 2006). We therefore wished to explore the consequences of having constrained all loci to share the same topology on our estimates of divergence dates within this group. First, under a species tree model, in which loci are constrained to coalesce within an inferred species topology (*BEAST, Heled and Drummond 2010) with $N_e$ fixed at $10^6$, we estimated the MRCA of *D. melanogaster* and *D. simulans* to be 1.3 (1.1–1.6) Ma, the MRCA of *D. erecta* and *D. yakuba* to be 3.0 (2.4–3.5) Ma,
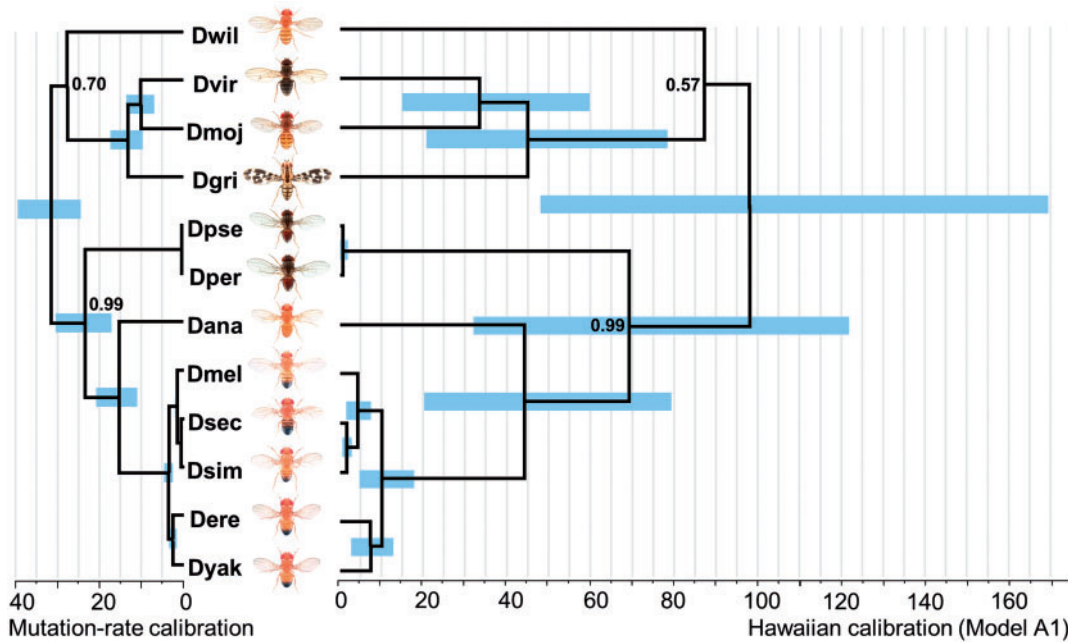
**Fig. 4.** A time-scaled phylogeny of *Drosophila*. Two alternative calibrations for a phylogenetic tree linking the 12 species of *Drosophila* for which complete genomes have been published, inferred from 50 one-to-one orthologs with low average codon usage bias. Trees were inferred under an uncorrelated log-normal relaxed clock, and node dates are scaled to the posterior median. The 95% highest posterior density date intervals (in blue) are shown for well-supported nodes and reflect uncertainty both in the rate estimate used to calibrate the tree and sampling error associated with the data. Nodes with less than 100% posterior support are labeled. The left-hand tree is based on 4-fold codons and was inferred by setting the prior distribution for the substitution rate for third positions to be normally distributed with a mean rate equal to a laboratory estimated neutral mutation rate and a variance that reflects the uncertainty in that estimate. The right-hand tree is based on all codons and was inferred by setting the prior distribution for the substitution rate of third positions to be an offset gamma distribution scaled to match the posterior distribution of this parameter inferred from Hawaiian calibration A1 (see main text). Note that the position of *D. willistoni* relative to the root differs from most previous studies. Given inferred fossil dates for *Drosophila* and Diptera as a whole, the timescale from Model A1 appears implausibly old (see Discussion). Species abbreviations: Dwil, *D. willistoni*; Dvir, *D. virilis*; D moj, *D. mojavensis*; Dgri, *D. grimshawi*; Dpse, *D. pseudoobscura*; Dper, *D. persimilis*; Dana, *D. annanassae*; Dmel, *D. melanogaster*; Dsec, *D. sechellia*; Dsim, *D. simulans*; Dere, *D. erecta*; and Dyak, *D yakuba*.
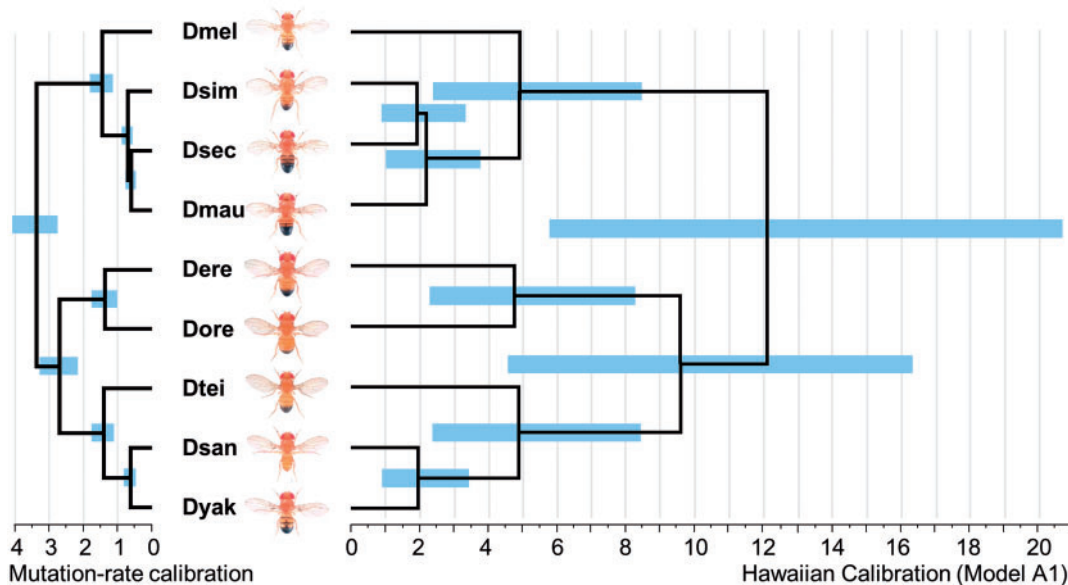


