

## INVITED REVIEWS AND SYNTHESSES

# Reanalysis suggests that genomic islands of speciation are due to reduced diversity, not reduced gene flow

TAMI E. CRUICKSHANK\* and MATTHEW W. HAHN\*†

*\*Department of Biology, Indiana University, Bloomington, IN 47405, USA, †School of Informatics and Computing, Indiana University, Bloomington, IN 47405, USA*

## Abstract

The metaphor of ‘genomic islands of speciation’ was first used to describe heterogeneous differentiation among loci between the genomes of closely related species. The biological model proposed to explain these differences was that the regions showing high levels of differentiation were resistant to gene flow between species, while the remainder of the genome was being homogenized by gene flow and consequently showed lower levels of differentiation. However, the conditions under which such differentiation can occur at multiple unlinked loci are restrictive; additionally, essentially, all previous analyses have been carried out using relative measures of divergence, which can be misleading when regions with different levels of recombination are compared. Here, we test the model of differential gene flow by asking whether absolute divergence is also higher in the previously identified ‘islands’. Using five species pairs for which full sequence data are available, we find that absolute measures of divergence are not higher in genomic islands. Instead, in all cases examined, we find reduced diversity in these regions, a consequence of which is that relative measures of divergence are abnormally high. These data therefore do not support a model of differential gene flow among loci, although islands of relative divergence may represent loci involved in local adaptation. Simulations using the program IMA2 further suggest that inferences of any gene flow may be incorrect in many comparisons. We instead present an alternative explanation for heterogeneous patterns of differentiation, one in which postspeciation selection generates patterns consistent with multiple aspects of the data.

**Keywords:**  $F_{ST}$ , hitchhiking, IM, recombination, sympatric speciation

Received 11 December 2013; revision received 5 May 2014; accepted 7 May 2014

## Introduction

Understanding the genetic basis of speciation is a major goal of evolutionary biology. However, identifying the loci responsible for reproductive isolation using linkage mapping is a challenging task, and only a handful of genes have been directly implicated in either prezygotic isolation (e.g. Hopkins & Rausher 2011) or postzygotic isolation (e.g. Wittbrodt *et al.* 1989; Barbash *et al.* 2003; Presgraves *et al.* 2003). Because traditional genetic mapping is so laborious – and may miss both loci of small effect and loci contributing to phenotypes that are not identified as important to reproductive isolation –

researchers have turned to large-scale sequencing of natural populations in order to identify regions involved in speciation. By considering naturally hybridizing taxa, studies of DNA polymorphism and divergence make it possible to quickly reveal many aspects of the speciation process. Such studies hold the promise of identifying the number and location of regions involved in isolation, and possibly even the individual mutations underlying important differences between populations (e.g. Turner *et al.* 2010). More intriguingly, despite the seemingly stringent requirements for evolving new species in the face of gene flow (Felsenstein 1981), these studies could support widespread sympatric or parapatric speciation.

Accumulating analyses of recently diverged species pairs have found highly heterogeneous patterns of

Correspondence: Matthew W. Hahn, Fax: 812-855-6705; E-mail: mwh@indiana.edu

genetic differentiation across the genome (e.g. Harr 2006; Carneiro *et al.* 2010; Nadeau *et al.* 2012). Some regions show little genetic differentiation, while others show high levels of differentiation and may even contain fixed differences that distinguish species. Although, in most cases, the highly differentiated regions consist of only a small proportion of the genome (e.g. Turner *et al.* 2005), in others a large fraction of the genome may either be differentiated (e.g. Garrigan *et al.* 2012) or show some association with key isolating traits (e.g. Michel *et al.* 2010).

A common inference when a small number of differentiated regions are found by these genome scans is that they represent loci involved in reproductive isolation or ecological specialization and that there is strong selection against their introgression between species (Wu 2001). Loci not involved in isolation will experience the homogenizing effects of migration and will show little-to-no differentiation. In this scenario, small regions of elevated divergence within a genome are surrounded by regions that have low divergence, and this difference is thought to reflect differences in the effective amount of gene flow experienced by the two types of loci. Because genes within regions of elevated divergence may be involved in reproductive isolation between species, these regions have been referred to as 'genomic islands of speciation' (Turner *et al.* 2005). This phrase has largely been replaced with the slightly less loaded terms 'genomic islands of differentiation' (Harr 2006) and 'genomic islands of divergence' (Nosil *et al.* 2009), although the latter term can itself connote patterns of species differences not actually present in the data (see below). While there are many conceptual predecessors to the ideas represented by these terms (cf. Harrison 2012), they do succinctly capture both the configuration of species differences and the presumed underlying model.

Here, we consider multiple biological models that can explain the observed patterns of heterogeneous divergence, as well as the unique predictions made by each. Based on a reanalysis of published data sets, we find little evidence to support models in which islands of divergence contain loci resistant to introgression. We therefore subsequently re-examine the basis for claims that the species under consideration are exchanging genes. Despite multiple lines of evidence clearly supporting gene flow between many species pairs, we find high rates of false positives when applying the most widely used tests for gene flow to closely related species. Finally, we elaborate on an alternative model that can explain patterns of heterogeneous genomic differentiation in terms of divergence without differential gene flow among loci. As more whole-genome data sets are collected, it will become possible to determine

the evolutionary mechanisms that generate islands of differentiation.

### Models of species and sequence divergence

For our purposes, models of species divergence can be usefully separated into those with and without gene flow. For those scenarios including gene flow, we must further distinguish between models where there is no initial period of allopatry or reduced gene flow (i.e. sympatric speciation) and those where gene flow occurs after a substantial period of independent evolution between diverging taxa (i.e. secondary contact). Although both models are often lumped together as 'speciation-with-gene-flow', secondary contact models do not pose the same challenges for the evolution of isolating barriers as do models of sympatric speciation (cf. Felsenstein 1981). In this study, we refer to these as 'primary' and 'secondary' models of speciation-with-gene-flow and separately consider the predictions of each model as well as the data supporting those predictions.

There is strong evidence supporting both primary and secondary models of speciation-with-gene-flow in nature. Convincing examples of sympatric speciation have been identified, complete with heterogeneous differentiation among the limited number of loci examined (e.g. Savolainen *et al.* 2006). Lower levels of differentiation in sympatric vs. allopatric populations of separate species have provided good examples of gene flow between species initially separated by isolating barriers (e.g. Kulathinal *et al.* 2009; Martin *et al.* 2013). In addition, studies of hybrid zones formed after secondary contact have provided major insights into the resistance of loci to introgression across such zones (e.g. Payseur *et al.* 2004; Kronforst *et al.* 2006; Teeter *et al.* 2008; Maroja *et al.* 2009). These studies generally find a small number of differentiated loci and a large number of loci with evidence of introgression. Results from taxa experiencing both primary and secondary speciation-with-gene-flow suggest that we may be able to identify the reproductive isolating genes underlying 'islands of speciation'.

Assuming introgression between species, differentiation is expected at loci where the coefficient of selection is greater than the rate of migration (Haldane 1930; Wright 1931; Bulmer 1972). While this would appear to be promising for the identification of narrow windows of differentiation, 'divergence hitchhiking' hypothesizes that divergent natural selection causes reduced effective gene flow near a selected locus during primary differentiation (Via & West 2008; Via 2009, 2012; Feder & Nosil 2010); this idea is supported by theoretical work on the effects of migration, selection and linkage (e.g.

Charlesworth *et al.* 1997; Nordborg 1997; Akerman & Bürger 2014). Overtime, larger regions are 'captured' by the selected locus, and therefore, larger regions begin to show high levels of differentiation. Divergence hitchhiking is expected to mainly affect small populations when migration is low (Feder & Nosil 2010), and it predicts that differentiation should decline with distance from the selected locus. As the speciation process proceeds, this model proposes that a progressively larger proportion of the genome will become differentiated, finally resulting in two completely distinct genomes (Feder *et al.* 2012). A similar model applies to cases of secondary contact, where neutral alleles with no effect on fitness but that are linked to selected alleles will also become, or remain, differentiated across a hybrid zone (Barton & Bengtsson 1986). Overall, for models of both primary and secondary speciation-with-gene-flow, patterns of heterogeneous differentiation have been interpreted as revealing loci resistant to introgression, with the size of differentiated regions correlated with the strength of selection and the stage of species divergence.

Heterogeneous patterns of genomic differentiation can also be explained by an alternative model in which there is no gene flow between recently diverged species (Noor & Bennett 2009; Turner & Hahn 2010; Hahn *et al.* 2012). Under this model, reproductive isolation is instantaneous and complete, although similar models do not require the complete absence of introgression. Species pairs show low overall levels of differentiation because they have only recently split, and shared alleles are due to ancestrally inherited variation, not introgression. Heterogeneity in the level of differentiation among loci in this model is due to stochastic variation in coalescent times (Barton 2006) and, more notably, natural selection: those loci experiencing strong selection will appear to be more differentiated and will share less ancestral variation (see next section). The targets of natural selection may be directly involved in the ecological, morphological or behavioural specialization of each species, or they may be unrelated to any trait involved in species divergence and simply represent the 'background' level of selection found in any organism. Most importantly, however, in this model, there is no differential gene flow among loci, and the regions of highest differentiation do not necessarily indicate the location of genes underlying isolating traits ('incidental islands': Turner & Hahn 2010; White *et al.* 2010).

Determining the contribution of these models – two with gene flow helping to outline the boundaries of regions involved in speciation and one with selection defining targets of adaptation that are not necessarily involved in the speciation process – has important consequences for the interpretation of the growing number

of genome scans conducted across taxa. Due to the ease with which patterns of molecular variation can be queried, scans for outlier loci (i.e. those loci with higher-than-expected differentiation) continue to be conducted in a wide variety of species, across different timescales of species divergence (e.g. Nadeau *et al.* 2013) and across varying geographic relationships (e.g. Renaut *et al.* 2013). Almost without exception, studies identifying outlier markers have concluded that they represent loci underlying isolating barriers or locally adapted alleles and that they are resistant to gene flow that homogenizes the majority of the genome. The conclusion reached by many researchers has therefore been that speciation-with-gene-flow (in this case, sympatric speciation) is not only possible, but potentially widespread (e.g. Nosil 2008; Via 2012). A growing body of theoretical work has arisen based on the idea that islands of differentiation are resistant to introgression (e.g. Feder & Nosil 2010; Flaxman *et al.* 2012, 2013), and wholly constructed histories of the speciation process have been proposed based on the same interpretation of the data (Feder *et al.* 2012). The apparent widespread acceptance of sympatric speciation based on patterns of molecular differentiation argues strongly for the close examination of the various evolutionary forces that could produce such patterns.

In the next section, we consider a further prediction of both speciation-with-gene-flow models, namely that absolute measures of species divergence (as opposed to relative measures, see next section) should also be higher in regions resistant to gene flow. Surprisingly, absolute measures have either not been considered by most researchers (but see Noor & Bennett 2009; Nachman & Payseur 2012), or relative measures were mistakenly thought to be absolute (a mistake one of us has made; Turner *et al.* 2005). We believe that comparisons among loci using absolute measures of divergence will be especially revealing with respect to the contribution of each mode of divergence between species.

## Interpreting measures of sequence divergence

### *Relative measures of divergence*

There are a number of widely used statistics that enable researchers to quantify the genetic distance between species, either at individual sites or averaged over multiple sites. One common class of measures is based on Wright's  $F_{ST}$  (Wright 1931, 1943), which measures the normalized difference in allele frequencies between populations. There are many different ways of calculating  $F_{ST}$  and its related statistics (Charlesworth 1998; Excoffier 2007), and many different ways of interpreting these statistics (Holsinger & Weir 2009). Regardless of

the specific way in which  $F_{ST}$  is calculated, the key aspect that all  $F_{ST}$ -like measures share in common is that their values are inflated when diversity within populations is low: loci with lower levels of within-species variation are expected to show higher values of  $F_{ST}$  (Nei 1973; Charlesworth *et al.* 1997; Charlesworth 1998; Jakobsson *et al.* 2013). Because of this dependence on levels of within-species polymorphism,  $F_{ST}$  is referred to as a 'relative' measure of divergence. Additional relative measures of divergence ( $d_a$  and  $d_f$ ) are described in Box 1.

#### Box 1. Measures of sequence divergence

There are multiple different ways to measure sequence divergence when more than just a single individual is sampled from each population or species. One of the most widely used statistics is  $F_{ST}$ , a normalized measure of allele frequency differences between populations (Wright 1931, 1943). There are many different ways to calculate  $F_{ST}$  and  $F_{ST}$ -like statistics (Charlesworth 1998; Excoffier 2007), although many of the subtle distinctions between them are not relevant here. A common way to calculate  $F_{ST}$  is:

$$F_{ST} = \frac{\pi_T - \pi_S}{\pi_T}$$

where  $\pi_T$  is the expected heterozygosity in the total sample and  $\pi_S$  is the average expected heterozygosity in each population. [Note that this statistic is sometimes called  $\gamma_{ST}$  (Nei 1982), and for single sites is equivalent to  $G_{ST}$  (Nei 1973).] One very important aspect of all  $F_{ST}$ -like statistics is that they are strongly influenced by within-population levels of variation – represented by  $\pi_S$  in the above formula – and because of this, we refer to them as *relative* measures of differentiation.

To get around the dependence on within-population variation, Nei & Li (1979) calculated the average number of pairwise differences between sequences from two populations, *excluding* all comparisons between sequences within populations. Here, we refer to this statistic as  $d_{XY}$ , although it is also referred to in the literature as  $\pi_{XY}$  (Nei & Li 1979),  $D_{XY}$  (Nei 1987) and  $\pi_B$  (Charlesworth 1998). This *absolute* measure of differentiation is independent of the levels of diversity within the two populations being compared and is calculated as:

$$d_{XY} = \sum_{ij} x_i y_j d_{ij}$$

where, in two populations, X and Y,  $d_{ij}$  measures the number of nucleotide differences between the  $i^{\text{th}}$

haplotype from X and the  $j^{\text{th}}$  haplotype from Y (Nei & Li 1979; Nei 1987, equation 10.20). The statistic  $d_{XY}$  is not affected by current levels of within-population diversity, although it is affected by both ancestral levels of diversity (see below) and the substitution rate. For instance, coding regions will often show lower levels of absolute divergence simply because there are fewer possible neutral mutations; conversely, at loci with a large number of substitutions fixed by positive selection,  $d_{XY}$  may actually be higher.

There are also a number of relative measures of differentiation based on  $d_{XY}$  that are in wide usage. Nei & Li (1979) defined a measure of the 'net' nucleotide differences between two populations,  $d_a$  (which they called  $\delta$ ), as:

$$d_a = d_{XY} - (\pi_X + \pi_Y)/2.$$

This statistic (also commonly called  $D_a$  or  $D_m$ ) is intended to capture only the differences that have accumulated between populations since they split (Fig. 1). It does so by subtracting out the differences that had accumulated before this split, assuming that the level of ancestral variation was equal to the average of the variation found in the two current populations.  $d_a$  is therefore a relative measure of divergence because it can be strongly affected by the amount of within-population variation. The number of fixed differences between populations or species,  $d_f$ , is a complex function of the time since divergence between populations and the time to coalescence within each population (Hey 1991) but is also a relative measure of divergence due to its reliance on levels of within-population variation.

