

# Emergent patterns of population genetic structure for a coral reef community

KIMBERLY A. SELKOE,\*† OSCAR E. GAGGIOTTI,‡ TOBO LABORATORY,§ BRIAN W. BOWEN\* and ROBERT J. TOONEN\*

\*Hawai'i Institute of Marine Biology, University of Hawai'i, Kāne'ohe, HI 97644, USA, †National Center for Ecological Analysis and Synthesis, 735 State St., Santa Barbara, CA 93101, USA, ‡Scottish Oceans Institute, School of Biology, University of St Andrews, St Andrews, Fife KY16 8LB, UK

## Abstract

What shapes variation in genetic structure within a community of codistributed species is a central but difficult question for the field of population genetics. With a focus on the isolated coral reef ecosystem of the Hawaiian Archipelago, we assessed how life history traits influence population genetic structure for 35 reef animals. Despite the archipelago's stepping stone configuration, isolation by distance was the least common type of genetic structure, detected in four species. Regional structuring (i.e. division of sites into genetically and spatially distinct regions) was most common, detected in 20 species and nearly in all endemics and habitat specialists. Seven species displayed chaotic (spatially unordered) structuring, and all were nonendemic generalist species. Chaotic structure also associated with relatively high global  $F_{ST}$ . Pelagic larval duration (PLD) was not a strong predictor of variation in population structure ( $R^2 = 0.22$ ), but accounting for higher  $F_{ST}$  values of chaotic and invertebrate species, compared to regionally structured and fish species, doubled the power of PLD to explain variation in global  $F_{ST}$  (adjusted  $R^2 = 0.50$ ). Multivariate correlation of eight species traits to six genetic traits highlighted dispersal ability, taxonomy (i.e. fish vs. invertebrate) and habitat specialization as strongest influences on genetics, but otherwise left much variation in genetic traits unexplained. Considering that the study design controlled for many sampling and geographical factors, the extreme interspecific variation in spatial genetic patterns observed for Hawai'i marine species may be generated by demographic variability due to species-specific abundance and migration patterns and/or seascape and historical factors.

**Keywords:** chaotic genetic heterogeneity, community genetics, Hawai'i, marine connectivity, pelagic larval duration, stepping stone dispersal

Received 6 March 2014; revision received 8 May 2014; accepted 9 May 2014

## Introduction

The structuring of species into genetically distinct populations has many impacts on a species' demography and evolution (Kokko & López-Sepulcre 2007). In turn, ecological and environmental factors influence population genetic structuring (Avice 2000; Storfer *et al.* 2007). Understanding linkages between ecological, genetic and

environmental patterns is central to many current challenges in organismal biology and conservation (Taberlet *et al.* 2012). Uncovering generalities about these linkages requires comparison across multiple species, habitats and scales. Meta-analyses can test for meaningful relationships between genetic structuring and ecological traits across many species, but are hindered by the large number of possible confounding variables. In fact, an early finding of the rapidly expanding field of landscape genetics is that genetic structuring is highly species specific, influenced by the individual's interaction with landscape features according to life history and

Correspondence: Kimberly A. Selkoe, Fax: 805-892-2510;

E-mail: selkoe@nceas.ucsb.edu

§Author names and affiliations listed in Appendix I.

demographic factors, such that generalities may be few (Manel *et al.* 2003). Marine systems are known for harbouring diverse and often surprising spatial population genetic patterns (Selkoe *et al.* 2008). Here, we characterize the variation in population genetic structure across species within a single marine community, which share a basic habitat array, environmental gradients and key study sampling design elements. Further, we examine whether life history traits associate with genetic patterns, perhaps pointing to mechanisms maintaining the diversity in genetic patterns across species.