**Fig. 5.** A time-scaled phylogeny of the *melanogaster* subgroup. Two alternatively calibrated phylogenetic trees for the *melanogaster* subgroup, inferred from 36 loci. Trees were inferred under an uncorrelated log-normal relaxed clock, and the topology presented is the maximum clade-credibility tree with node dates scaled to the posterior median. The 95% highest posterior density intervals (in blue) are shown for each node, and reflect uncertainty both in the rate estimate used to calibrate the tree and sampling error associated with the data. Nodes with less than 100% posterior support are labeled. The alternative time calibrations (mutation rate vs. Hawai'i) are the same as those used in figure 4 (see main text). Note that the inferred topology of the *D. simulans* complex differs between the two data sets. Species name abbreviations are given in figure 4.

and the MRCA of the melanogaster subgroup to be 3.5 (2.9–4.1) Ma. However, we recommend caution when making interpretations from the tails of these distributions (i.e., the bounds) as some parameters represent an effective sample size of only 80–100, due to inefficient mixing in the MCMC. Second, using six loci for which there were multiple sequence accessions per species, we are able to coestimate $N_e$ from the data in the analysis (*BEAST Heled and Drummond 2010). Using this approach, we estimated the MRCA of *D. melanogaster* and *D. simulans* to be 1.3 (0.8–1.8) Ma, the MRCA of *D. erecta* and *D. yakuba* to be 2.9 (2.0–3.9) Ma, and the MRCA of the melanogaster subgroup to be 3.5 (2.4–4.6) Ma. Remarkably, these different estimates for speciation dates within the melanogaster subgroup are in very close agreement with the estimates obtain in the earlier analyses, despite the fact that some equate sequence divergence with speciation and others do not.