Finally, note that  $d_{XY}$  and  $d_a$  are more informative for full sequence data; that is, for data in which at least moderate length sequences have been collected from homologous loci in two species.  $d_{XY}$  can be calculated for a single site, but it no longer has the property of being an absolute measure and offers little precision. This is most easily demonstrated by considering the expectation and variance for  $d_{XY}$  under a strictly allopatric model (Nei 1987, equation 13.83):

$$E(d_{XY}) = 2\mu t + \theta_{Anc}$$

$$Var(d_{XY}) = \frac{(2\mu t + \theta_{Anc})}{L} + \theta_{Anc},$$

where  $\mu$  is the neutral mutation rate,  $t$  is the time since the species split,  $\theta_{Anc}$  is the level of diversity in the ancestral (presplit) population, and  $L$  is the length of locus being considered. We can see that the



variance in  $d_{XY}$  goes down with increasing sequence length, such that we can have more confidence in estimates that come from longer loci. This means that we can calculate  $d_{XY}$  from reduced representation methods that produce sequences of hundreds of nucleotides, such as RAD-seq (Baird *et al.* 2008). However, these expressions also indicate that for estimates of  $d_{XY}$  from single SNPs, the variance is necessarily greater than the expectation, so that it will be difficult to differentiate values of  $d_{XY}$  among loci. All of this implies that for many kinds of markers – SNPs, AFLPs, microsatellites – only relative measures of divergence (i.e. only  $F_{ST}$ ) will be informative, and therefore, that no comparison between relative and absolute measures can be made for these markers. Unfortunately, SNPs, AFLPs and microsatellites are still the most common types of markers used in studies of between-species divergence. When full sequence data are collected, relative measures including  $F_{ST}$ ,  $d_a$  and  $d_f$  can be compared to  $d_{XY}$ .

Because relative measures of divergence rely on within-species diversity, markers sampled from parts of the genome with more or less diversity will provide very different views on levels of differentiation. For instance,  $F_{ST}$  from regions of reduced recombination – which often have reduced diversity due to linked selection – will be higher than  $F_{ST}$  from regions of normal recombination for no other reason than that total levels of nucleotide diversity are different (Charlesworth *et al.* 1997; Charlesworth 1998; Noor & Bennett 2009; Nachman & Payseur 2012). As a consequence, Charlesworth (1998) recommended that relative measures of differentiation ‘are not necessarily appropriate if we wish to compare loci with very different levels of within-population variation’. Differences in levels of variation are not always associated with rates of recombination, although this is a common correlate (Hahn 2008; Cutter & Payseur 2013).

With respect to the interpretation of genomic islands, these considerations suggest that regions may differ in values of  $F_{ST}$  (or other relative measures of divergence) not because of differences in the amount of gene flow, but simply because of differences in levels of within-species polymorphism. Extreme reductions in within-species polymorphism at specific loci can be due to selection directly on adaptive traits encoded by such regions, or due to linked selection on either advantageous (‘hitchhiking’; Maynard Smith & Haigh 1974) or deleterious mutations (‘background selection’; Charlesworth *et al.* 1993). Both background selection and hitchhiking reduce within-population diversity and lead

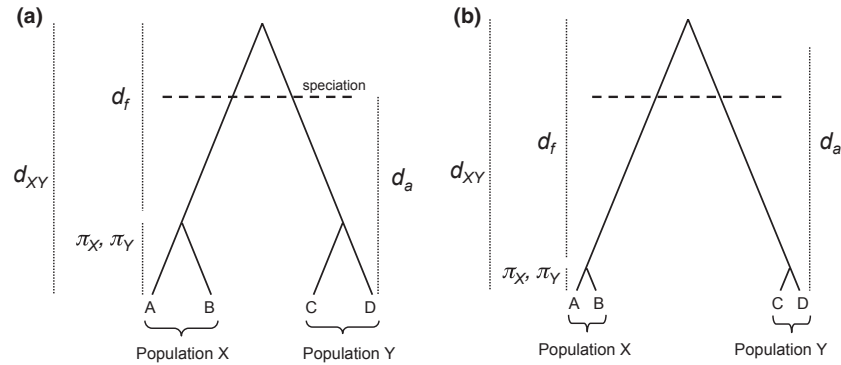
to higher between-population differentiation (Kaplan *et al.* 1989; Nordborg *et al.* 1996; Slatkin & Wiehe 1998). As a consequence, variation in the effects of selection along a chromosome can produce islands of relative divergence even in the complete absence of gene flow. Accordingly, presumed speciation islands are often found in centromeric or rearranged regions where recombination is reduced (e.g. Turner *et al.* 2005; Kulathinal *et al.* 2009). Because higher levels of relative divergence in regions of low recombination are a prediction of both speciation-with-gene-flow models and models without gene flow (Nachman & Payseur 2012), neither the presence nor the location of islands alone can support one model over the other. Instead, we can use predictions about the behaviour of alternative measures of species divergence.

### Absolute measures of divergence

To circumvent the dependence of relative measures of divergence on within-species polymorphism, we can calculate the average number of pairwise differences between sequences from two species, excluding all comparisons between sequences within species (Fig. 1; Box 1); we refer to this statistic as  $d_{XY}$ . It does not matter exactly how many chromosomes are sampled in each species, as all chromosomes from one species will show a similar level of divergence to any chromosome sampled from the other species. In fact, even if only a single sequence is sampled from each species – as is commonly done when calculating species divergence using the statistic  $d$  – this is still an estimate of  $d_{XY}$  (Nei 1987). Because  $d_{XY}$  is independent of the levels of diversity within the two populations or species being compared, we refer to it as an ‘absolute’ measure of divergence.

Absolute measures of divergence capture the number of sequence differences at a locus since the most recent common ancestor (MRCA) of the sampled chromosomes (Fig. 1). Because the MRCA almost certainly existed in the ancestral population before the split into the two descendent species,  $d_{XY}$  measures both the differences that have accumulated since the species split and the differences that were present at the time of the split (Gillespie & Langley 1979). As such, variation in  $d_{XY}$  can be due to variation in levels of ancestral polymorphism, but is not affected by variation in levels of current polymorphism. Consequently, this means that linked selection in extant populations will not cause the appearance of ‘islands’ of high divergence in  $d_{XY}$ , although linked selection in ancestral populations can actually lower the value of  $d_{XY}$  (Begun *et al.* 2007).

To see how absolute and relative measures can give dissimilar results in the presence of linked selection,



**Fig. 1** Relative and absolute measures of divergence, with the effect of linked selection. Demonstrating the differences between  $\pi_X$ ,  $\pi_Y$ ,  $d_a$ ,  $d_f$  and  $d_{XY}$ . Panels (a) and (b) both show example genealogies relating four sampled chromosomes (A, B, C and D) from two populations or species (X and Y). The statistics  $\pi_X$  and  $\pi_Y$  measure the average number of nucleotide differences between samples in population X and each sample in population Y, with no comparisons made within a population.  $d_a$  uses the average current levels of polymorphism as a measure of ancestral polymorphism, and subtracts this value from the total divergence ( $= d_{XY} - [\pi_X + \pi_Y]/2$ ).  $d_f$  represents the total number of fixed differences between the two populations. For ease of comparison, the genealogies in the two panels have the same height, as do the genealogies for populations X and Y within each figure. The important distinction between panels a) and b) is that due to linked selection, there is a difference in  $\pi_X$ ,  $\pi_Y$ ,  $d_a$  and  $d_f$  between the two panels, but no difference in  $d_{XY}$ .

consider the two genealogies shown in Fig. 1. Figure 1a shows one hypothetical genealogical history of two species for a locus in a region that has not experienced any selection. Levels of within-species diversity are unaffected by linked selection and are therefore high relative to the between-species portion of their history. In this case,  $\pi$  (a measure of within-population diversity) is high in both species, and  $d_a$  and  $d_f$ , two measures of between-species relative divergence similar to  $F_{ST}$  (Box 1; Fig. 1), are commensurately low because they are dependent on within-species variation. Figure 1b, on the other hand, shows a hypothetical genealogical history for a locus affected by linked selection in both species (either hitchhiking or background selection), sampled from exactly the same individuals as in Fig. 1a. The time to the most recent common ancestor within populations in this scenario is extremely recent due to linked selection. This means that  $\pi$  in each species is low and that  $d_a$  and  $d_f$  are commensurately high ( $F_{ST}$  would show the same pattern). Most importantly,  $d_{XY}$  is exactly the same in the two panels. In this way, relative and absolute measures can provide very different pictures about the level of between-species divergence from locus to locus, simply due to linked selection.

#### The effect of gene flow on measures of divergence

To demonstrate the clear predictions of models with gene flow on both relative and absolute divergence at loci that are resistant to gene flow (possibly because alternative alleles confer a locally adapted phenotype),

we take advantage of the theory laid out in Charlesworth *et al.* (1997). This mathematical framework allows us to make predictions not only about divergence at the selected site, but also at linked neutral alleles for models in which there is gene flow between species, via either primary divergence or secondary contact. Surprisingly, even though the theoretical machinery presented in Charlesworth *et al.* (1997) has been used in multiple studies to elaborate on patterns of relative divergence expected at equilibrium (e.g. Via & West 2008; Nosil *et al.* 2009; Feder & Nosil 2010), to our knowledge it has not been used to clarify predictions about absolute divergence.

It is relatively easy to derive expectations for  $d_{XY}$  under this model directly from the equations presented in Charlesworth *et al.* (1997). The within- and among-population components of sequence diversity at neutral polymorphisms linked to a locus under selection in a subdivided population are summarized in Table 1 of Charlesworth *et al.* (1997) and in Nordborg (1997). We make a number of the same simplifying assumptions as in the original papers (including large equally sized demes with symmetric migration) and sum the components [*i.e.* the variation within populations (Nordborg 1997; equation 65) and the expected neutral site differentiation among populations (Charlesworth *et al.* 1997; Table 1)] to find that:

$$E(d_{XY}) = (1 - q)\alpha + \frac{r + 4\tilde{m}}{8N_e\tilde{m}r}$$

Here,  $q$  is the frequency of a locally deleterious allele, and variation is reduced due to background selection in

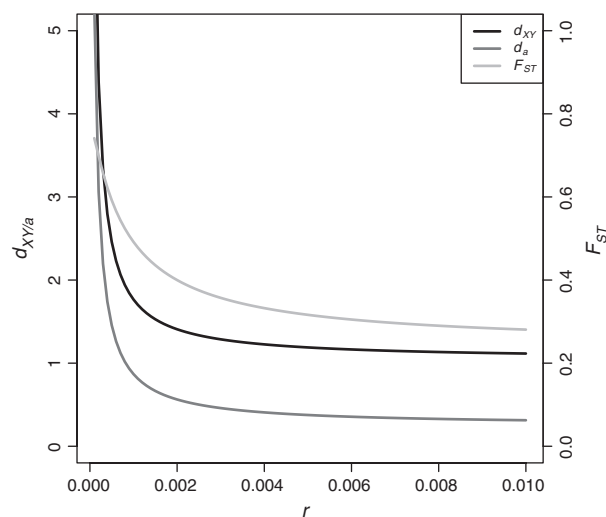
**Table 1** Nucleotide divergence in islands and nonislands for five pairs of recently diverged taxa

		Islands	Nonislands		Bin type	References for sequence data
<i>Oryctolagus cuniculus</i>	$d_{XY}$	0.0368	0.0409	$N = 27,17; P = 0.40$	X/A	Carneiro <i>et al.</i> (2009, 2010);
<i>cuniculus</i> and <i>O. c. algirus</i>	$F_{ST}$	0.4401	0.1456	$P = 0.0002$		Geraldes <i>et al.</i> (2006)
	$d_{XY}$	0.0358	0.0411	$N = 23,21; P = 0.91$	Migration*	
	$F_{ST}$	0.615	0.103	$P < 0.0001$		
<i>Mus musculus musculus</i> and	$d_{XY}$	0.007	0.0079	$N = 8,19; P = 0.75$	High/low $r$	Geraldes <i>et al.</i> (2011)
<i>M. m. domesticus</i>	$F_{ST}$	0.7762	0.6931	$P = 0.21$		
	$d_{XY}$	0.0105	0.008	$N = 8,6; P = 0.47$	High/low $F_{ST}$	Harr (2006)
	$F_{ST}$	0.6387	0.2417	$P = 0.0002$		
<i>Heliconius melpomene aglaope</i>	$d_{XY}$	0.0145	0.0155 (B/D region)	5 islands; $P = 0.80$	High/low $F_{ST}$	Nadeau <i>et al.</i> (2012)
and <i>H. m. amaryllis</i>	$F_{ST}$	0.4284	0.184	$P < 0.0001$		
	$d_{XY}$	0.0197	0.0140 (Yb region)	7 islands; $P = 0.87$	High/low $F_{ST}$	
	$F_{ST}$	0.3422	0.132	$P < 0.0001$		
<i>Anopheles coluzzii</i> (M form)	$d_{XY}$	0.006	0.0096	$N = 17,12; P = 0.17$	high/low $F_{ST}$	Turner <i>et al.</i> (2005);
and <i>A. gambiae</i> (S form)	$F_{ST}$	0.5242	0.2138	$P = 0.008$		Turner & Hahn (2007);
						White <i>et al.</i> (2010)
<i>Ficedula albicollis</i> and	$d_{XY}$	0.0036	0.0044	~50 islands; N/A	High/low $F_{ST}$	Ellegren <i>et al.</i> (2012)
<i>F. hypoleuca</i>	$F_{ST}$	0.742	0.357	N/A		

\*Based on results Sousa *et al.* (2013).

proportion to  $\alpha$  (see equation 24 in Nordborg 1997). The effective migration rate,  $\hat{m}$ , accounts for differences in fitness between migrants and nonmigrants (see equation A3 in Charlesworth *et al.* 1997). This expression for  $d_{XY}$  is equivalent to the average pairwise divergence between chromosomes sampled from two subpopulations. Expectations for  $d_a$  in this scenario can be easily derived from the expectation for  $d_{XY}$ . Equations for calculating all statistics from data can be found in Box 1, as can the expectation for  $d_{XY}$  under a purely allopatric model.