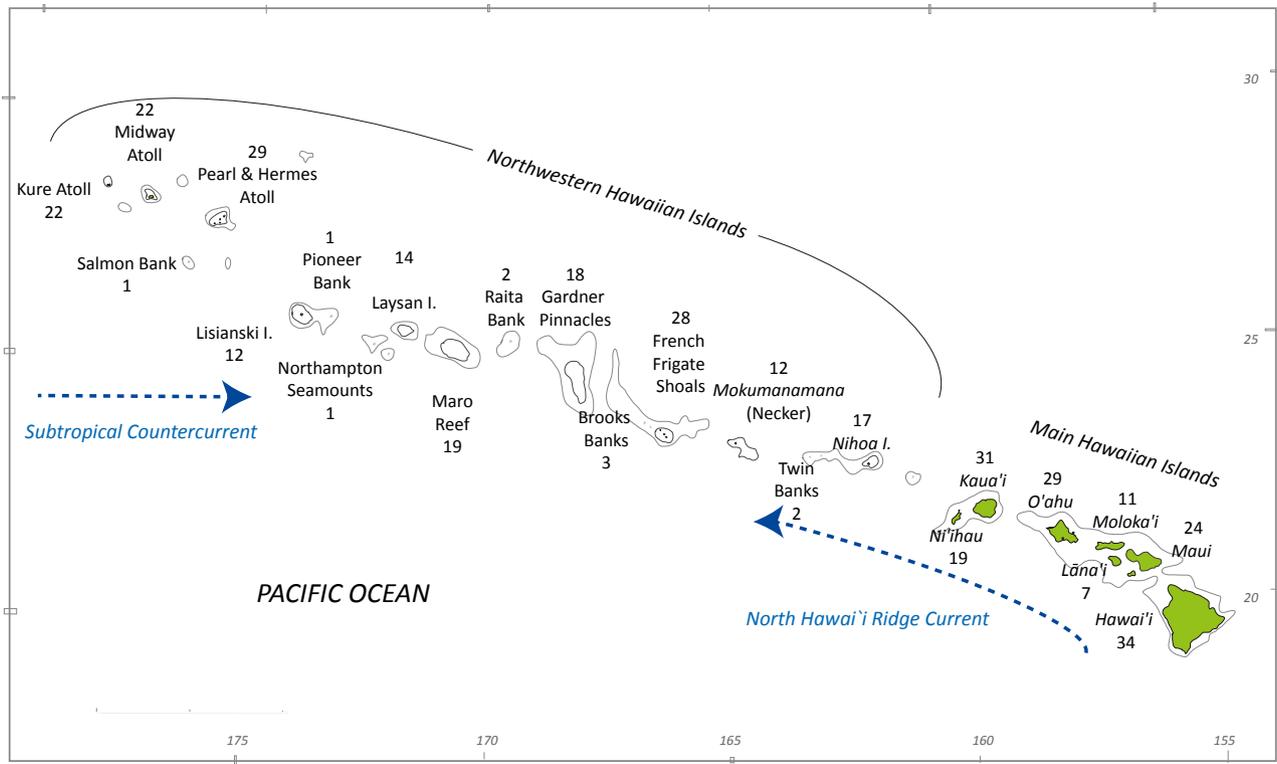
There is great interest in determining what drives spatial patterns of population genetics for marine species and the extent to which life history traits associate with particular types of genetic structuring. Theory suggests that the scale and pattern of genetic structure reflect long-term rates of gene flow driven primarily by migration, drift and selection. A longstanding focus of the field of population genetics is the relationship of dispersal potential to gene flow, because dispersal is a difficult trait to study directly but central to many basic and applied questions in ecology. Across studies of marine species, dispersal traits show significant correlation with genetic structuring, albeit often only weakly (Bradbury *et al.* 2008; Weersing & Toonen 2009; Riginos *et al.* 2011; Selkoe & Toonen 2011; Faurby & Barber 2012), leaving open the question of whether the remaining variation in genetic differentiation between populations may be explainable by factors such as taxonomy, life history, sampling design or historical effects. Despite hundreds of single-species marine population genetics studies across the globe, it is still unclear whether stronger or more coherent links between genetic and life history traits might emerge if variables such as history, habitat array or taxonomy could be constrained.

Two basic categories of population genetic structure are historically recognized the island model (discrete structuring in which individuals exist in genetically homogeneous 'islands' with limited gene flow between them) and the stepping stone model (continuously increasing differentiation along spatial gradients; Wright 1943). A third possibility is extensive dispersal and low genetic drift, whereby genetic differentiation is statistically insignificant (i.e. panmixia) and the entire geographic domain is genetically homogeneous. This is particularly likely in the ocean, where populations can be very large and migration is favoured by the high-dispersal medium (e.g. Theisen *et al.* 2008). A fourth model has emerged out of empirical marine genetics, 'chaotic' population structuring (Johnson & Black 1982; Hedgecock & Pudovkin 2011) in which the level of genetic structuring has been found to be highly variable with no obvious spatial patterning, possibly indicating nonequi-