Finally, we emphasize the importance of avoiding genes with high codon usage bias in these analyses. To explore the consequences of codon usage bias, we repeated the 12-genome analysis using 50 genes with high codon-usage bias in place of genes with low codon usage. We found that estimates for the MRCA for *Drosophila* and *Sophophora* are about twice as recent (14 [9–18] Ma) and that those for *D. melanogaster-D. simulans* are ∼25% more recent (1.0 [0.7–1.4] Ma). This suggests that, as expected, the bias introduced by weak constraint is larger over longer timescales.

## Discussion

It is important for many evolutionary analyses to place an absolute timescale on a phylogeny, and for *Drosophila*, this has often been done by associating the speciation of Hawaiian species with the formation of Hawaiian islands (e.g., Carson 1976; Easteal and Oakeshott 1985; Thomas and Hunt 1993; Russo et al. 1995; Tamura et al. 2004). In our analyses, we have revisited these estimates with an expanded data set and improved evolutionary models and compared these results with estimates based on laboratory mutation rates. We found that the different methods gave widely varying results and that there is often great uncertainty associated with the estimates. Below we discuss why these differences arise and make recommendations regarding which estimates are likely to be the most reliable.

### Phylogeographic versus Mutation Rate Calibration

The experimentally estimated mutation rate was considerably higher than the substitution rate estimated from the Hawaiian *Drosophila* using our a priori favored model of the colonization (model A1), resulting in more recent estimates of divergence dates. For example, the mutation rate calibration suggests that the MRCA of the genus *Drosophila* lived ∼32 (25–40) Ma, whereas the Hawaiian A1 model, in which islands are colonized shortly after emergence, suggests this occurred 103 (47–170) Ma. This assumes that the precolonization delay is short (thousands to tens of thousands of years), but if we instead assume that on-going volcanic activity is incompatible with speciation—such that

islands are not colonized until volcanic activity has completely ceased at the end of the shield-building phase (Hawaiian model A2)—we get an estimate of only 26 (21–32) Ma.

Weigmann et al. (2011) recently estimated divergence dates in the Diptera using a fossil calibration, and they estimated that the Drosophilidae (including genera basal to *Drosophila*) had a common ancestor ∼50 Ma and that the Schizophora (a major group of flies that includes many families) arose less than 100 Ma. The limited fossil data available within the Drosophilidae suggest a minimum date of 20–50 Ma for the origin of the family and minimum of ∼20 Ma for the genus *Drosophila* (Grimaldi 1987, 1988). Comparing our results to the fossil-calibrated estimated divergence dates in the Diptera as a whole, it therefore seems that the Hawaiian calibration using island emergence (model A1) results in implausibly ancient divergence dates (>100 Ma for *Drosophila*), whereas the shield-completion (model A2) results in unexpectedly recent divergence dates (<30 Ma). Surprisingly, the mutation rate calibration is actually the closest to being compatible with fossil evidence (just >30 Ma).

### Uncertainty in the Mutation Rate Calibration

Our conclusion that laboratory estimates of the mutation rate lead to potentially realistic estimates of ancient divergence dates runs contrary to work on other taxa, where it has been found that this approach can lead to erroneously recent divergence dates, and substitution rates that are too high (see references in Ho et al. 2007). One important reason why substitution rates may often be lower than mutation rates is the action of purifying selection on the sites being studied. We have attempted to mitigate this effect by the use of 4-fold sites and genes with low codon bias and that the good match between our estimates and fossil dates may superficially suggest our approach has been successful.

However, a potentially much larger source of error in the mutation rate calibration derives from uncertainty in the generation time. Experimental estimates of the mutation rate are per generation, so that absolute dates can only be obtained by making strong assumptions about generation time. In this study, we have assumed 10 generations per year in the wild (as used previously in this context, e.g., Cutter 2008). However, if the true value is different, then the mutation rate calibration would need to be altered proportionately. For example, 20 generations per year (which would be more consistent with generation times seen in the laboratory) implies a date of only 16 Ma for the origin of the genus *Drosophila*, much more recent than would be plausible based on fossil data.