As Fig. 2 shows, under this model at equilibrium  $F_{ST}$ ,  $d_a$  and  $d_{XY}$  are expected to be high at loci that do not introgress and to decline sharply as one moves away from the selected locus. Gene flow reduces both absolute and relative measures of species divergence, leading to a clear difference in  $d_{XY}$  at loci that experience different levels of effective gene flow. Notably, the magnitude of differences in both  $F_{ST}$  and  $d_{XY}$  between introgressing and nonintrogressing regions is of the same scale when there is differential gene flow, so that loci with relative divergence that is ten times higher than background are also expected to have absolute divergence that is this much higher. The slight difference observed in Fig. 2 in the rate at which  $F_{ST}$  and  $d_{XY}$  drop off with distance from the selected site is not consistent across all combinations of parameters and cannot be expected to be a distinguishing characteristic of these two measures. When populations are not at equilibrium – such as shortly after speciation – expectations for  $F_{ST}$  and  $d_{XY}$  are more complex. Box 2 contains simulation



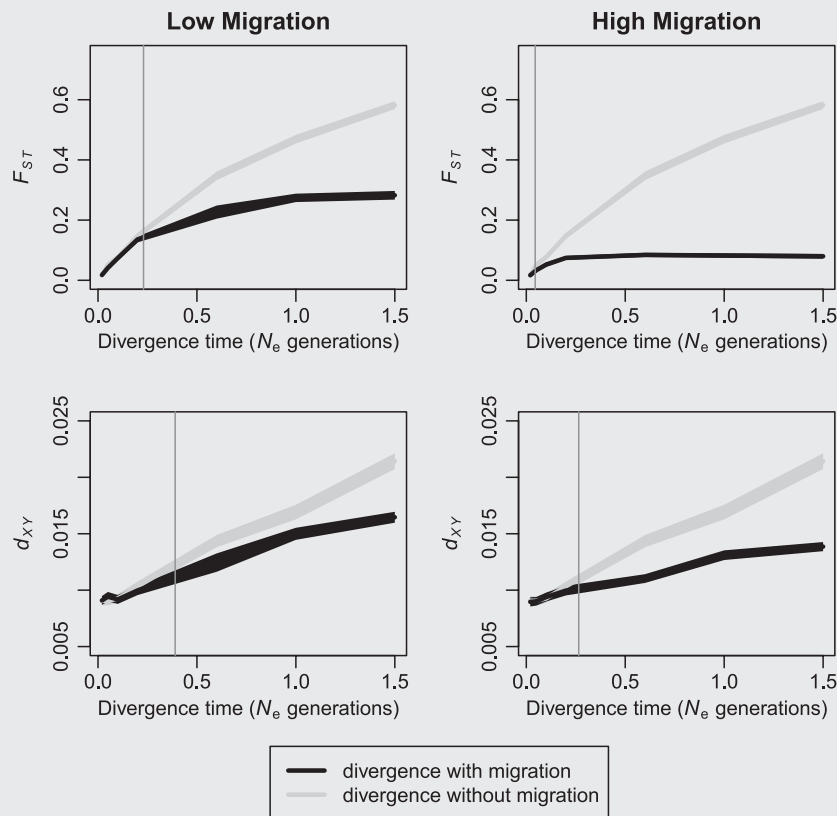
**Fig. 2** Analytical predictions for divergence with gene flow. Analytical predictions for absolute divergence ( $d_{XY}$ ) and relative divergence ( $F_{ST}$  and  $d_a$ ) for neutral sites linked to a selected site in a subdivided population. The selected site lies all the way at the left end of the x-axis, with increasing recombination distance ( $r$ ) moving away from this site. The populations are assumed to have a total size of  $N_e = 5000$ , with migration among demes occurring at rate  $m = 0.0001$ . The left-hand y-axis measures the values of  $d_{XY}$  and  $d_a$ , while the right-hand y-axis measures the values of  $F_{ST}$ . All measures are expected to decrease quickly with distance from a selected locus ( $s = 0.01$ ).

results showing the differences in these statistics with and without gene flow when loci are sampled shortly after a species split.

**Box 2.** Statistical power to detect differences in effective migration among loci.

An important assumption of the comparisons made in the main text is that we have the same statistical power to detect regions with elevated relative and absolute divergence. While Fig. 2 suggests that this is true for populations that have reached an equilibrium between migration, drift and selection, this will not necessarily be the case shortly after a speciation event because populations will not be at equilibrium. Lacking theoretical expectations for non-equilibrium conditions of primary speciation-with-gene-flow, Fig. B1 shows the results for different measures of divergence from simulated populations shortly after speciation (where the species split happened instantaneously at time 0). In each simulation, we measured relative and absolute divergence at loci with and without migration between the two species. The results from the simulations lead to several important conclusions about our ability to detect differences in the effective migration rate among loci.

First, in all conditions, those loci not experiencing gene flow between species diverge faster than those loci that are being homogenized by gene flow, as expected. However, in all cases, we also lack power to distinguish between these two classes of loci for a period of time after speciation (Fig. B1; the vertical line in each panel marks the position where 95% confidence intervals between these classes do not overlap). In the best-case scenario (using  $F_{ST}$  when there is high background migration), we are unable to identify those loci that are not introgressing between species until  $\sim 0.1 N_e$  generations after the speciation event. This can be a substantial amount of time, even for species with relatively small effective population sizes.



**Fig. B1** Divergence, estimated by the statistics  $F_{ST}$  (top) and  $d_{XY}$  (bottom), is plotted as a function of time as two populations split (measured in  $N_e$  generations). Loci are sampled from populations simulated with zero ( $N_e m = 0$ ) and nonzero migration rates ( $N_e m = 1$  on the left, or  $N_e m = 5$  on the right). Within each condition and for each measure, divergence with and without migration are indicated by grey and black lines, respectively. The width of each line represents the 95% confidence interval for 500 simulated data sets, and the vertical line within each panel shows the point in time at which there is no overlap between the confidence intervals with and without migration. At this point, loci with and without migration would be expected to show significantly different values of the relevant statistic.



Second, and as noted in the main text, in non-equilibrium populations, there is lower power to detect differences among loci using  $d_{XY}$  than there is using  $F_{ST}$  for very short divergence times (Fig. B1; compare top and bottom panels in each column). This means that  $F_{ST}$  is more sensitive to a lack of gene flow than is  $d_{XY}$  and will begin to differ between loci with and without migration more quickly after speciation. The lag is due to the fact that mutations must arise in order for  $d_{XY}$  to increase, while  $F_{ST}$  only requires changes in allele frequencies. If there are loci that do show an increase in  $F_{ST}$  due to a lack of gene flow, these results imply that they will not necessarily also show a statistically significant difference in  $d_{XY}$ . Note, however, that on average, these loci will show a higher value of  $d_{XY}$  (significant or not) and are not expected to have significantly lower levels of polymorphism (cf. Tables 1 and 2).

Finally, the simulations also demonstrate that there is more power to distinguish between loci with and without migration when background migration rates are high (Fig. B1; compare left and right panels in each row). That is, differences between these classes of loci can be detected much more quickly after the species split with high gene flow. In addition, the difference in the power of  $F_{ST}$  and  $d_{XY}$  is magnified under conditions of high migration, with a greater time period expected in which the relative measure ( $F_{ST}$ ) is significant and the absolute measure ( $d_{XY}$ ) is not.

In the previous section, we demonstrated that there are different expectations for the levels of relative and absolute measures of divergence in a strictly allopatric model with linked selection. In contrast, models of primary and secondary speciation-with-gene-flow make similar predictions about all measures of divergence (see Table 1 in Nachman & Payseur 2012). If islands represent regions resistant to introgression that are surrounded by loci homogenized by gene flow because they are not involved in species isolation, then regions of high absolute divergence should co-occur with regions of high relative divergence. The convergence of measures that are and are not affected by within-species diversity is a clear prediction of the speciation-with-gene-flow models, one that eliminates the possibility that reduced intraspecific variation has created peaks of divergence independent of the speciation process.

## Re-analysis of published sequence data

### *Patterns of sequence divergence do not support differential introgression among loci*

The fact that linked selection affects relative measures of divergence but does not affect absolute measures of divergence – and that the expectations are the same for all measures under a model with differential gene flow across loci – suggests a straightforward way to distinguish which model has had the greatest contribution to species divergence. With few exceptions (Noor & Bennett 2009; Nachman & Payseur 2012), this difference has not been examined in the many studies assessing levels of differentiation between closely related species. In part this is due to the fact that absolute divergence is not informative for many kinds of markers (see Box 1), and may also simply be due to the growing acceptance of speciation-with-gene-flow models. Here, we evaluate the

two contrasting explanations for islands of differentiation by examining both relative and absolute measures of divergence in five systems that have become widely cited as examples of speciation-with-gene-flow (Table 1). Three of the species pairs are clear examples of gene flow after secondary contact (*Ficedula*, *Mus*, *Oryctolagus*), while two are hypothesized to be examples of primary speciation-with-gene-flow (*Anopheles*, *Heliconius*). Although recent studies have also reviewed the evidence for islands using absolute divergence in a subset of these species (Noor & Bennett 2009; Nachman & Payseur 2012), here, we have used additional data and additional methods for identifying islands.

The data sets vary in the technologies used for sequencing and in the proportion of the genome covered. For one system (*Ficedula*), a whole-genome data set is available, with multiple individual birds sequenced separately from each species (Ellegren *et al.* 2012). For another system (*Heliconius*), the sequences come from a few large regions of the genome that were obtained by sequence-capture technology (Nadeau *et al.* 2012). Sequences from moderate numbers of independently sequenced loci were available for *Anopheles* (Turner *et al.* 2005; Turner & Hahn 2007; White *et al.* 2010), *Mus* (Harr 2006; Geraldès *et al.* 2011) and *Oryctolagus* (Geraldès *et al.* 2006; Carneiro *et al.* 2009, 2010).

In each data set considered, previous analyses have identified islands of differentiation (or simply distinct loci that were highly differentiated) based on relative measures of divergence – either  $F_{ST}$ ,  $d_{ar}$ ,  $d_f$ , or all three together – even if absolute measures of divergence were also used (e.g. Carneiro *et al.* 2010; Geraldès *et al.* 2011; Ellegren *et al.* 2012). Using a definition of ‘islands’ based on the boundaries or divisions identified in the original studies, we calculated absolute divergence ( $d_{XY}$ ) for the same regions. In none of the five species pairs was absolute divergence significantly higher in

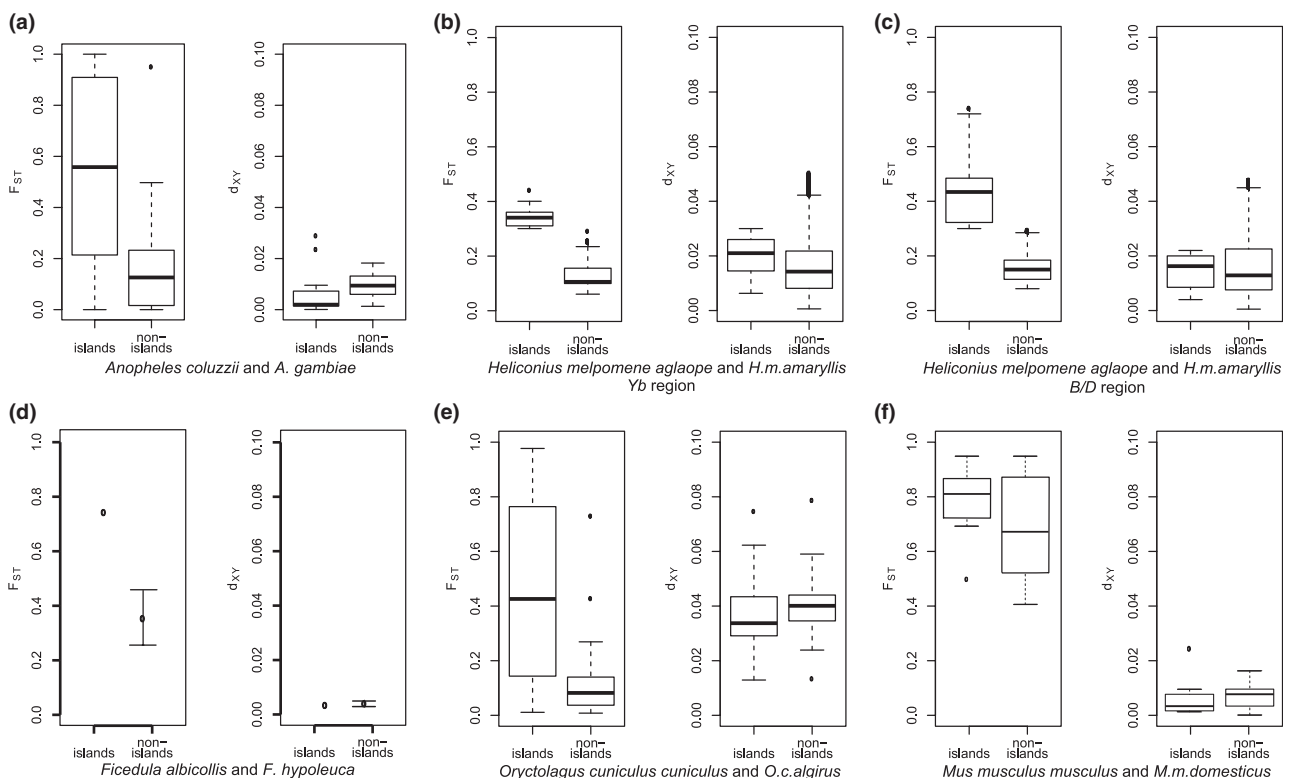
the islands of relative divergence than it was in matched flanking regions or the genome as a whole (Table 1). Despite the many-fold higher values of relative divergence reported for these regions (Table 1), in four of the five studies absolute divergence was in fact lower in the islands than in nonislands, a result that is never predicted by models that interpret the islands simply as regions resistant to introgression between hybridizing species. On the contrary, lower values of  $d_{XY}$  likely indicate that the regions identified as islands have been experiencing recurrent bouts of selection since before the species split (see below). The fact that many of these islands of relative divergence also coincide with regions of low recombination only serves to emphasize this expected pattern of low absolute divergence (see below, and Begun *et al.* 2007).

While a global view of sequence divergence is informative about the contributions of models of speciation (or at least models of introgression among loci), examining individual species pairs may offer a more nuanced understanding of what drives species divergence. Below, we discuss what we consider to be key features – or important caveats – in interpreting patterns of sequence evolution in each of the systems re-analysed here.

### *Anopheles (mosquitoes)*

Patterns of differentiation between the 'M' and 'S' form mosquitoes in the species *Anopheles gambiae* (since renamed *A. coluzzii* and *A. gambiae*, respectively; Coetzee *et al.* 2013) were some of the first to be examined genomewide (Wang *et al.* 2001; Stump *et al.* 2005; Turner *et al.* 2005). Based on relative measures of divergence, at least three islands of divergence are consistently found across populations, and these coincide with pericentromeric regions (Stump *et al.* 2005; Turner *et al.* 2005; Turner & Hahn 2007; Neafsey *et al.* 2010; White *et al.* 2010).

Here, we have focused on reanalyzing the data from loci sequenced in three earlier studies (Table 1; Turner *et al.* 2005; Turner & Hahn 2007; White *et al.* 2010). [Upon reanalysis of whole-genome data from *Anopheles* (Lawniczak *et al.* 2010), we did not find patterns of relative sequence divergence supporting the previously identified islands and therefore did not use these data.] As shown in Fig. 3a, while  $F_{ST}$  in the islands is higher than nonislands – as expected by the definition of islands – the values of  $d_{XY}$  are not. Absolute divergence is lower in the pericentromeric islands (Table 1), as is



**Fig. 3** Reanalysis of published data sets. Box plots are drawn from two measures of divergence –  $F_{ST}$  (left) and  $d_{XY}$  (right) – for genes sampled from pairs of recently diverged taxa. For each measure and each species pair, a comparison is made between islands and nonislands. Islands are regions of high  $F_{ST}$  and are defined differently for each system (see text for details). For panel d), only a single genome-wide measure was available for islands. Central horizontal lines represent medians.

expected in regions of low recombination that experience recurrent selection. As shown here (Table 2) and reported earlier from SNP-chip data (Neafsey *et al.* 2010), there is a strong reduction in polymorphism in islands of relative divergence. This decrease in within-species polymorphism is expected under models of linked selection and will cause increased values of relative divergence measures unrelated to levels of effective gene flow.