librium conditions, sweepstakes recruitment (Hedgecock 1994; Arnaud-Haond *et al.* 2008), drift-dominated structuring (Johnson & Black 1982; Puritz & Toonen 2011; Broquet *et al.* 2012; Yearsley *et al.* 2013), unaccounted for seascape drivers or selection (Baums *et al.* 2006; Galindo *et al.* 2010; Selkoe *et al.* 2010; White *et al.* 2010; Foster *et al.* 2012), or some combination of all these factors (Toonen & Grosberg 2011). The relative frequencies of these four types of genetic structuring and their main drivers are unknown for marine ecosystems.

The present study leverages recent genetic studies of Hawaiian coral reef species to examine the range of population genetic patterns across reef animals within a single community and investigate whether species traits covary with metrics and models of genetic structure. The geography of the Hawaiian Archipelago provides an especially tractable system in which to study marine population structure, because it is the remotest archipelago in the world, composed of a nearly linear 2400-km-long array of discrete habitat patches of islands, atolls and seamounts. Insofar as possible in a natural system, these factors constrain the patterns of population structure, connectivity across patches and genetic history of populations (Fig. 1). We began by categorizing data sets for 37 coral reef species sampled with standard genetic markers at >5 islands. The focal species represent diverse taxa, but all are connected through the near-shore food web and share the same benthic and/or pelagic spaces as larvae and/or adults, thus can be considered a single reef community. Based on habitat array, we hypothesized that a stepping stone pattern of dispersal producing an isolation-by-distance pattern of spatial genetic structure would be prevalent. There are no obvious known physical barriers or strong oceanographic discontinuities that might lead to hierarchical genetic structuring. However, a precursor to this study found that locations of significant pairwise  $F_{ST}$  were highly variable across 27 species in Hawai'i, but occurred most commonly in the main Hawaiian Islands (MHI) where islands are more closely spaced (Toonen *et al.* 2011).

Riginos *et al.* (2011) outlined two approaches to studying life history effects on population structure: planned multispecies comparisons using a common sampling regime and geography, and post hoc compilation of published studies. This study represents a hybrid, in which compilation of raw genetic data sets for codistributed species within a single discrete study region allowed a large degree of standardization. One would expect that by controlling for basic habitat array, environmental gradients and many shared historical influences, much of the noise in this relationship might be eliminated. In this way, more nuanced multivariate influences on gene flow could emerge, enhancing our



**Fig. 1** Map of the Hawaiian archipelago. The number of species sampled per island or atoll is indicated next to each. Major currents are represented by arrows. 1000 and 2000 m isobaths delineated.

understanding of the feedbacks between life history, ecology and genetics.

By comparing life history and genetic data across a broad taxonomic range of species, we hope to gain insights into the mechanisms that drive geographic population structuring in marine systems. We characterize variation in population structure across species in two ways. First, we constrain the question by theory, evaluating how the established models of genetic structuring are represented by the 37 species. Second, we use unconstrained ordination to determine natural divisions among the data sets based on a suite of genetic metrics of spatial structure. These two approaches are complimentary in that the first is using significance testing of how spatial distributions of genetic diversity fit with a priori models, whereas the second is based on combinations of the metrics themselves, perhaps revealing additional divisions in the database which may not map well to the theory-based categories (i.e. if the values of the genetic metrics in each category show different ranges or variances).

Next, we test whether taxonomic, life history or sampling factors contribute to the observed variation in genetic structure. Using canonical analysis, we estimate how much variation in genetic metrics across species can be explained by available life history and taxonomic

traits. This broad-brush assessment is followed by alternative model testing of relationships between particular genetic and life history traits, to get insight into mechanisms driving structuring in this system. Previous empirical studies comparing large numbers of marine genetic data sets have reported that pelagic larval duration (PLD) shows positive correlation with genetic differentiation (Weersing & Toonen 2009; Selkoe & Toonen 2011) and that body size and depth preference show negative correlations with genetic differentiation (Bradbury *et al.* 2008; Kelly & Palumbi 2010; Riginos *et al.* 2011). We test each of these relationships here and also compare species with strong vs. weak genetic structuring, and Hawaiian endemics vs. widespread species, as endemism occurs at a high rate in Hawaii and could be associated with distinct genetic characteristics.