### Uncertainty in the Phylogeographic Calibration

The data strongly support a topology consistent with the island formation/speciation model (fig. 2 and supplementary fig. S1, Supplementary Material online), suggesting that the use of Hawaiian island formation to infer speciation rates is viable. However, the model in which the Hawaiian *Drosophila* colonize islands immediately after island emergence (model A1) gives dates that are implausibly ancient when compared

with fossil data and have low precision (table 3). There are two good candidate explanations for why the island emergence model (A1) appears to seriously underestimate the substitution rate. First, the molecular clock correlates with generation time in invertebrates (Thomas et al. 2010), and it is therefore possible that our result is a consequence of systematic differences in generation times between Hawaiian *Drosophila* and most *Drosophila* in our phylogeny. The Hawaiian species often have long generation times—*D. silvestris* and *D. heteroneura* typically only have about four generations per year in the laboratory (Boake et al. 1998). In contrast, most other species reproduce far more rapidly—*D. melanogaster* will go from egg to adult in ∼11 days at 22°C in the laboratory—and are consequently expected to have higher substitution rates. Therefore, if a per-year substitution rate is estimated in slowly reproducing Hawaiian species and applied to a group with a shorter generation time, without accounting for this difference, this will lead to overestimates of divergence dates—as we appear to have observed. Indeed, if the per-year mutation rate of Hawaiian species in the wild is 3–4 times lower than the rest of the genus, then the divergence dates estimated from the Hawaiian calibration (model A1) and the mutation rate calibration would be very similar.

Second, our assumption that flies colonize new islands shortly after they break the ocean surface may be incorrect. In particular, if colonization is delayed until the end of the shield-building phase (several hundred thousand years), then the alternative calibration (model A2) may be appropriate. This model is close to being compatible with *Drosophila* fossil calibrations and with the analysis by Weigmann et al. (2011) based on multiple Dipteran fossils. However, at 26 (21–32) Ma for the common ancestor of *Drosophila*, it is still much more recent than previous estimates.

## Comparison with Previous Estimates

The original motivation for this reanalysis of *Drosophila* divergence dates was that previous estimates from the Hawaiian *Drosophila* were based on small amounts of data, had used simple models of island colonization and gene coalescence, and had tended to underestimate the geological age of the islands. Despite this, our estimates are surprisingly similar to Russo et al. (1995), who applied a Hawaiian-derived substitution rate estimate of $1 \times 10^{-8}\,\text{bp}^{-1}\,\text{year}^{-1}$ (based on Adh) to infer dates for several nodes in the *Drosophila* phylogeny (table 3). This agreement may partly be a matter of chance, as the biases due to the assumptions about the age of the islands and the colonization process go in opposite directions and therefore have little net effect. However, this does not mean that they are individually unimportant. For example, if we assume the "upper limit" colonization and coalescent model adopted by some previous studies (fig. 1) in our analyses, we estimate that the common ancestor of subgenera *Sophophora* and *Drosophila* existed 192 Ma as opposed to 103 Ma in model A1 (supplementary table S1, Supplementary Material online).

The most sophisticated and data-rich molecular dating for *Drosophila* currently available is that of Tamura et al. (2004) (see table 3 for dates), who applied a Hawaiian calibration derived from *Adh* sequences to other *Drosophila* species. A primary aim of that study was to account for the effects of constraint (i.e., substitution at less than the neutral rate at some synonymous sites). It is particularly striking that for model A1 many of our recent dates are similar to theirs, but our estimates of the age of the root of the tree tend to be much older (table 3). We have explored this discrepancy by applying the method of Tamura et al. (2004) to our 250-gene data set derived from the 12 genomes, and we find that the difference may be explained by substitution saturation of 4-fold sites (supplementary fig. S2, Supplementary Material online). Specifically, single-locus TN93 (Tamura and Nei 1993) divergence estimates based only on 4-fold sites, as used by Tamura et al. (2004), substantially underestimate divergence when it approaches or exceeds 200% (supplementary fig. S2, Supplementary Material online). Our own approach attempts to mitigate the problem of saturation by including both low-constraint (4-fold or third position) and high-constraint sites (first and second codon positions) and modeling the relative rates of evolution at different sites (model SRD06; Shapiro et al. 2006). If the substitution model is a good description of the evolutionary process, the loss of information through saturation should be reflected in the resulting credibility intervals. Nevertheless, this approach makes the strong assumption that the relative rates of different codon positions are constant across the tree and this will not be the case if effective population sizes (and thus the efficacy of selection) vary, or if there are a substantial number of weakly deleterious variants that remain polymorphic for an extended period before selective removal.