Conclusions about the level of ongoing gene flow between M and S are mixed, with clear evidence of hybridization (though at very low levels; Tripet *et al.* 2001), including the introgression of highly advantageous insecticide resistance alleles (Djogbénou *et al.* 2008; Etang *et al.* 2009), but also a large deficit of recombinant genotypes (White *et al.* 2010). Whether M and S represent largely nonhybridizing species at an advanced stage of reproductive isolation (Lawniczak *et al.* 2010; Turner & Hahn 2010; White *et al.* 2010; Hahn *et al.* 2012) or are still exchanging genes across most of the genome (Reidenbach *et al.* 2012; Lee *et al.* 2013), these results do not support a model in which islands of differentiation represent loci that have been resistant to introgression over long periods of time.

### *Heliconius* (butterflies)

Butterflies in the genus *Heliconius* have undergone a recent radiation that has created a large number of closely related species and races (Jiggins 2008). *Heliconius* butterflies have conspicuous colour patterns and behaviours that differ between species, but previous phylogenetic analyses have revealed the adaptive introgression of colour-pattern genes even between nonsister species within the genus (Pardo-Diaz *et al.* 2012; The *Heliconius* Genome Consortium 2012). Peaks of relative differentiation between two races in the species *H. melpomene*, *H. m. amaryllis* and *H. m. aglaope* were found in two genomic regions that contain these colour-pattern genes (Nadeau *et al.* 2012, 2013). The low values of  $F_{ST}$  in flanking regions surrounding these loci and across the genome have been interpreted as evidence for gene flow, and the small size of selected regions indicates particularly high gene flow according to the divergence hitchhiking hypothesis (Nadeau *et al.* 2012).

Using targeted resequencing data from two known colour-pattern loci (Nadeau *et al.* 2012), we calculated  $d_{XY}$  across both regions. As can be seen in Fig. 3b,  $c$ ,  $d_{XY}$  is not higher within the peaks of  $F_{ST}$  relative to

**Table 2** Nucleotide diversity in islands and nonislands for five pairs of recently diverged taxa

	$\pi$ islands	$\pi$ nonislands		Bin type
<i>Oryctolagus c. cuniculus</i>	0.004238	0.006328	$N = 27,17$	X/A
<i>O. c. algirus</i>	0.003772	0.006171	$P = 0.039$	
			$P = 0.013$	
<i>O. c. cuniculus</i>	0.003459	0.006436	$N = 23,21$	Migration*
<i>O. c. algirus</i>	0.002993	0.006404	$P = 0.0018$	
			$P = 0.0002$	
<i>Mus m. musculus</i>	0.00087	0.00199	$N = 8,19;$	High/low $r$
<i>M. m. domesticus</i>	0.00139	0.00244	$P = 0.025$	
			$P = 0.16$	
<i>M. m. musculus</i>	0.000625	0.002916	$N = 8,6;$	High/low $F_{ST}$
<i>M. m. domesticus</i>	0.0015	0.003867	$P = 0.022$	
			$P = 0.04$	
<i>Heliconius m. aglaope</i> (Yb)			5 islands;	High/low $F_{ST}$
<i>H. m. amaryllis</i> (Yb)			$P = \text{n.s.}^\dagger$	
			$P < 0.05^\dagger$	
<i>H. m. aglaope</i> (B/D)			7 islands;	High/low $F_{ST}$
<i>H. m. amaryllis</i> (B/D)			$P < 0.05^\dagger$	
			$P < 0.05^\dagger$	
<i>Anopheles coluzzii</i> (M form)	0.00338	0.01267	$N = 17,12;$	High/low $F_{ST}$
<i>A. gambiae</i> (S form)	0.00529	0.01866	$P = 0.195$	
			$P = 0.123$	
<i>Ficedula albicollis</i>	0.00067	0.00219	50 islands,	High/low $F_{ST}$
<i>F. hypoleuca</i>	0.00132	0.0037	$P < 0.001^\dagger$	
			$P < 0.001^\dagger$	

All data come from the same references as in Table 1.

\*Based on results from Sousa *et al.* (2013).

<sup>†</sup>As reported in the original paper.

flanking regions. Although the precise genomic coordinates of the divergence islands were not specified – they were defined as regions with  $F_{ST} > 0.3$  – we attempted to match the regions of maximal relative divergence shown as grey shaded bars in Fig. 2 of Nadeau *et al.* (2012). In neither the *HmB/D* region nor the *HmYb* region is  $d_{XY}$  significantly higher in islands than nonislands (Table 1), and in fact  $d_{XY}$  is lower in the islands of relative divergence for the *B/D* region. Within both races, nucleotide diversity is lower in islands compared with nonislands for the *HmB/D* regions; in *HmYb*, diversity is lower in islands in only one of the two races considered (*H. m. amaryllis*; Table 2). Absolute divergence was higher overall in these two regions compared to three unlinked regions (Nadeau *et al.* 2012), but it is not clear exactly what this indicates.  $F_{ST}$  ranges from  $\sim 0.1$  to  $>0.7$  in the *HmB/D* region (which is more than 700 kb long), and the interpretation of such a pattern in a speciation-with-gene-flow context is that there is differential gene flow among loci within this region. Differences in absolute divergence between very large regions of the genome may indicate different mutational environments experienced by each region, but further exploration of this pattern will require genome-wide data.

One further issue with interpreting the data from these two races is whether this comparison relates to speciation at all. There is strong geographic structure involving the wing colour patterns that define these morphs as races, largely due to selection determined by colour morphs in the Müllerian mimic, *H. erato* (Mallet *et al.* 1990). But the races are not separate species: they do not show evidence of hybrid sterility or inviability and appear to be randomly mating in the narrow zone where the colour morphs overlap (Mallet *et al.* 1990). This raises the possibility that the colour-patterning loci contain locally adapted alleles within a largely panmictic (or at least continuously distributed) population and that gene flow outside of these regions represents nothing more than the normal movement of alleles within a species. In this case, there should be high  $F_{ST}$  and  $d_{XY}$  between allelic types determining colour, as the locus is behaving like a balanced polymorphism within a species (Charlesworth 2006). It is therefore possible that the higher value of  $d_{XY}$  observed in the region overall is consistent with this expectation, but the physical extent and level of divergence will be determined by the age and amount of selection on the balanced polymorphism rather than selection against migrants.

There may be other pairs of taxa within the genus *Heliconius* for which comparisons of relative and absolute divergence can be made in order to test the speciation-with-gene-flow model, but they will have to come from

pairs that are further along in the speciation process. Even if the races of *H. melpomene* analysed here are considered separate species, they are almost certainly too closely related to show significant signs of restricted gene flow by either  $F_{ST}$  or  $d_{XY}$  (Box 2).

### *Ficedula* (flycatchers)

Divergence between the collared flycatcher (*F. albicollis*) and the pied flycatcher (*F. hypoleuca*) was previously assessed by sequencing the genomes of 10 males from each species (Ellegren *et al.* 2012). These species are only very recently split, with genome-wide estimates of the average pairwise divergence between species ( $d_{XY} = 0.0046$ ) only marginally higher than the average pairwise differences within species ( $\pi = 0.0036$  and  $\pi = 0.0021$  for *F. albicollis* and *F. hypoleuca*, respectively; Ellegren *et al.* 2012). Despite this recent split, there are strong premating isolating barriers and almost complete postzygotic isolation resulting from 'prolonged periods of allopatric isolation' (Sætre & Sæther 2010). Any gene flow would therefore result from secondary contact between these species.

Using relative measures of divergence, there are strong signals of islands of differentiation between these species, with more than 50 islands where  $F_{ST} > 0.5$  (Ellegren *et al.* 2012). Values of  $F_{ST}$  are more than twice as high in islands of divergence as the genomic background (Table 1; Fig. 3d), and values of  $d_f$  are more than six times as high in the islands (Ellegren *et al.* 2012). Extreme islands show up to 50-fold higher levels of relative divergence than the background, and there are many signatures of hitchhiking events. However, the high relative divergence is accompanied by commensurately strong reductions in polymorphism within islands (Table 2), which means that these statistical measures of divergence may be misleading. Indeed, the original study reported that ' $d_{XY}$  did not exceed background levels in divergence islands' (Ellegren *et al.* 2012). Although the genotypes for each individual were not available for analysis, the mean value of  $d_{XY}$  in islands of divergence was 0.0036 (H. Ellegren, personal communication). This level of absolute divergence is much lower than the value for the genomic background in equivalently sized windows ( $d_{XY} = 0.0044$ ) and is equal to levels of within-species polymorphism in *F. albicollis*. Unfortunately, no formal statistical test can be performed for these reported values of absolute divergence between islands and nonislands.

As with previous studies in *Anopheles* (Turner *et al.* 2005; Turner & Hahn 2007), the number of fixed differences between flycatcher species was mistakenly used as an alternative to relative measures of divergence like  $F_{ST}$ . In fact, as is discussed above (see Fig. 1),  $d_f$  is

highly sensitive to levels of within-species diversity. Reduced levels of polymorphism due to linked selection will cause large shifts in the proportion of fixed differences in a region for closely related species, although this will not change  $d_{XY}$ . There are multiple signatures of strong selection in the *Ficedula* divergence islands – including reduced diversity, a skewed allele frequency spectrum, and higher linkage disequilibrium (LD) – and they largely coincide with centromeres, regions known to have reduced recombination in syntenic bird genomes (e.g. Backström *et al.* 2010). Together with the lower values of absolute divergence, the data better support a model in which regions of 50-fold higher ‘divergence’ instead indicate regions of 50-fold lower polymorphism and that many of these islands have been subject to recurrent selection since before the species split.

#### *Oryctolagus* (rabbits) and *Mus* (mice)

For both the rabbit subspecies (*O. cuniculus cuniculus* and *O. c. algirus*) and the mouse subspecies (*M. musculus musculus* and *M. m. domesticus*) considered here, the incipient species meet in a hybrid zone after periods of allopatry (Biju-Duval *et al.* 1991; Branco *et al.* 2000; Ferrand & Branco 2007; Teeter *et al.* 2008). Both cases therefore represent scenarios of secondary speciation-with-gene-flow, similar to the situation in *Ficedula*. Nevertheless, despite the relative ease with which isolating barriers may arise in allopatry, introgression upon secondary contact is expected to leave similar signatures on relative and absolute divergence measures.

We used loci that were Sanger sequenced in multiple individuals, from either rabbits (Geraldès *et al.* 2006; Carneiro *et al.* 2009, 2010) or mice (Harr 2006; Geraldès *et al.* 2011) for the analysis. Because multiple definitions are used in previous studies, we alternatively defined ‘islands’ as regions that contained more fixed differences than expected by chance based on the genome-wide average (Harr 2006); regions with low recombination relative to the genome-wide average (Geraldès *et al.* 2011); or highly differentiated regions on the X chromosome (Carneiro *et al.* 2010). Despite values of  $F_{ST}$  that ranged from 0.008 to 0.977 in rabbits, we found no significant difference in  $d_{XY}$  between the two classes of loci (Table 1; Fig. 3e). As with several previous data sets,  $d_{XY}$  was lower in the islands of relative divergence than it was in nonislands, at least for some subsets of the data (Table 1). Comparisons between other pairs of *Mus* subspecies have found higher values of absolute measures of divergence in islands (Nachman & Payseur 2012), although these were not significant. Likewise, we found no significant difference in  $d_{XY}$  between regions that did differ significantly in relative divergence in

mice (Table 1; Fig. 3f), but we did find significantly lower levels of polymorphism in the islands (Table 2).

One further caveat with using  $d_{XY}$  in both the *Mus* and *Oryctolagus* systems must be mentioned, largely due to the fact that many of the islands are on the X chromosome in both species. Because  $d_{XY}$  is affected by both the substitution rate and levels of ancestral variation, comparisons of divergence between X-linked and autosomal loci suffer from several inequalities. First, in mammals, there is an increased rate of mutation in males (Shimmin *et al.* 1993), leading to an increased rate of substitution on the autosomes (which spend half their time in males) relative to the X (which spends one-third its time in males). To avoid this problem, one can rescale levels of absolute divergence relative to an outgroup (e.g. Feder *et al.* 2005), which controls for unequal rates of substitution among loci or chromosomes. However, when such methods were applied to the *Mus* and *Oryctolagus* data sets, no significant difference in divergence levels was found between islands and nonislands (Nachman & Payseur 2012). Second, the X chromosome is expected to have three-quarters the effective population size of the autosomes under a neutral model, resulting in an expectation of three-quarters the level of ancestral polymorphism. Given the very small difference in current levels of X vs. autosomal polymorphism (e.g. 0.1% in *Oryctolagus*; Carneiro *et al.* 2010), this is unlikely to contribute to a lack of difference in  $d_{XY}$ . For both species pairs, therefore, it again appears as though a model invoking postspeciation linked selection – especially on genes in regions of reduced recombination – can explain the results at least as well as a model with differential gene flow among loci.

#### Reconciling studies of introgression with estimates of absolute divergence

The results presented above suggest that there is little evidence for differential introgression among the loci considered, for scenarios involving both primary and secondary speciation-with-gene-flow. Although the power to detect differential gene flow using absolute measures of divergence can be low for taxa that have recently split (Box 2), the same should not be true for situations of secondary contact between species because there has been enough time for mutations to accumulate. Under secondary contact, the divergence accumulated during periods of allopatry makes it easier to distinguish linked selection from introgression, although determining the minimum time in allopatry necessary requires further investigation.

These considerations lead to an apparent contradiction: for many examples of secondary contact, including



those analysed here, there is strong evidence for differential introgression of loci across hybrid zones (e.g. Sætre *et al.* 2003; Payseur *et al.* 2004; Teeter *et al.* 2008, 2010; Carneiro *et al.* 2009). How then can it be that analyses of absolute sequence divergence fail to detect this differential introgression?

We believe that the apparently contradictory results are a function of the type of data collected from different studies of introgression and not a failure of any single method for detecting introgression. Absolute measures of sequence divergence have been used to uncover differential gene flow among loci, most often as 'islands of introgression' (Garrigan *et al.* 2012) – loci with low absolute divergence set against a genomic background of high absolute divergence. Examples of such introgressed regions are found in many species, including *Drosophila* (Garrigan *et al.* 2012; Brand *et al.* 2013), *Anopheles* (Djogbénou *et al.* 2008; Etang *et al.* 2009), *Heliconius* (The Heliconius Genome Consortium 2012; Smith & Kronforst 2013) and *Mus* (Song *et al.* 2011). There is therefore little reason to believe that we could not also detect differential introgression using absolute divergence in the data analysed here.

However, for the vast majority of studies examining gene flow across hybrid zones – the main source of evidence for differential introgression – introgression is detected by using markers that are fixed (or nearly fixed) between species in populations far from the zone of overlap (e.g. Rieseberg *et al.* 1999; Sætre *et al.* 2003; Payseur *et al.* 2004; Nolte *et al.* 2006; Teeter *et al.* 2008, 2010). Such markers are generally necessary to study gene flow across a hybrid zone (for an exception, see Yatabe *et al.* 2007), but obviously cannot – almost by definition – represent loci that have freely introgressed across species boundaries and spread across species ranges. The loci that have introgressed across species and that would therefore show reduced absolute divergence when estimated from sequences collected in allopatry, would not have any fixed differences between species and would not show any difference between allopatric and sympatric comparisons. While recent secondary contact may not have allowed enough time for truly neutral markers to have introgressed across an entire species' range, it is also possible that weak selection against even the least differentiated loci that still show fixed differences in allopatry has not allowed them to introgress fully. Similar arguments apply to contrasts of differentiation among sympatric and allopatric populations of different species, at least for those species in which there is no isolation by distance.