## Methods

### Data set preparation format

Data sets were assembled primarily from collections made on NOAA expeditions throughout the Hawaiian archipelago from 2005 to 2012 and subsequent publications of mitochondrial DNA (mtDNA) and nuclear DNA (usually microsatellite) data sets. Species chosen

tended to be abundant and easy to sample and identify. Genetic data sets were contributed in ARLEQUIN format (version 3.5.1.2, Excoffier & Lischer 2010) or an Excel format that allowed easy conversion to ARLEQUIN format. A modified version of PGDSPIDER (version 2.0.5.1, Lischer & Excoffier 2012) was used to convert between file formats for genetic analyses. For coral species only, GENETIX (version 4.05.2, Belkhir *et al.* 2002) was used to estimate and filter out clonal replicates within sites. ARLEQUIN files were modified to give all sites standardized four letter name codes and standardized ordering from SE to NW along the island chain. Because we have data at the scale of the island/atoll, we focus hypotheses at this spatial resolution. In most cases, allele/haplotype frequencies at adjacent islands are statically indistinguishable, indicating that island/atoll is an appropriate spatial scale for our study.

Any distinct subisland or subatoll samples were kept separate, with distinct names, except when  $F_{ST}$  was statistically indistinguishable from 0, in which case sublocalities were lumped. Several species showed samples collected in the vicinity of Kona to be distinct from those near Hilo on Hawai'i Island, and *Acanthaster planci* showed two distinct populations at Pearl and Hermes Atoll.

### Sampling filters

For inclusion in the analyses, a data set required at least five sites sampled with at least 10 individuals per site. For data sets that meet these criteria, sites with fewer than 10 individuals were also excluded. We also analysed results for a sample size minimum of 20 because  $F_{ST}$  can be inflated at small sample size, and allele frequency estimates are less reliable for low frequency alleles at highly polymorphic loci. Using a minimum sample size of 20 individuals per site excluded 90 of 533 samples in the data set using 10 or more samples per site (17% of samples). We comment below on how the two sampling filters affect results.

### Summary statistics

Nuclear loci with significant deviation from Hardy–Weinberg equilibrium were excluded before analysis. GENODIVE (version 2.0b23, Meirmans & van Tienderen 2004) was used to calculate estimates of global and pairwise  $F_{ST}$  based on Weir & Cockerham's (1984)  $\theta$ , with AMOVA using 9999 permutations. ARLEQUIN was used to calculate AMOVA based on  $\phi_{ST}$ , using AIC from jMODELTEST (version 2.1.4, Durrant *et al.* 2012) to choose the most appropriate mutational model in ARLEQUIN. SMOGD online calculator (Crawford 2010) was used to calculate  $D_{EST}$  and effective alleles (Jost 2008). GENOD-

IVE's  $K$ -means clustering was run for number of clusters ( $K$ ) from 1 to  $N-2$  using AMOVA-based simulated annealing with 50 000 steps and 20 repeats. Cluster membership was examined to determine whether adjacent sampling sites clustered together, highlighting where genetic boundaries (i.e. genetic discontinuities) between regions might exist. Genetic boundaries were considered where AMOVA estimation of  $F_{CT}$  across the boundary was statistically significant. The largest number of clusters of spatially discrete samples that returned significant  $F_{CT}$  results with AMOVA was recorded. In most cases, this was  $K = 2$  or 3. In some cases,  $K$ -means clustering showed slightly spatially mixed clustering. For example if Midway, a Northwestern Hawaiian Islands (NWHI) site, grouped with the MHI but otherwise MHI and NWHI sites were in two distinct clusters, a 'spatially strict' version of the clusters (e.g. Midway was placed in the NWHI cluster) was tested with AMOVA to confirm that  $F_{CT}$  values were significant after the regrouping. This procedure was only used when 1–2 samples were geographically incongruent in the clustering results. Clusters made up of spatially mixed samples were considered evidence that genetic structuring was not regionally organized. Pairwise geographic distance between sites based on coordinates was generated using GENODIVE. Isolation-by-distance analyses were generated using linearized  $F_{ST}$  [ $F_{ST}/(1-F_{ST})$ ] vs. Euclidean distance. Significance testing was based on Mantel tests with 999 replicates performed in GENODIVE.