## Conclusion

We have found that divergence dates in the genus *Drosophila* derived from laboratory estimates of the mutation rate are broadly compatible with evidence from fossils and suggest that these (figs. 4 and 5, tables 3 and 4) might be used until more data (e.g., generation time in the wild, dated fossils, and mutation rate estimates from multiple species) and better models (e.g., the inclusion of covariate traits such body size, mutation rate, and generation time) are available. Surprisingly, our results seem to suggest that laboratory-derived mutation rates may be a viable proxy for the substitution rate at 4-fold degenerate sites in genes with low codon bias, though this is highly conditional on our estimate of generation time in the field. In contrast, divergence dates based on the colonization of the Hawaiian Islands suffer from considerable uncertainty, and our a priori preferred model (colonization associated with island emergence) results in implausibly ancient divergence dates when compared with fossil estimates. Consequently, we believe estimates of absolute divergence dates in *Drosophila* should still be treated with some caution.

## Supplementary Material

## Acknowledgments

## References

Andolfatto P, Przeworski M. 2000. A genome-wide departure from the standard neutral model in natural populations of *Drosophila*. *Genetics* 156:257–172.

Ashburner M, Bodmer M, Lemeunier F. 1984. On the evolutionary relationships of *Drosophila melanogaster*. *Dev Genet.* 4:295–312.

Atkinson IAE. 1970. Successional trends in the coastal and lowland forest of Mauna Loa and Kilauea volcanoes, Hawaii. *Pacific Sci.* 24:387–400.

Begun DJ, Holloway AK, Stevens K, et al. 2007. Population genomics: whole-genome analysis of polymorphism and divergence in *Drosophila simulans*. *PLoS Biol.* 5:2534–2559.

Boake CRB, Price DK, Andreadis DK. 1998. Inheritance of behavioural differences between two interfertile, sympatric species, *Drosophila silvestris* and *D. heteroneura*. *Heredity* 80:642–650.

Bonacum J, O'Grady PM, Kambysellis M, DeSalle R. 2005. Phylogeny and age of diversification of the *planitibia* species group of the Hawaiian *Drosophila*. *Mol Phylogenet Evol.* 37:73–82.

Bromham L, Penny D. 2003. The modern molecular clock. *Nat Rev Genet.* 4:216–224.

Cao H, Wang X, Gao J, Prigent SR, Watabe H, Zhang Y, Chen H. 2011. Phylogeny of the African and Asian *Phortica* (Drosophilidae) deduced from nuclear and mitochondrial DNA sequences. *Mol Phylogenet Evol.* 61:677–685.

Carson HL. 1976. Inference of the time of origin of some *Drosophila* species. *Nature* 259:395–396.

Carson HL, Lockwood JP, Craddock EM. 1990. Extinction and recolonisation of local populations on a growing shield volcano. *Proc Natl Acad Sci USA.* 87:7055–7057.

Charlesworth D. 2010. Don't forget the ancestral polymorphisms. *Heredity* 105:509–510.

Clague DA, Dalrymple GB. 1987. Tectonics, geochronology and origin of the Hawaiian-Emperor volcanic chain. In: Decker RW, Wright TL, Stauffer PH, editors. US Geological Survey Professional Paper 1350. Washington (DC): US Government Printing Office. p. 1–55.

Clark AG, Eisen MB, Smith DR, et al. (417 co-authors). 2007. Evolution of genes and genomes on the *Drosophila* phylogeny. *Nature* 450:203–218.

Cutter AD. 2008. Divergence times in *Caenorhabditis* and *Drosophila* inferred from direct estimates of the neutral mutation rate. *Mol Biol Evol.* 25:778–786.

Da Lage JL, Kergoat GJ, Maczkowiak F, Silvain JF, Cariou ML, Lachaise D. 2007. A phylogeny of Drosophilidae using the *Amyrel* gene: questioning the *Drosophila melanogaster* species group boundaries. *J Zool Syst Evol Res.* 45:47–63.