The above examples generally come from cases of secondary contact, where there may be more power to distinguish introgressing from nonintrogressing loci,

and where it is clear that there is gene flow between lineages after an initial period of isolation. However, for cases of primary speciation-with-gene-flow, it is often very difficult to determine whether lineages are exchanging migrants at all. One popular way to infer whether gene flow is occurring is to apply a coalescent-based model for gene flow, referred to as the isolation-with-migration (IM) model (Nielsen & Wakeley 2001; Hey & Nielsen 2004, 2007; Hey 2010). In the next section, we examine the statistical properties of this model, finding that it has unacceptably high false positive rates for recently diverged taxa.

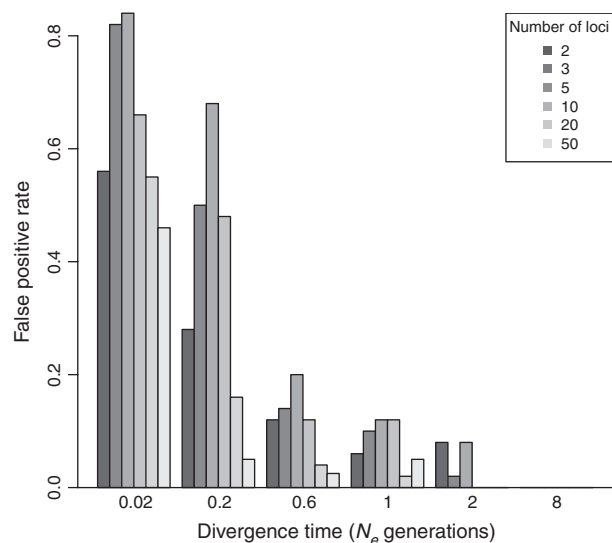
### Questions about inferring gene flow using the isolation-with-migration (IM) model

One of the most challenging tasks in molecular population genetics is to distinguish between the causes of shared polymorphisms among populations or species. Because both migration and recent splitting events result in shared polymorphisms, both of these processes can explain a wide range of values of differentiation. The IM program and its successors, IMA (Hey & Nielsen 2007) and IMA2 (Hey 2010), are intended to be used to estimate fundamental evolutionary parameters between recently diverged species, including migration rates. Comparisons between the probabilities of models with and without migration can be used to conduct a likelihood-ratio test of the null hypothesis of zero gene flow (Nielsen & Wakeley 2001; Pinho & Hey 2010).

The inference of parameters in the isolation-with-migration model using IM offers a powerful way to infer recent evolutionary processes, but is not without problems. Most difficulties arise from the assumptions made by the methods, including selective neutrality at all loci (i.e. no positive or balancing selection), no recombination within genes, no population structure within the ancestral population, independence between genes and no migration from unsampled populations. Although some of these assumptions can be relaxed without affecting the inferences made in many cases (Becquet & Przeworski 2009; Strasburg & Rieseberg 2010), we suspected that violations of these assumptions might lead to false signals of migration, even between species that were completely reproductively isolated. For instance, the IM model favours a scenario with ongoing gene flow over one with strict allopatric divergence both between *Drosophila melanogaster* and *D. simulans* (Wang & Hey 2010) and between humans and chimpanzees (Wang 2009; Mailund *et al.* 2012). Given the high degree of isolation between these species pairs, it appears that the IM model may be erroneously inferring gene flow when there is none.

A common approach to applying IM is to consider large numbers of individuals sampled at a small number of loci. The parameters of the IM model can then be estimated by Markov chain Monte Carlo (MCMC) sampling, such as is implemented in IMa2 (Hey 2010) and MIMAR (Becquet & Przeworski 2007). Studies using these methods commonly find evidence for gene flow between closely related species (e.g. Geraldès *et al.* 2006; Nadachowska & Babik 2009; Runemark *et al.* 2012; Muñoz *et al.* 2013), although the number of loci used in such studies is often far less than the number recommended in order to get accurate results (see supplementary table 1 in Pinho & Hey 2010). We used simulations to ask whether recent split times between species or small numbers of sampled loci could lead to erroneous inferences of gene flow, even when all other IM model assumptions are met. While a number of studies have examined the consequences of violating IM's assumptions on parameter estimates (e.g. Becquet & Przeworski 2009; Hey 2010; Strasburg & Rieseberg 2010), to our knowledge, the effect of recent split times – those that are most relevant to possible cases of primary speciation-with-gene-flow – on inferring whether there is any migration at all has not been reported.

We simulated populations with no migration using the program *ms* (Hudson 2002). Sequences for 15 individuals were generated from each of two species (30 chromosomes total) with splitting times ranging from less than 1 up to  $8N_e$  generations (Fig. 4). We simulated



**Fig. 4** Assessing the accuracy of IMa2 in inferring gene flow. Percentage of false positives observed for simulated populations with no gene flow. False positives occur when the likelihood-ratio statistic indicates a significant nonzero migration rate in one or two directions using IMa2 (Hey 2010). 100 simulations were run for all conditions except those representing 50 loci, where 50 simulations were run.

2–50 loci, each 2500 bp in length with no recombination. We simulated a total of 36 parameter combinations: six different split times each simulated with six different numbers of loci. For each set of parameters, 100 independent simulations were run, each producing a data set of sequences by running the *ms* output through Seq-Gen (Rambaut & Grassly 1997) with an infinite sites model specified.

Using IMa2, we attempted to estimate migration and divergence time parameters for all simulated data sets. Despite the fact that all simulations were carried out without gene flow, IMa2 consistently reported significant nonzero migration rates (Fig. 4). A data set was considered to be a 'false positive' if either of the two migration parameters (one into each species) was identified by IMa2 to be significant by a likelihood-ratio test; in many cases, we found that the inferred migration rates were highly asymmetric, often with only one parameter significant (data not shown). Inconsistencies between simulated and estimated migration rates were more common in samples with few independent loci ( $n = 2$ –10) and recent population split times ( $N_e < 1$ ). In extreme cases of few loci and very recent split times, false positive rates exceeded 80% (Fig. 4).

Our results are also consistent with previous simulation results. For some conditions tested (two to three populations, 10 loci,  $t = 2$ ), Hey (2010) reported false positive rates of 30% for individual migration parameters despite the low proportion of false positives when averaging over all parameters (see supplementary table 1 in Hey 2010). The same general pattern is reported for MIMAR: results are not reliable for fewer than 20 sampled loci (Becquet & Przeworski 2007). However, neither of these studies considered very recent split times, the factor that appears to have the greatest effect on our results.

Recently, Sousa *et al.* (2013) have relaxed the assumption within IM that all loci share a single migration history. In order to identify loci at which there is selection against gene flow, they developed a new version of IMa2 that classifies loci by their inferred effective migration rates. When applied to the same *Oryctolagus* data as analysed above, the authors identified two groups of loci: those where the data are consistent with high migration, and those with low migration. Even though the latter group is enriched for highly differentiated loci (i.e. those with high  $F_{ST}$ ),  $d_{XY}$  does not differ between the two categories and is actually lower for the low-migration set of loci (mean  $d_{XY} = 0.0358$  [low migration] vs. 0.0411 [high migration],  $P = 0.91$ ; Table 1). Clearly, support for differential gene flow among loci using the isolation-with-migration model can also be misleading.

Although it should be noted that we found the most serious problems with IM to occur at very short divergence times – more recently, for instance, than the

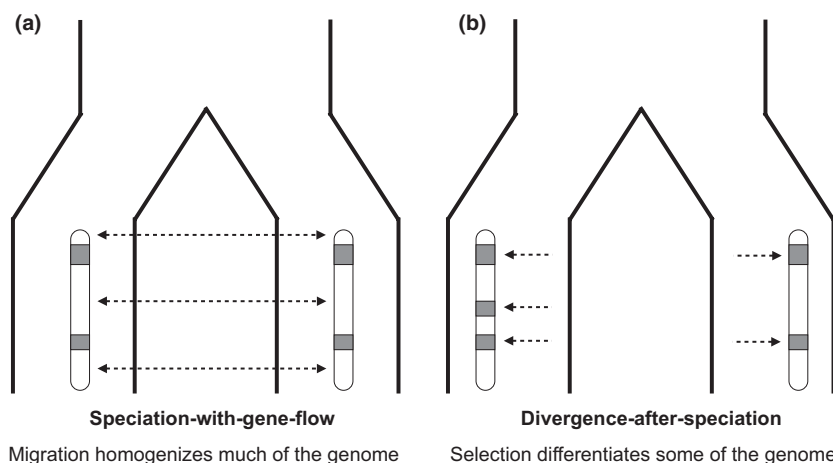
cases of secondary contact considered above – our simulations have actually met all of the assumptions of the IM model. When considering natural populations undergoing selection and recombination, it may be that these violations of the underlying model cause false signals of gene flow regardless of the divergence times. This is likely to be the problem in the primate and *Drosophila* examples cited above. Similar problems may affect any methods that do not consider the effects of selection on population differentiation when fitting models of gene flow (e.g. Duvaux *et al.* 2011; Nadachowska-Brzyska *et al.* 2013).

### An alternative model: Divergence after speciation

As discussed in the Introduction, the observed patterns of heterogeneous differentiation across genomes of closely related species have been interpreted as evidence for speciation-with-gene-flow, either primary or secondary. To be specific, the model that is often invoked proposes that low levels of differentiation are due to ongoing gene flow homogenizing a majority of the genome, while regions of high differentiation ('islands') contain loci refractory to introgression because they underlie isolating traits between the species (Fig. 5a). Under these models, we also expect higher absolute divergence in islands shortly after speciation, a pattern that was not found in the data sets considered here. The lack of support for speciation-with-gene-flow models begs the question: what model can explain the data? Although a number of alternative explanations are surely possible, we emphasize one model in which population divergence is not determined by differential migration across the genome. The outlines of this model have been proposed previously (Noor & Bennett 2009; Turner & Hahn 2010; White *et al.* 2010; Hahn *et al.* 2012;

Nachman & Payseur 2012), so here, we attempt to clarify specific predictions of the model and to assess the fit of various aspects of the data to this alternative model.

The main feature of the alternative model – variously referred to as the 'incidental island' model (Turner & Hahn 2010), the 'selection at linked sites' model (Nachman & Payseur 2012) or the 'low-gene-flow' model (Hahn *et al.* 2012) – is that heterogeneous levels of relative divergence are largely due to heterogeneous natural selection across the genome and not heterogeneous migration rates (Fig. 5b). Because any type of selection that lowers levels of linked neutral diversity will result in increased relative divergence (Charlesworth *et al.* 1997; Charlesworth 1998; Slatkin & Wiehe 1998), differences in the effects of selection across the genome will result in different levels of relative divergence. Importantly, linked selection will not increase absolute divergence (Birky & Walsh 1988), and therefore, this model predicts no islands of absolute divergence that overlap with the islands of relative divergence (though other processes can increase absolute divergence; Box 1). The model further proposes that low levels of differentiation across the majority of the genome are largely due to shared ancestral polymorphism: the more recent the split between the taxa being compared, the more similar allele frequencies will be at ancestrally inherited polymorphisms. The lineage-sorting process is a slow one, and in the absence of selection, two species will not be reciprocally monophyletic at 95% of loci (that is, they will no longer share ancestral variation) until  $9-12N_e$  generations after speciation (Hudson & Coyne 2002). At those loci where natural selection has acted in one or both taxa, however, lineage sorting is accelerated and ancestral variation is removed quickly. These regions of reduced diversity also then have increased relative divergence. The alternative model can include no gene flow or a low level of gene flow between the species



**Fig. 5** Speciation-with-gene-flow model vs. divergence-after-speciation model. Panel (a) shows a cartoon of the speciation-with-gene-flow model, with migration (arrows) homogenizing most of the genome (white); grey areas represent those loci resistant to introgression. Panel (b) shows a cartoon of the divergence-after-speciation model, with selection (arrows) driving relative divergence in a small number of regions (grey); white regions are similar between the two daughter species because of shared ancestral polymorphism.

being considered, and in fact the power to detect regions resistant to introgression will be higher with higher levels of gene flow across the genome (Box 2). The key feature of the alternative model is not the presence or absence of gene flow or selection against migrant genotypes – instead, it is that the heterogeneity in levels of divergence among loci is not due to differential effective migration, but rather differential selection.

One further elaboration to the alternative model may explain many of the patterns previously associated only with the speciation-with-gene-flow model. Although much of the selection acting to reduce levels of diversity in each descendant lineage may be due to adaptation unrelated to the speciation process, or to selection against recurrent deleterious mutations, it is also likely that immediately after speciation, many new alleles that increase adaptation are selected for. Any alleles that increase niche specialization, decrease direct competition with congeners, or that decrease hybrid matings will be immediately advantageous and are therefore likely to be some of the first targets of selection postspeciation. There are multiple signals of positive (adaptive) selection in sequence data collected from islands (e.g. Turner *et al.* 2005; Ellegren *et al.* 2012), and genes suggestive of species-specific adaptations are often found in islands (see below). The alternative model therefore suggests that islands of relative divergence should not be dismissed as nonexistent or unimportant, but instead that some of these islands may contain loci closely tied to postspeciation adaptation.

### What does the alternative model explain?

In the study thus far, we have largely concerned ourselves with examining predictions concerning levels of relative and absolute sequence differences under alternative models of speciation and divergence. For this particular set of predictions, we have argued that the data do not support models of speciation-with-gene-flow, where islands of high relative divergence represent loci that are not introgressing between species. However, lack of support for speciation-with-gene-flow models does not immediately suggest that our favoured alternative model – divergence after speciation – is strongly supported. There are multiple reasons why we may not have found patterns of sequence divergence supporting the original model (including low statistical power; Box 2), and possibly, there are alternative models that we have not considered. But we believe that the alternative model described in the previous section is a well-supported one and that it can explain multiple aspects of the data collected on genomic divergence between closely related species to date. While both

introgression and selection after speciation are likely to be occurring in many taxa, below we consider a number of different features of such data that are explained equally well, or better, by a model of divergence after speciation as compared to models of speciation-with-gene-flow.

### *Islands most often occur in regions of low recombination*

Since it was first found in *Drosophila* (Begun & Aquadro 1992), one of the most ubiquitous patterns in molecular evolution is the reduction in diversity found in regions of reduced recombination (Hahn 2008; Sella *et al.* 2009; Cutter & Payseur 2013). This reduction in diversity has been ascribed to linked selection on both advantageous (Kaplan *et al.* 1989) and deleterious mutations (Charlesworth *et al.* 1993), and also lower mutation rates in these regions (e.g. Hellmann *et al.* 2003; Kulathinal *et al.* 2008). However, close analyses of the relationship between recombination and nucleotide divergence (e.g. McGaugh *et al.* 2012), and related predictions about the effects of linked selection in reducing  $N_e$  in these regions (Pease & Hahn 2013), all indicate that linked selection (or processes that mimic selection, such as meiotic drive) is the major driver of reduced variability.

In light of the universality of these patterns, it should not be surprising that regions of reduced recombination also show signals of increased relative divergence (e.g. Begun & Aquadro 1993; Stephan *et al.* 1998; Keinan & Reich 2010). For four of the data sets analysed here (*Anopheles*, *Ficedula*, *Mus*, *Oryctolagus*), the islands of divergence are found in pericentromeric or peritelomeric regions. Such regions often show reduced levels of recombination (i.e. crossing-over) and can have reduced diversity across megabase-size regions of individual chromosomes (Cutter & Payseur 2013). This same pattern is also seen in other data sets examining genome-wide levels of relative divergence, at least when the recombination rate associated with individual markers is known (e.g. Renaut *et al.* 2013).