DePaolo DJ, Stolper EM. 1996. Models of Hawaiian volcano growth and plume structure: implications of results from the Hawaii Scientific Drilling Project. *J Geophys Res Solid Earth.* 101:11643–11654.

DeSalle R, Grimaldi DA. 1991. Morphological and molecular systematics of the Drosophilidae. *Annu Rev Ecol Syst.* 22:447–475.

Drummond AJ, Ho SY, Phillips MJ, Rambaut A. 2006. Relaxed phylogenetics and dating with confidence. *PLoS Biol.* 4:699–710.

Drummond AJ, Rambaut A. 2007. BEAST: Bayesian evolutionary analysis by sampling trees. *BMC Evol Biol.* 7:214.

Easteal S, Oakeshott JG. 1985. Estimating divergence times of *Drosophila* species from DNA sequence comparisons. *Mol Biol Evol.* 2:87–91.

Fleischer RC, McIntosh CE, Tarr CL. 1998. Evolution on a volcanic conveyor belt: using phylogeographic reconstructions and K-Ar-based ages of the Hawaiian Islands to estimate molecular evolutionary rates. *Mol Ecol.* 7:533–545.

Fridriksson S. 2000. Vascular plants on Surtsey, Iceland, 1991–1998. *Surtsey Res.* 11:21–28.

Gao JJ, Hu YG, Toda MJ, Katoh T, Tamura K. 2011. Phylogenetic relationships between Sophophora and Lordiphosa, with proposition of a hypothesis on the vicariant divergences of tropical lineages between the Old and New Worlds in the family Drosophilidae. *Mol Phylogenet Evol.* 60:98–107.

Gao JJ, Watabe HA, Aotsuka T, Pang JF, Zhang YP. 2007. Molecular phylogeny of the *Drosophila* obscura species group, with emphasis on the Old World species. *BMC Evol Biol.* 7.

Grimaldi DA. 1987. Amber Fossil Drosophilidae (Diptera), with Particular Reference to the Hispaniolan taxa. *Am Museum Novitates.* 2880:1–23.

Grimaldi DA. 1988. Relicts in the Drosophilidae (Diptera). In: Liebherr JK, editor. Zoogeography of Caribbean insects. Ithaca (NY): Cornell University Press. p. 183–213.

Grimaldi DA. 1990. A phylogenetic, revised classification of the genera in the Drosophilidae (Diptera). Bulletin of the American Museum of Natural History. New York: American Museum of Natural History. p. 139.

Guillou H, Sinton J, Laj C, Kissel C, Szeremeta N. 2000. New K-Ar ages of shield lavas from Waianae Volcano, Oahu, Hawaiian Archipelago. *J Volcanol Geotherm Res.* 96:229–242.

Haag-Liautard C, Dorris M, Maside X, Macaskill S, Halligan DL, Charlesworth B, Keightley PD. 2007. Direct estimation of per nucleotide and genomic deleterious mutation rates in *Drosophila*. *Nature* 445:82–85.

Halligan DL, Keightley PD. 2006. Ubiquitous selective constraints in the *Drosophila* genome revealed by a genome-wide interspecies comparison. *Genome Res.* 16:875–884.

Hardy DE. 1965. Diptera: Cyclorrhapha II, series Schizophora, section Acalypterae, family Drosophilidae. Honolulu (HI): University of Hawaii Press.

Hasegawa M, Kishino H, Yano TA. 1985. Dating of the the human-ape splitting by a molecular clock of mitochondrial DNA. *J Mol Evol.* 22: 160–174.

Heads M. 2011. Old Taxa on young Islands:a critique of the use of island age to date island-endemic clades and calibrate phylogenies. *Syst Biol.* 60:204–218.

Heled J, Drummond AJ. 2010. Bayesian inference of species trees from multilocus data. *Mol Biol Evol.* 27:570–580.

Ho SY, Lanfear R, Bromham L, Phillips MJ, Soubrier J, Rodrigo AG, Cooper A. 2011. Time-dependent rates of molecular evolution. *Mol Ecol.* 20:3087–3101.