The effects of linked selection are not limited to only regions of low recombination. Although the length of the region affected scales negatively with the recombination rate (Kaplan *et al.* 1989), individual adaptive substitutions in regions of normal recombination can also reduce linked neutral variation (e.g. Kern *et al.* 2002; Sattath *et al.* 2011; McGaugh *et al.* 2012); this reduction in variation will affect measures of relative divergence. As expected, when islands of relative divergence have been identified in regions of normal recombination, they have been on average smaller than those regions linked to centromeres and telomeres (e.g. Turner *et al.* 2005). However, it is also the case that the length of the



region affected by selection is strongly dependent on the strength of selection, so that even in regions of typical recombination selective sweeps may reduce variation over a wide area (e.g. Tishkoff *et al.* 2007). Larger regions of higher relative divergence are therefore perfectly consistent with selection after divergence, and 'divergence hitchhiking' does not have to be invoked (cf. Via & West 2008). The effects of linked selection are also relatively independent of standing levels of linkage disequilibrium, and as such, measures of LD cannot be expected to be predictive of the length of a region affected by linked selective events (Kaplan *et al.* 1989).

The role for regions of reduced recombination in models of primary speciation-with-gene-flow is to ensure linkage disequilibrium between alleles for local adaptation and alleles for nonrandom mating (Felsenstein 1981). In this model, therefore, there are predicted to be more than one gene involved in isolating barriers within each individual region of reduced recombination. The fact that there are multiple, unlinked regions of reduced recombination across the genome is not necessarily expected under this model, as the free recombination among islands works against coadapted allelic combinations (see below). By contrast, in a model of divergence after speciation, there is nothing necessarily special about such regions. They may or may not be involved with divergent selection between the species and do not have to contain multiple targets of divergent selection. It is simply that repeated bouts of linked selection (of any kind) guarantee that signals of relative divergence are greatly increased.

#### *Colocalization of islands and species-defining traits*

In systems where genetic crosses are possible, it has been noted that there is a statistically significant correspondence between the location of quantitative trait loci (QTL) underlying traits that differ between closely related species and the location of peaks of  $F_{ST}$  (Rogers & Bernatchez 2005, 2007; Via & West 2008; Via 2012). In some cases, the traits are species-defining behaviours or phenotypes (e.g. Via 2012), while in others, they are differences in gene expression between species (Derome *et al.* 2006; Cassone *et al.* 2008; Whiteley *et al.* 2008; Renaut *et al.* 2009). There is also a coarse correspondence between QTL for reproductive isolation and regions of high differentiation (e.g. Payseur *et al.* 2004), especially involving inversions (e.g. Rieseberg *et al.* 1999; Yatabe *et al.* 2007; Kulathinal *et al.* 2009). Among the systems considered here, the peaks in relative divergence between races of *Heliconius* butterflies were found within regions known to be involved in wing colour patterning (Nadeau *et al.* 2012), a trait important for mate choice (Jiggins *et al.* 2001). In such cases,

speciation-with-gene-flow models have been invoked to explain this concordance, with the implication that there is a lack of introgression of alternative alleles found at loci controlling 'speciation phenotypes' (Shaw & Mullen 2011).

As argued above, however, it appears equally plausible to argue that many of these traits – and ones like them – would be just as important if they arose post-speciation, when reproductive isolation is largely complete. In this scenario selection on the underlying loci occurred after postzygotic or postmating reproductive isolation has already evolved, as any alleles increasing local adaptation or decreasing nonproductive matings (i.e. reinforcing already existing reproductive barriers) would be favoured. Behaviours increasing host-plant fidelity when there is local adaptation (e.g. Via 2012), or that increase within-species preferences for mates (e.g. Pennetier *et al.* 2010), both fall into this category. Selection would again act to reduce within-species polymorphism, consequently increasing between-species relative divergence.

In other cases, no known QTL or overrepresented class of genes has been found in islands of relative divergence (e.g. Ellegren *et al.* 2012), or only classes of genes that are generally found to be rapidly evolving are overrepresented in these regions (e.g. Harr 2006). It therefore seems likely that there are at least two processes that lead to islands of divergence. In the first, specific adaptive alleles that increase the fit of newly evolved species to new environments are favoured and fix within one or both species. There is no a priori reason to expect that selection in these regions has affected both daughter species, but even reductions in polymorphism in one of them will cause an increase in relative divergence. It may also be the case that such adaptive alleles are only favoured in a subset of populations within a species, leading to population-specific sweeps and population-specific islands (e.g. Turner & Hahn 2007). Such adaptive alleles may arise multiple times, leading to 'parallel' islands of divergence in multiple species (e.g. Nadeau *et al.* 2012). However, pairwise comparisons of  $F_{ST}$  values that use polymorphism from a single population in all comparisons are also nonindependent, as reduced polymorphism in one species will cause increased relative divergence in all comparisons; truly, independent comparisons are needed to investigate whether true hotspots of selection exist (e.g. Baxter *et al.* 2010). The second process creating islands represents the general, nonspecies-specific reduction in polymorphism seen in all genomes, especially the reductions associated with regions of low recombination. In these cases, it may even be background selection that acts as the major force reducing polymorphism, and no association between adaptive



traits and islands is necessarily expected. Even when hitchhiking of adaptive substitutions is causing the decrease in diversity, this adaptation may be unrelated to the speciation process.

Regardless of which process is acting, postspeciation selection does not require multiple tightly linked alleles in any single region of low recombination, as migration is not continuously introducing alternative arrangements via which recombination can break up locally favourable combinations of alleles. Therefore, no clustering of selected loci is expected, whether such loci arise in situ or by rearrangement (cf. Yeaman 2013).

#### *Near-perfect linkage disequilibrium among unlinked islands*

One of the most compelling, and puzzling, patterns of divergence between closely related species is the existence of multiple, unlinked islands. Many of the systems considered here (e.g. *Anopheles* and *Ficedula*) show high numbers of fixed differences between species at many loci throughout the genome (Turner *et al.* 2005; Ellegren *et al.* 2012). These islands of divergence segregate independently, but show near-perfect linkage disequilibrium: that is, alleles associated with one species in one island are always found in the presence of alleles associated with that same species at all the other islands. In these cases, the linkage disequilibrium is being measured in the two populations or species taken together; no linkage disequilibrium is expected within either population, except under scenarios of incomplete introgression after secondary contact.

In the speciation-with-gene-flow models, introgression in islands is prevented by selection, while recombination between the islands allows for the homogenization of the rest of the genome. However, while it is fairly easy to imagine how this can occur for a single island, if alleles at multiple islands must remain associated with each other throughout hybridization, it is much more difficult to get introgression across the genome, at least for models of primary speciation-with-gene-flow (Turner & Hahn 2010; White *et al.* 2010). In order to maintain the strong association among alleles at unlinked loci in a primary speciation-with-gene-flow model, it must be the case that either a) a large fraction of offspring with recombinant genotypes does not survive or b) some form of very strong transmission ratio distortion favours gametes with 'matching' alleles at all islands. With more than even a few islands that do not introgress because of selection, the number of recombinant individuals containing the correct combination of parental alleles at all islands – and therefore that can pass material between species – becomes vanishingly small. The fact that there is low recombination within individual islands does not

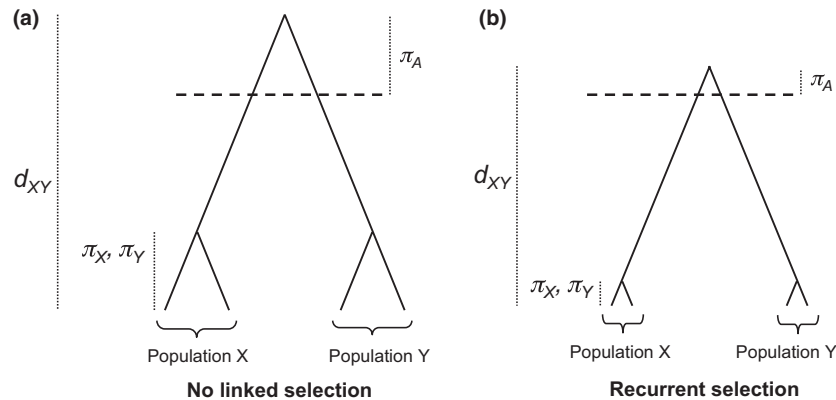
help to explain this pattern, as there is free recombination between islands on different chromosomes. The prospect of strong transmission ratio distortion has been examined in *Anopheles*, but none was found (Hahn *et al.* 2012). In the case of secondary speciation-with-gene-flow, standing hybrid zones make it much easier for neutral alleles unlinked to selectively disadvantageous alleles to introgress between species. The constant homogenization of such neutral loci by gene flow will inevitably lead to LD between alleles involved in isolating barriers. These associations will continue at least until the number of genes involved in the barriers to exchange becomes so great that effectively all loci are linked to them (Barton & Bengtsson 1986); at this point, gene flow is substantially reduced across the genome (cf. Flaxman *et al.* 2013).

In the divergence-after-speciation model, it is quite straightforward to explain perfect LD among unlinked loci: alternate alleles are fixed by selection in each species, at one locus or at many loci. If these species are already completely reproductively isolated, then there is no opposing force that generates large numbers of recombinant individuals. The alleles at each island are in LD simply because of the population structure generated by the absence of gene flow, and no selection maintaining these specific associations must be invoked. Again, the shared variation found across the rest of the genome is due to shared ancestral variation and not recent introgression between species. No strong selection or strong transmission ratio distortion is needed to explain these patterns.

#### *Lower absolute divergence in islands due to recurrent selection*

In our above analyses, we found that absolute divergence was lower in most of the islands of increased relative divergence (Table 1; Fig. 3). While linked selection in extant populations should increase relative divergence, it is not expected to decrease absolute divergence. The question therefore arises: why is absolute divergence lower in these regions?

It is important to recall that the value of  $d_{XY}$  at the time of population splitting is equal to the amount of variation present in the single prespeciation population,  $\theta_{Anc}$ . In other words, the absolute divergence between two species is made up of the accumulation of substitutions post-speciation plus the amount of variation that existed at speciation (Box 1; Gillespie & Langley 1979). This implies that while linked selection in current-day populations will not cause a decrease in  $d_{XY}$ , any form of selection that reduced variation in the ancestral population can decrease  $d_{XY}$  (Fig. 6). The correlation found between low levels of current diversity and low absolute divergence in islands therefore suggests that there have been



**Fig. 6** The effect of recurrent selection on levels of absolute divergence. Panel (a) shows a scenario where there is no linked selection in either the prespeciation ancestral population or current-day populations. In this case, measures of ancestral polymorphism ( $\pi_A$ ) and current polymorphism ( $\pi_X$  and  $\pi_Y$ ) are at their neutral expectations. Panel (b) shows a scenario where there is recurrent linked selection, such that selection has affected both ancestral and current-day populations. In this case, measures of both ancestral polymorphism and current polymorphism are lower than expected. Note that, overall, there is a correlation between  $d_{XY}$  and ancestral polymorphism levels.

recurrent bouts of linked selection at these loci, including at times pre-dating speciation (cf. Begun *et al.* 2007).

Prespeciation reductions in polymorphism also have an unfortunate effect on the power of our test for reduced gene flow in islands – they could result in even lower statistical power to detect high absolute divergence at these loci. In the worst-case scenario, a very strong selective sweep preceding speciation removes all variation, and  $d_{XY}$  at the time of speciation is 0. In order to detect the increase in  $d_{XY}$  that is indicative of a reduction in gene flow at such a locus, there would have to be very high levels of gene flow across the rest of the genome. Even less-severe reductions in ancestral diversity will have some effect on the power of our test (Box 2). These results apply solely to primary speciation-with-gene-flow, as the ability to detect differences in  $d_{XY}$  among loci after secondary contact will be determined predominantly by the time spent in allopatry.

One additional caveat concerning the identification of islands of divergence is suggested by these considerations, as well as the case of the *Heliconius* colour morphs discussed earlier. Because  $d_{XY}$  is affected by ancestral levels of diversity, loci that experience balancing selection in the ancestral population can show *increased* values of absolute divergence. While this likely will not lead to incorrect inferences when multiple balanced alleles are all inherited by both species after splitting, scenarios in which the balanced alleles segregate themselves between species can generate a false signal of an island of divergence. Examples of this would involve locally adapted alleles defining morphs that eventually become species (as could happen between *Heliconius* colour morphs), in which case the ‘island of divergence’ is already present at speciation, and need not have been formed by differential

introgression after the evolution of reproductive isolation. Higher absolute divergence between species within inversions may constitute one stereotypical example of this phenomenon, especially as inversions can show extreme patterns of local adaptation and sequence divergence even within species (e.g. Cheng *et al.* 2012).

#### *Species trees are recapitulated in islands*

Many studies that have conducted comparisons of relative divergence across the genomes of closely related organisms have also been able to construct phylogenetic trees describing relationships among sampled populations (reviewed in Nosil *et al.* 2009). The main conclusion from examining these trees is that those made from loci within islands of high relative divergence often reflect the known species relationships, while those made from loci outside islands do not (Via & West 2008; Nosil *et al.* 2009; Keller *et al.* 2013). In the context of speciation-with-gene-flow models, this difference is thought to reflect differences in gene flow at the two types of loci: selection maintains the ‘true’ species tree at loci within islands, while migration obscures the true relationships at all other loci.

However, we may just as easily interpret this difference in tree topology as due to selection after speciation, with no effect at all of migration. As mentioned earlier, it takes a long time after an initial split for species to become reciprocally monophyletic (Hudson & Coyne 2002). This is because the lineage-sorting process due to drift alone is quite slow, and ancestral variation can be shared among daughter lineages for thousands to millions of generations after speciation. As a result, gene trees will not necessarily reflect current species

relationships until well after speciation has occurred and full reproductive isolation has evolved (Avice 2004). In contrast, when selection acts at a locus, it greatly accelerates lineage sorting, removing ancestral variation and making species reciprocally monophyletic. Therefore, in the context of the divergence-after-speciation model, trees constructed from islands of high relative divergence more accurately reflect species relationships because selection has reduced the effective population size, accelerating the lineage-sorting process (cf. Pease & Hahn 2013). Trees constructed from outside islands do not necessarily reflect species relationships because they still contain large amounts of ancestral polymorphism; however, this does not indicate anything about the differential effects of migration at such loci.

In general, divergence leads to lineage sorting. Over time the genomes of closely related species diverge, with the accumulation of species-specific genetic changes and the sorting out of ancestral variation among lineages. Both of these processes are expected to be somewhat heterogeneous, with variation in ancestral coalescence times, differential selection and differential mutation all causing loci to diverge at different rates. Species pairs examined after the initial split will show heterogeneous patterns of divergence – and divergence will increase with the time since the split – but no novel processes or models of speciation must be invoked to explain the data (e.g. Feder *et al.* 2012).

## Conclusions

We have argued here that a fuller understanding of heterogeneous genomic divergence requires that we use both relative and absolute measures of sequence divergence. In particular, although relative measures have been used in the vast majority of studies examining the possibility of speciation-with-gene-flow, elevated values of these statistics (e.g.  $F_{ST}$ ,  $d_a$  and  $d_f$ ) can be due to selection rather than reduced gene flow. Relative and absolute measures of sequence divergence are both appropriate summaries of nucleotide differences, and either may be preferred under certain circumstances. However, in weighing the claims regarding differential migration among loci, we expect the two types of measures to agree; in all of the cases examined here, they do not agree. Linked selection therefore appears to be a more likely driver of the observed patterns.