Ho SY, Larson G. 2006. Molecular clocks: when times are a-changin'. *Trends Genet.* 22:79–83.

Ho SY, Phillips MJ. 2009. Accounting for calibration uncertainty in phylogenetic estimation of evolutionary divergence times. *Syst Biol.* 58: 367–380.

Ho SY, Shapiro B, Phillips MJ, Cooper A, Drummond AJ. 2007. Evidence for time dependency of molecular rate estimates. *Syst Biol.* 56: 515–522.

Keightley PD, Trivedi U, Thomson M, Oliver F, Kumar S, Blaxter ML. 2009. Analysis of the genome sequences of three *Drosophila melanogaster* spontaneous mutation accumulation lines. *Genome Res.* 19: 1195–1201.

Kellermann V, Loeschcke V, Hoffmann AA, Kristensen TN, Fløjgaard C, David JR, Svenning J-C, Overgaard J. Forthcoming 2012. Phylogenetic constraints in key functional traits behind species' climate niches: patterns of dessication and cold resistance across 95 *Drosophila* species. *Evolution.* Advance Access published May 28, 2012, doi:10.1111/j.1558-5646.2012.01685.x

Kopp A. 2006. Basal relationships in the *Drosophila melanogaster* species group. *Mol Phylogenet Evol.* 39:787–798.

Langley CH, Stevens K, Cardeno C, et al. (17 co-authors). 2012. Genomic variation in natural populations of *Drosophila melanogaster*. *Genetics.* Advance Access published June 5, 2012, doi:10.1534/genetics.112.142018

MacDonald GA, Katsura T. 1964. Chemical composition of Hawaiian lavas1. *J Petrol.* 5:82–133.

Moore JG, Clague DA. 1992. Volcano growth and evolution of the island of Hawaii. *Geol Soc Am Bull.* 104:1471–1484.

Moore JG, Clague DA, Normark WR. 1982. Diverse basalt types from Loihi seamount, Hawaii. *Geology* 10:88–92.

Nei M. 1971. Interspecific gene difference and evolutionary time estimated from electrophoretic data on protein identity. *Am Nat.* 105: 385.

O'Grady P, DeSalle R. 2008. Out of Hawaii: the origin and biogeography of the genus *Scaptomyza* (Diptera: Drosophilidae). *Biol Lett.* 4: 195–199.

O'Grady PM, Lapoint RT, Bonacum J, Lasola J, Owen E, Wu Y, DeSalle R. 2011. Phylogenetic and ecological relationships of the Hawaiian *Drosophila* inferred by mitochondrial DNA analysis. *Mol Phylogenet Evol.* 58:244–256.

Oliveira DC, Almeida FC, O'Grady PM, Armella MA, DeSalle R, Etges WJ. Forthcoming 2012. Monophyly, divergence times, and evolution of host plant use inferred from a revised phylogeny of the *Drosophila* repleta species group. *Mol Phylogenet Evol.* Advance Access published May 24, 2012, doi:10.1016/j.ympev.2012.05.012

Pelandakis M, Solignac M. 1993. Molecular phylogeny of *Drosophila* based on ribosomal RNA sequences. *J Mol Evol.* 37:525–543.

Peterson DW, Moore RB. 1987. Geologic history and evolution of geologic concepts, Island of Hawaii. In: Decker RW, Wright TL, Stauffer PH, editors. Volcanism in Hawaii. U.S. Geological Survey Professional Paper 1350. Reston (VA): U.S. Geological Survey. p. 149–190.

Peterson GI, Masel J. 2009. Quantitative prediction of molecular clock and Ka/Ks at short timescales. *Mol Biol Evol.* 26:2595–2603.

Pollard DA, Iyer VN, Moses AM, Eisen MB. 2006. Widespread discordance of gene trees with species tree in *Drosophila*: evidence for incomplete lineage sorting. *PLoS Genet.* 2:1634–1647.

Price JP, Clague DA. 2002. How old is the Hawaiian biota? Geology and phylogeny suggest recent divergence. *Proc R Soc London Ser B Biol Sci.* 269:2429–2435.

Raghwani J, Thomas X, Koekkoek S, et al. (11 co-authors). 2012. The origin and evolution of the unique HCV circulating recombinant form 2k/1b. *J Virol.* 86:2212–2220.

Remsen J, DeSalle R. 1998. Character congruence of multiple data partitions and the origin of the Hawaiian Drosophilidae. *Mol Phylogenet Evol.* 9:225–235.

Remsen J, O'Grady P. 2002. Phylogeny of Drosophilinae (Diptera: Drosophilidae), with comments on combined analysis and character support. *Mol Phylogenet Evol.* 24:249–264.

Robe LJ, Loreto EL, Valente VL. 2010. Radiation of the "*Drosophila*" subgenus (Drosophilidae, Diptera) in the Neotropics. *J Zool Syst Evol Res.* 48:310–321.

Russo CA, Takezaki N, Nei M. 1995. Molecular phylogeny and divergence times of Drosophilid species. *Mol Biol Evol.* 12:391–404.

Schawaroch V. 2002. Phylogeny of a paradigm lineage: the *Drosophila melanogaster* species group (Diptera: Drosophilidae). *Biol J Linnean Soc.* 76:21–37.

Shapiro B, Rambaut A, Drummond AJ. 2006. Choosing appropriate substitution models for the phylogenetic analysis of protein-coding sequences. *Mol Biol Evol.* 23:7–9.

Sharp WD, Renne PR. 2005. The Ar-40/Ar-39 dating of core recovered by the Hawaii Scientific Drilling Project (phase 2), Hilo, Hawaii. *Geochem Geophys Geosyst.* 6:18.

Sharp WD, Turrin BD, Renne PR, Lanphere MA. 1996. The 40Ar/39Ar and K/Ar dating of lavas from the Hilo 1-km core hole, Hawaii Scientific Drilling Project. *J. Geophys Res.* 101:11607–11616.

Sherrod DR, Sinton JM, Watkins SE, Brunt KM. 2007. Geologic map of the state of Hawai'i. U.S. Geological Survey Open-File Report 2007–1089. Reston (VA): U.S. Geological Survey.

Stearns HT. 1946. Geology of the Hawaiian islands. Hawaii (Terr.) Division of Hydrography. Bulletin. 8:106.

Tamura K, Nei M. 1993. Estimation of the the number of nucleotide substitutions in the control region of mitochondrial DNA in humans and chimpanzees. *Mol Biol Evol.* 10:512–526.

Tamura K, Subramanian S, Kumar S. 2004. Temporal Patterns of Fruit Fly (*Drosophila*) Evolution Revealed by Mutation Clocks. *Mol Biol Evol.* 21:36–44.

Thomas JA, Welch JJ, Lanfear R, Bromham L. 2010. A Generation Time Effect on the Rate of Molecular Evolution in Invertebrates. *Mol Biol Evol.* 27:1173–1180.

Thomas RH, Hunt JA. 1993. Phylogenetic relationships in *Drosophila*—a conflict between molecualr and morphological data. *Mol Biol Evol.* 10:362–374.

Throckmorton LH. 1975. The phylogeny, ecology and geography of *Drosophila*. In: King RC, editor. Handbook of genetics. New York: Plenum. p. 421–469.

van der Linde K, Houle D, Spicer GS, Steppan SJ. 2010. A supermatrix-based molecular phylogeny of the family Drosophilidae. *Genet Res.* 92:25–38.

Vicario S, Moriyama EN, Powell JR. 2007. Codon usage in twelve species of *Drosophila*. *BMC Evol Biol.* 7:226.

Wiegmann BM, Trautwein MD, Winkler IS, et al. 2011. Episodic radiations in the fly tree of life. *Proc Natl Acad Sci USA.* 108:5690–5695.

Yang Y, Hou Z-C, Qian Y-H, Kang H, Zeng Q-T. 2012. Increasing the data size to accurately reconstruct the phylogenetic relationships between nine subgroups of the *Drosophila melanogaster* species group (Drosophilidae, Diptera). *Mol Phylogenet Evol.* 62:214–223.

Zuckerkandl E, Pauling L. 1965. Evolutionary divergence and convergence in proteins. In: Bryson V, Vogel HJ, editors. Evolving genes and proteins. New York: Academic Press. p. 97–166.