Our analyses should raise concerns about the existing evidence for primary speciation-with-gene-flow (i.e. sympatric speciation), especially for systems without clear evidence of secondary contact. For most of the systems considered here, there is good external evidence that the two species are truly exchanging genes, but in many other cases, more limited evidence

supports inferences about gene flow. We have shown that a particular application of the isolation-with-migration model to sequence data has unacceptably high false positive rates, such that migration is inferred even when none exists. We caution the further use of such models with recently diverged taxa and/or small numbers of loci. A further piece of data often cited as evidence for gene flow is the presence of hybrids, often only  $F_1$  hybrids. If these hybrids, or later-generation individuals, exhibit strong extrinsic or intrinsic incompatibilities, then their existence does not always indicate actual gene flow. Most importantly, even in those cases in which there is demonstrable introgression between recently separated species – whether or not there are different levels of introgression among loci – this does not constitute evidence for sympatric speciation. While gene flow shortly after speciation is certainly a challenge to the existence of two independently evolving taxa, it is the rare case in which the conditions necessary for truly sympatric speciation have been met (e.g. Savolainen *et al.* 2006). In many cases, gene flow may simply be occurring after isolation has evolved in allopatry, and secondary contact should not be confused with primary ‘speciation-with-gene-flow’.

Finally, it should be noted that none of our results preclude the existence of true islands of divergence in other systems that have yet to be examined using absolute divergence (e.g. sticklebacks; Hohenlohe *et al.* 2010; Jones *et al.* 2012). If two ‘good’ species are exchanging genes, then it would be expected that loci underlying ecological and/or reproductive isolation would not flow freely between them. Although there may be further alternative explanations for observed patterns of divergence (e.g. Bierne *et al.* 2013), we expect that multiple systems will be found in the near future with the expected pattern. Evidence that individual loci do show reduced gene flow will simply require that absolute measures of sequence divergence agree with the patterns so often shown by  $F_{ST}$ .

## Acknowledgements

We thank N. Besansky, M. Lawniczak, S. Emrich, K. Reidenbach, H. Ellegren and B. Harr for sharing data; L. Moyle, G. Coop, M. Whitlock, J. Mallet, J. Hey, R. Guerrero, S. Yeaman, J.B.W. Wolf, B. Payseur and R. Harrison for helpful discussions; and L. Bernatchez and three anonymous reviewers for their constructive comments. This work is funded by National Science Foundation Grant MCB-1127059.

## References

- Akerman A, Bürger R (2014) The consequences of gene flow for local adaptation and differentiation: a two-locus two-deme model. *Journal of Mathematical Biology*, **68**, 1135–1198.

- Avise JC (2004) *Molecular Markers, Natural History, and Evolution*. Sinauer Associates, Sunderland, Massachusetts.
- Backström N, Forstmeier W, Schielzeth H *et al.* (2010) The recombination landscape of the zebra finch *Taeniopygia guttata* genome. *Genome Research*, **20**, 485–495.
- Baird NA, Etter PD, Atwood TS *et al.* (2008) Rapid SNP discovery and genetic mapping using sequenced RAD markers. *PLoS ONE*, **3**, e3376.
- Barbash DA, Siino DF, Tarone AM, Roote J (2003) A rapidly evolving MYB-related protein causes species isolation in *Drosophila*. *Proceedings of the National Academy of Sciences*, **100**, 5302–5307.
- Barton NH (2006) Evolutionary biology: how did the human species form? *Current Biology*, **16**, R647–R650.
- Barton N, Bengtsson BO (1986) The barrier to genetic exchange between hybridising populations. *Heredity*, **57**, 357.
- Baxter SW, Nadeau NJ, Maroja LS *et al.* (2010) Genomic hot-spots for adaptation: the population genetics of Müllerian mimicry in the *Heliconius melpomene* clade. *PLoS Genetics*, **6**, e1000794.
- Becquet C, Przeworski M (2007) A new approach to estimate parameters of speciation models with application to apes. *Genome Research*, **17**, 1505–1519.
- Becquet C, Przeworski M (2009) Learning about modes of speciation by computational approaches. *Evolution*, **63**, 2547–2562.
- Begun DJ, Aquadro CF (1992) Levels of naturally occurring DNA polymorphism correlate with recombination rates in *Drosophila melanogaster*. *Nature*, **356**, 519–520.
- Begun DJ, Aquadro CF (1993) African and North American populations of *Drosophila melanogaster* are very different at the DNA level. *Nature*, **365**, 548–550.
- Begun DJ, Holloway AK, Stephens K *et al.* (2007) Population genomics: whole-genome analysis of polymorphism and divergence in *Drosophila simulans*. *PLoS Biology*, **5**, e310.
- Bierne N, Gagnaire P-A, David P (2013) The geography of introgression in a patchy environment and the thorn in the side of ecological speciation. *Current Zoology*, **59**, 72–86.
- Biju-Duval C, Ennafaa H, Dennebouy N *et al.* (1991) Mitochondrial DNA evolution in lagomorphs: origin of systematic heteroplasmy and organization of diversity in European rabbits. *Journal of Molecular Evolution*, **33**, 92–102.
- Birky CW, Walsh JB (1988) Effects of linkage on rates of molecular evolution. *Proceedings of the National Academy of Sciences*, **85**, 6414–6418.
- Branco M, Ferrand N, Monnerot M (2000) Phylogeography of the European rabbit (*Oryctolagus cuniculus*) in the Iberian peninsula inferred from RFLP analysis of the cytochrome b gene. *Heredity*, **85**, 307–317.
- Brand CL, Kingan SB, Wu L, Garrigan D (2013) A selective sweep across species boundaries in *Drosophila*. *Molecular Biology and Evolution*, **30**, 2177–2186.
- Bulmer MG (1972) Multiple niche polymorphism. *The American Naturalist*, **106**, 254–257.
- Carneiro M, Ferrand N, Nachman MW (2009) Recombination and speciation: loci near centromeres are more differentiated than loci near telomeres between subspecies of the European rabbit (*Oryctolagus cuniculus*). *Genetics*, **181**, 593–606.
- Carneiro M, Blanco-Aguilar JA, Villafuerte R, Ferrand N, Nachman MW (2010) Speciation in the European rabbit (*Oryctolagus cuniculus*): Islands of differentiation on the X chromosome and autosomes. *Evolution*, **64**, 3443–3460.
- Cassone BJ, Mouline K, Hahn MW *et al.* (2008) Differential gene expression in incipient species of *Anopheles gambiae*. *Molecular Ecology*, **17**, 2491–2504.
- Charlesworth B (1998) Measures of divergence between populations and the effect of forces that reduce variability. *Molecular Biology and Evolution*, **15**, 538–543.
- Charlesworth D (2006) Balancing selection and its effects on sequences in nearby genome regions. *PLoS Genetics*, **2**, e64.
- Charlesworth B, Morgan MT, Charlesworth D (1993) The effect of deleterious mutations on neutral molecular variation. *Genetics*, **134**, 1289–1303.
- Charlesworth B, Nordborg M, Charlesworth D (1997) The effects of local selection, balanced polymorphism and background selection on equilibrium patterns of genetic diversity in subdivided populations. *Genetical Research*, **70**, 155–174.
- Cheng C, White BJ, Camdem C *et al.* (2012) Ecological genomics of *Anopheles gambiae* along a latitudinal cline in Cameroon: a population resequencing approach. *Genetics*, **190**, 1417–1432.
- Coetzee M, Hunt RH, Wilkerson R *et al.* (2013) *Anopheles coluzzii* and *Anopheles amharicus*, new members of the *Anopheles gambiae* complex. *Zootaxa*, **3619**, 246–274.
- Cutter AD, Payseur BA (2013) Genomic signatures of selection at linked sites: unifying the disparity among species. *Nature Reviews Genetics*, **14**, 262–274.
- Derome N, Duchesne P, Bernatchez L (2006) Parallelism in gene transcription among sympatric lake whitefish (*Coregonus clupeaformis* Mitchill) ecotypes. *Molecular Ecology*, **15**, 1239–1249.
- Djogbénou L, Chandre F, Berthomieu A *et al.* (2008) Evidence of introgression of the *ace-1<sup>R</sup>* mutation and of the *ace-1* duplication in West African *Anopheles gambiae* s. s. *PLoS ONE*, **3**, e2172.
- Duvaux L, Belkhir K, Boulesteix M, Boursot P (2011) Isolation and gene flow: inferring the speciation history of European house mice. *Molecular Ecology*, **20**, 5248–5264.
- Ellegren H, Smeds L, Burri R *et al.* (2012) The genomic landscape of species divergence in *Ficedula* flycatchers. *Nature*, **491**, 756–760.
- Etang J, Vicente JL, Nwane P *et al.* (2009) Polymorphism of intron-1 in the voltage-gated sodium channel gene of *Anopheles gambiae* s.s. populations from Cameroon with emphasis on insecticide knockdown resistance mutations. *Molecular Ecology*, **18**, 3076–3086.
- Excoffier L (2007) Analysis of population subdivision. In: *Handbook of Statistical Genetics*, vol. 2 (eds Balding DJ, Bishop M, Cannings C), pp. 980–1020. Wiley, West Sussex.
- Feder JL, Nosil P (2010) The efficacy of divergence hitchhiking in generating genomic islands during ecological speciation. *Evolution*, **64**, 1729–1747.
- Feder JL, Xie X, Rull J *et al.* (2005) Mayr, Dobzhansky, and Bush and the complexities of sympatric speciation in *Rhagoletis*. *Proceedings of the National Academy of Sciences*, **102**, 6573–6580.
- Feder JL, Egan SP, Nosil P (2012) The genomics of speciation-with-gene-flow. *Trends in Genetics*, **28**, 342–350.
- Felsenstein J (1981) Skepticism towards Santa Rosalia, or why are there so few kinds of animals? *Evolution*, **35**, 124–138.
- Ferrand N, Branco M (2007) The evolutionary history of the European rabbit (*Oryctolagus cuniculus*): major patterns of population differentiation and geographic expansion inferred



- from protein polymorphism. In: *Phylogeography of Southern European Refugia* (eds Weiss S, Ferrand N), pp. 207–235. Springer.
- Flaxman S, Feder J, Nosil P (2012) Spatially explicit models of divergence and genome hitchhiking. *Journal of Evolutionary Biology*, **25**, 2633–2650.
- Flaxman SM, Feder JL, Nosil P (2013) Genetic hitchhiking and the dynamic buildup of genomic divergence during speciation with gene flow. *Evolution*, **67**, 2577–2591.
- Garrigan D, Kingan SB, Geneva AJ *et al.* (2012) Genome sequencing reveals complex speciation in the *Drosophila simulans* clade. *Genome Research*, **22**, 1499–1511.
- Geraldes A, Ferrand N, Nachman MW (2006) Contrasting patterns of introgression at X-linked loci across the hybrid zone between subspecies of the European rabbit (*Oryctolagus cuniculus*). *Genetics*, **173**, 919–933.
- Geraldes A, Basset P, Smith KL, Nachman MW (2011) Higher differentiation among subspecies of the house mouse (*Mus musculus*) in genomic regions with low recombination. *Molecular Ecology*, **20**, 4722–4736.
- Gillespie JH, Langley CH (1979) Are evolutionary rates really variable? *Journal of Molecular Evolution*, **13**, 27–34.
- Hahn MW (2008) Toward a selection theory of molecular evolution. *Evolution*, **62**, 255–265.
- Hahn MW, White BJ, Muir CD, Besansky NJ (2012) No evidence for biased co-transmission of speciation islands in *Anopheles gambiae*. *Philosophical Transactions of the Royal Society of London. Series B, Biological Sciences*, **367**, 374–384.
- Haldane JBS (1930) A mathematical theory of natural and artificial selection (Part VI. Isolation). *Proceedings of the Cambridge Philosophical Society*, **26**, 220–230.
- Harr B (2006) Genomic islands of differentiation between house mouse subspecies. *Genome Research*, **16**, 730–737.
- Harrison RG (2012) The language of speciation. *Evolution*, **66**, 3643–3657.
- Hellmann I, Ebersberger I, Ptak SE, Pääbo S, Przeworski M (2003) A neutral explanation for the correlation of diversity with recombination rates in humans. *The American Journal of Human Genetics*, **72**, 1527–1535.
- Hey J (1991) The structure of genealogies and the distribution of fixed differences between DNA sequence samples from natural populations. *Genetics*, **128**, 831–840.
- Hey J (2010) Isolation with migration models for more than two populations. *Molecular Biology and Evolution*, **27**, 905–920.
- Hey J, Nielsen R (2004) Multilocus methods for estimating population sizes, migration rates and divergence time, with applications to the divergence of *Drosophila pseudoobscura* and *D. persimilis*. *Genetics*, **167**, 747–760.
- Hey J, Nielsen R (2007) Integration within the Felsenstein equation for improved Markov chain Monte Carlo methods in population genetics. *Proceedings of the National Academy of Sciences*, **104**, 2785–2790.
- Hohenlohe PA, Bassham S, Etter PD *et al.* (2010) Population genomics of parallel adaptation in threespine stickleback using sequenced RAD tags. *PLoS Genetics*, **6**, e1000862.
- Holsinger KE, Weir BS (2009) Genetics in geographically structured populations: defining, estimating and interpreting  $F_{ST}$ . *Nature Reviews Genetics*, **10**, 639–650.
- Hopkins R, Rausher MD (2011) Identification of two genes causing reinforcement in the Texas wildflower *Phlox drummondii*. *Nature*, **469**, 411–414.
- Hudson RR (2002) Generating samples under a Wright-Fisher neutral model of genetic variation. *Bioinformatics*, **18**, 337–338.
- Hudson RR, Coyne JA (2002) Mathematical consequences of the genealogical species concept. *Evolution*, **56**, 1557–1565.
- Jakobsson M, Edge MD, Rosenberg NA (2013) The relationship between  $F_{ST}$  and the frequency of the most frequent allele. *Genetics*, **193**, 515–528.
- Jiggins CD (2008) Ecological speciation in mimetic butterflies. *BioScience*, **58**, 541–548.
- Jiggins CD, Naisbit R, Coe R, Mallet J (2001) Reproductive isolation caused by colour pattern mimicry. *Nature*, **411**, 302–305.
- Jones FC, Grabherr MG, Chan YF *et al.* (2012) The genomic basis of adaptive evolution in threespine sticklebacks. *Nature*, **484**, 55–61.
- Kaplan NL, Hudson R, Langley C (1989) The “hitchhiking effect” revisited. *Genetics*, **123**, 887–899.
- Keinan A, Reich D (2010) Human population differentiation is strongly correlated with local recombination rate. *PLoS Genetics*, **6**, e1000886.
- Keller I, Wagner CE, Greuter L *et al.* (2013) Population genomic signatures of divergent adaptation, gene flow and hybrid speciation in the rapid radiation of Lake Victoria cichlid fishes. *Molecular Ecology*, **22**, 2848–2863.
- Kern AD, Jones CD, Begun DJ (2002) Genomic effects of nucleotide substitutions in *Drosophila simulans*. *Genetics*, **162**, 1753–1761.
- Kronforst MR, Young LG, Blume LM, Gilbert LE (2006) Multilocus analyses of admixture and introgression among hybridizing *Heliconius* butterflies. *Evolution*, **60**, 1254–1268.
- Kulathinal RJ, Bennett SM, Fitzpatrick CL, Noor MA (2008) Fine-scale mapping of recombination rate in *Drosophila* refines its correlation to diversity and divergence. *Proceedings of the National Academy of Sciences*, **105**, 10051–10056.
- Kulathinal RJ, Stevison LS, Noor MA (2009) The genomics of speciation in *Drosophila*: diversity, divergence, and introgression estimated using low-coverage genome sequencing. *PLoS Genetics*, **5**, e1000550.
- Lawniczak MKN, Emrich SJ, Holloway AK *et al.* (2010) Widespread divergence between incipient *Anopheles gambiae* species revealed by whole genome sequences. *Science*, **330**, 512–514.
- Lee Y, Marsden CD, Norris LC *et al.* (2013) Spatiotemporal dynamics of gene flow and hybrid fitness between the M and S forms of the malaria mosquito, *Anopheles gambiae*. *Proceedings of the National Academy of Sciences*, **110**, 19854–19859.
- Mailund T, Halager AE, Westergaard M *et al.* (2012) A new isolation with migration model along complete genomes infers very different divergence processes among closely related great ape species. *PLoS Genetics*, **8**, e1003125.
- Mallet J, Barton N, Lamas G *et al.* (1990) Estimates of selection and gene flow from measures of cline width and linkage disequilibrium in *Heliconius* hybrid zones. *Genetics*, **124**, 921–936.
- Maroja LS, Andrés JA, Harrison RG (2009) Genealogical discordance and patterns of introgression and selection across a cricket hybrid zone. *Evolution*, **63**, 2999–3015.
- Martin SH, Dasmahapatra KK, Nadeau NJ *et al.* (2013) Genome-wide evidence for speciation with gene flow in *Heliconius* butterflies. *Genome Research*, **23**, 1817–1828.
- Maynard Smith J, Haigh J (1974) The hitch-hiking effect of a favourable gene. *Genetical Research*, **23**, 23–35.



- McGaugh SE, Heil CS, Manzano-Winkler B *et al.* (2012) Recombination modulates how selection affects linked sites in *Drosophila*. *PLoS Biology*, **10**, e1001422.
- Michel AP, Sim S, Powell THQ *et al.* (2010) Widespread genomic divergence during sympatric speciation. *Proceedings of the National Academy of Sciences*, **107**, 9724–9729.
- Muñoz MM, Crawford NG, McGreevy TJ *et al.* (2013) Divergence in coloration and ecological speciation in the *Anolis marmoratus* species complex. *Molecular Ecology*, **22**, 2668–2682.
- Nachman MW, Payseur BA (2012) Recombination rate variation and speciation: theoretical predictions and empirical results from rabbits and mice. *Philosophical Transactions of the Royal Society of London. Series B, Biological Sciences*, **367**, 409–421.
- Nadachowska K, Babik W (2009) Divergence in the face of gene flow: the case of two newts (Amphibia: Salamandridae). *Molecular Biology and Evolution*, **26**, 829–841.
- Nadachowska-Brzyska K, Burri R, Olason PI *et al.* (2013) Demographic divergence history of pied flycatcher and collared flycatcher inferred from whole-genome re-sequencing data. *PLoS Genetics*, **9**, e1003942.
- Nadeau NJ, Whibley A, Jones RT *et al.* (2012) Genomic islands of divergence in hybridizing *Heliconius* butterflies identified by large-scale targeted sequencing. *Philosophical Transactions of the Royal Society of London. Series B, Biological Sciences*, **367**, 343–353.
- Nadeau NJ, Martin SH, Kozak KM *et al.* (2013) Genome-wide patterns of divergence and gene flow across a butterfly radiation. *Molecular Ecology*, **22**, 814–826.
- Neafsey D, Lawniczak M, Park D *et al.* (2010) SNP genotyping defines complex gene-flow boundaries among African malaria vector mosquitoes. *Science*, **330**, 514–517.
- Nei M (1973) Analysis of gene diversity in subdivided populations. *Proceedings of the National Academy of Sciences*, **70**, 3321–3323.
- Nei M (1982) Evolution of human races at the gene level. In: *Human Genetics, Part A: The Unfolding Genome* (ed. Bonne-Tamir B), pp. 167–181. Alan R. Liss, New York.
- Nei M (1987) *Molecular Evolutionary Genetics*. Columbia University Press, New York.
- Nei M, Li W-H (1979) Mathematical model for studying genetic variation in terms of restriction endonucleases. *Proceedings of the National Academy of Sciences*, **76**, 5269–5273.
- Nielsen R, Wakeley J (2001) Distinguishing migration from isolation: a Markov chain Monte Carlo approach. *Genetics*, **158**, 885–896.
- Nolte AW, Freyhof J, Tautz D (2006) When invaders meet locally adapted types: rapid moulding of hybrid zones between sculpins (*Cottus*, Pisces) in the Rhine system. *Molecular Ecology*, **15**, 1983–1993.
- Noor MAF, Bennett SM (2009) Islands of speciation or mirages in the desert? Examining the role of restricted recombination in maintaining species. *Heredity*, **103**, 439–444.
- Nordborg M (1997) Structured coalescent processes on different time scales. *Genetics*, **146**, 1501–1514.
- Nordborg M, Charlesworth B, Charlesworth D (1996) The effect of recombination on background selection. *Genetical Research*, **67**, 159–174.
- Nosil P (2008) Speciation with gene flow could be common. *Molecular Ecology*, **17**, 2103–2106.
- Nosil P, Funk DJ, Ortiz-Barrientos D (2009) Divergent selection and heterogeneous genomic divergence. *Molecular Ecology*, **18**, 375–402.
- Pardo-Diaz C, Salazar C, Baxter SW *et al.* (2012) Adaptive introgression across species boundaries in *Heliconius* butterflies. *PLoS Genetics*, **8**, e1002752.
- Payseur BA, Krenz JG, Nachman MW (2004) Differential patterns of introgression across the X chromosome in a hybrid zone between two species of house mice. *Evolution*, **58**, 2064–2078.
- Pease JB, Hahn MW (2013) More accurate phylogenies inferred from low-recombination regions in the presence of incomplete lineage sorting. *Evolution*, **67**, 2376–2384.
- Pennetier C, Warren B, Dabiré KR, Russell IJ, Gibson G (2010) “Singing on the wing” as a mechanism for species recognition in the malarial mosquito *Anopheles gambiae*. *Current Biology*, **20**, 131–136.
- Pinho C, Hey J (2010) Divergence with gene flow: models and data. *Annual review of ecology, evolution, and systematics*, **41**, 215–230.
- Presgraves DC, Balagopalan L, Abmayr SM, Orr HA (2003) Adaptive evolution drives divergence of a hybrid inviability gene between two species of *Drosophila*. *Nature*, **423**, 715–719.
- Rambaut A, Grassly NC (1997) Seq-Gen: an application for the Monte Carlo simulation of DNA sequence evolution along phylogenetic trees. *Computer applications in the biosciences: CABIOS*, **13**, 235–238.
- Reidenbach KR, Neafsey DE, Costantini C *et al.* (2012) Patterns of genomic differentiation between ecologically differentiated M and S forms of *Anopheles gambiae* in West And Central Africa. *Genome Biology and Evolution*, **4**, 1202–1212.
- Renaut S, Nolte A, Bernatchez L (2009) Gene expression divergence and hybrid misexpression between lake whitefish species pairs (*Coregonus* spp. Salmonidae). *Molecular Biology and Evolution*, **26**, 925–936.
- Renaut S, Grassa C, Yeaman S *et al.* (2013) Genomic islands of divergence are not affected by geography of speciation in sunflowers. *Nature Communications*, **4**, 1827.
- Rieseberg LH, Whitton J, Gardner K (1999) Hybrid zones and the genetic architecture of a barrier to gene flow between two sunflower species. *Genetics*, **152**, 713–727.
- Rogers S, Bernatchez L (2005) Integrating QTL mapping and genome scans towards the characterization of candidate loci under parallel selection in the lake whitefish (*Coregonus clupeaformis*). *Molecular Ecology*, **14**, 351–361.
- Rogers S, Bernatchez L (2007) The genetic architecture of ecological speciation and the association with signatures of selection in natural lake whitefish (*Coregonus* sp. Salmonidae) species pairs. *Molecular Biology and Evolution*, **24**, 1423–1438.
- Runemark A, Hey J, Hansson B, Svensson EI (2012) Vicariance divergence and gene flow among islet populations of an endemic lizard. *Molecular Ecology*, **21**, 117–129.
- Sætre G-P, Sæther SA (2010) Ecology and genetics of speciation in *Ficedula* flycatchers. *Molecular Ecology*, **19**, 1091–1106.
- Sætre G-P, Borge T, Lindroos K *et al.* (2003) Sex chromosome evolution and speciation in *Ficedula* flycatchers. *Proceedings of the Royal Society of London. Series B: Biological Sciences*, **270**, 53–59.
- Sattath S, Elyashiv E, Kolodny O, Rinott Y, Sella G (2011) Pervasive adaptive protein evolution apparent in diversity

- patterns around amino acid substitutions in *Drosophila simulans*. *PLoS Genetics*, **7**, e1001302.
- Savolainen V, Anstett M-C, Lexer C *et al.* (2006) Sympatric speciation in palms on an oceanic island. *Nature*, **441**, 210–213.
- Sella G, Petrov DA, Przeworski M, Andolfatto P (2009) Pervasive natural selection in the *Drosophila* genome? *PLoS Genetics*, **5**, e1000495.
- Shaw KL, Mullen SP (2011) Genes versus phenotypes in the study of speciation. *Genetica*, **139**, 649–661.
- Shimmin LC, Chang BH, Li W-H (1993) Male-driven evolution of DNA sequences. *Nature*, **362**, 745–747.
- Slatkin M, Wiehe T (1998) Genetic hitch-hiking in a subdivided population. *Genetical Research*, **71**, 155–160.
- Smith J, Kronforst MR (2013) Do *Heliconius* butterfly species exchange mimicry alleles? *Biology Letters*, **9**, 20130503.
- Song Y, Endepols S, Klemann N *et al.* (2011) Adaptive introgression of anticoagulant rodent poison resistance by hybridization between Old World mice. *Current Biology*, **21**, 1296–1301.
- Sousa VC, Carneiro M, Ferrand N, Hey J (2013) Identifying loci under selection against gene flow in isolation-with-migration models. *Genetics*, **194**, 211–233.
- Stephan W, Xing L, Kirby DA, Braverman JM (1998) A test of the background selection hypothesis based on nucleotide data from *Drosophila ananassae*. *Proceedings of the National Academy of Sciences*, **95**, 5649–5654.
- Strasburg JL, Rieseberg LH (2010) How robust are “Isolation with Migration” analyses to violations of the IM model? A simulation study. *Molecular Biology and Evolution*, **27**, 297–310.
- Stump AD, Fitzpatrick MC, Lobo NF *et al.* (2005) Centromere-proximal differentiation and speciation in *Anopheles gambiae*. *Proceedings of the National Academy of Sciences*, **102**, 15930–15935.
- Teeter KC, Payseur BA, Harris LW *et al.* (2008) Genome-wide patterns of gene flow across a house mouse hybrid zone. *Genome Research*, **18**, 67–76.
- Teeter KC, Thibodeau LM, Gompert Z *et al.* (2010) The variable genomic architecture of isolation between hybridizing species of house mice. *Evolution*, **64**, 472–485.
- The Heliconius Genome Consortium (2012) Butterfly genome reveals promiscuous exchange of mimicry adaptations among species. *Nature*, **487**, 94–98.
- Tishkoff SA, Reed FA, Ranciaro A *et al.* (2007) Convergent adaptation of human lactase persistence in Africa and Europe. *Nature Genetics*, **39**, 31–40.
- Tripet F, Toure Y, Taylor C *et al.* (2001) DNA analysis of transferred sperm reveals significant levels of gene flow between molecular forms of *Anopheles gambiae*. *Molecular Ecology*, **10**, 1725–1732.
- Turner TL, Hahn MW (2007) Locus- and population-specific selection and differentiation between incipient species of *Anopheles gambiae*. *Molecular Biology and Evolution*, **24**, 2132–2138.
- Turner TL, Hahn MW (2010) Genomic islands of speciation or genomic islands and speciation? *Molecular Ecology*, **19**, 848–850.
- Turner T, Hahn M, Nuzhdin SV (2005) Genomic islands of speciation in *Anopheles gambiae*. *PLoS Biology*, **3**, e285.
- Turner TL, Bourne EC, Von Wettberg EJ, Hu TT, Nuzhdin SV (2010) Population resequencing reveals local adaptation of *Arabidopsis lyrata* to serpentine soils. *Nature Genetics*, **42**, 260–263.
- Via S (2009) Natural selection in action during speciation. *Proceedings of the National Academy of Sciences*, **106**, 9939–9946.
- Via S (2012) Divergence hitchhiking and the spread of genomic isolation during ecological speciation-with-gene-flow. *Philosophical Transactions of the Royal Society of London. Series B, Biological Sciences*, **367**, 451–460.
- Via S, West J (2008) The genetic mosaic suggests a new role for hitchhiking in ecological speciation. *Molecular Ecology*, **17**, 4334–4345.
- Wang Y (2009) *Estimating the Process of Speciation for Humans and Chimpanzees*. PhD Thesis, Rutgers University.
- Wang Y, Hey J (2010) Estimating divergence parameters with small samples from a large number of loci. *Genetics*, **184**, 363–379.
- Wang R, Zheng L, Touré YT, Dandekar T, Kafatos FC (2001) When genetic distance matters: measuring genetic differentiation at microsatellite loci in whole-genome scans of recent and incipient mosquito species. *Proceedings of the National Academy of Sciences*, **98**, 10769–10774.
- White BJ, Cheng C, Simard F, Costantini C, Besansky NJ (2010) Genetic association of physically unlinked islands of genomic divergence in incipient species of *Anopheles gambiae*. *Molecular Ecology*, **19**, 925–939.
- Whiteley AR, Derome N, Rogers SM *et al.* (2008) The phenomics and expression quantitative trait locus mapping of brain transcriptomes regulating adaptive divergence in lake whitefish species pairs (*Coregonus* sp.). *Genetics*, **180**, 147–164.
- Wittbrodt J, Adam D, Malitschek B *et al.* (1989) Novel putative receptor tyrosine kinase encoded by the melanoma-inducing *Tu* locus in *Xiphophorus*. *Nature*, **341**, 415–421.
- Wright S (1931) Evolution in Mendelian populations. *Genetics*, **16**, 97–159.
- Wright S (1943) Isolation by distance. *Genetics*, **28**, 114–138.
- Wu C-I (2001) The genic view of the process of speciation. *Journal of Evolutionary Biology*, **14**, 851–865.
- Yatabe Y, Kane NC, Scotti-Saintagne C, Rieseberg LH (2007) Rampant gene exchange across a strong reproductive barrier between the annual sunflowers, *Helianthus annuus* and *H. petiolaris*. *Genetics*, **175**, 1883–1893.
- Yeaman S (2013) Genomic rearrangements and the evolution of clusters of locally adaptive loci. *Proceedings of the National Academy of Sciences*, **110**, E1743–E1751.

---

M.W.H. conceived of the project. T.E.C. carried out the analyses. T.E.C. and M.W.H. wrote the paper.

---

## Data accessibility

Data available from the Dryad Digital Repository: <http://dx.doi.org/10.5061/dryad.cc2sk>.