Nine species were represented by two data sets, one using a mtDNA marker and a second using one or more nuclear markers (e.g. microsatellite panels or nuclear intron sequence). Genetic summary statistics were calculated for each marker class independently and then compared to gauge congruence. The mtDNA data set was preferentially chosen to represent the species in subsequent ordination analyses (which required one data set per species to avoid double counting), except where sampling power of the nuclear data set was superior, see results for details.

### Categories of spatial genetic structure

Based on the above summary statistics, data sets were placed in the following categories of spatial genetic structuring, summarized in Table 1:

- 1 Panmixia – defined as a lack of spatial genetic structuring, indicated here when global  $F_{ST}$ ,  $\phi_{ST}$  and  $D_{EST}$   $P > 0.05$ , spatial groupings based on  $K$ -means clustering show  $F_{CT}$   $P > 0.05$ , and isolation by distance (IBD) testing shows Mantel  $r$   $P > 0.05$ .
- 2 Chaotic genetic heterogeneity – defined as genetic differentiation of samples with no apparent spatial

**Table 1** Summary of the criteria used to categorize data sets by type of spatial genetic structuring. Significance tests used  $P < 0.05$  without correction for multiple tests

	Global $F_{ST}$ , $\phi_{ST}$ or $D_{EST}$ test significant?	Spatial clustering ( $F_{CT}$ ) significant?	IBD test significant?
1. Panmixia	No	No	No
2. Chaotic	Yes	No	No
4. IBD	Yes or no	No	Yes
5. Regional groups	Yes or no	Yes	Yes or no

organization, indicated here when global  $F_{ST}$ ,  $\phi_{ST}$  and/or  $D_{EST}$   $P < 0.05$ , but neither IBD nor any spatial clustering is statistically significant.

- 3 IBD – a significant IBD Mantel correlation ( $P < 0.05$ ) without significant spatial clustering, or within clusters, indicates auto-correlated spatial variation, regardless of the global tests of differentiation.
- 4 Regional genetic structure – when  $K$ -means clustering identified groupings of adjacent populations with  $F_{CT}$   $P < 0.05$ , regardless of IBD, and global differentiation.

It is possible that a species could conform to more than one category in different regions of the archipelago or when evaluated at different scales. Most data sets could not be properly evaluated for this possibility due to limited sampling. However, for every case of regional structure, we tested for the joint presence of IBD and regional groups (i.e. IBD for the sites within a cluster). As illustrated by Meirmans (2012), IBD and hierarchical structure can be confounded. Hierarchical structure can mimic IBD when differentiation within regions is low and distant pairs are cross-regional comparisons, whereas IBD can mimic hierarchical population structure if sampling is sparse and uneven. These scenarios were distinguished (albeit with low power in our case) using stratified Mantel tests in GENODIVE to permute the locations of populations within the clusters.

#### Clustering data sets by genetic summary statistics

The above categorizations are based on labelling data sets according to their fit with existing models of genetic structure derived from genetic theory. The designations are based on the statistical significance at  $\alpha = 0.05$  of a small number of genetic metrics. This approach ignores possibly useful information contained in the continuous range of values of the metrics themselves. It is also sensitive to sample size, which influences statistical significance. As an unconstrained

alternative, we conducted a principle components analysis (PCA) with JMP version 10 (SAS). These included all genetic summary statistics ( $F_{ST}$ ,  $\phi_{ST}$ ,  $D_{EST}$ ,  $F_{CT}$ , IBD  $r$  and the number of genetic regions; Table S1, Supporting information) to find natural divisions in the data sets that are unconstrained by any pre-existing labels or theory. Genetic metrics were linearized and log-transformed to homogenize scales prior to all analyses. Negative values of  $F_{ST}$ ,  $\phi_{ST}$  and  $D_{EST}$  were set to 0 to avoid a confounding influence on ordinations. PCA allowed us to visualize the main trends in summary statistics, ascertain redundancy in summary statistics and visualize natural breaks or clusters of data sets by genetic traits.

#### Life history data

Published literature and FishBase were searched for each species to gather basic life history data (Table 2). All life history traits available for a great majority of species were included, producing nine variables in the initial analyses. Estimates of mean PLD were available in the literature for 32 of the 37 species. To fill in missing values, a mean based on congeners ( $n = 16$ ) was used for the two *Chaetodon* spp. lacking PLD data and a mean based on confamilials ( $n = 7$ ) was used for the two groupers (family Serranidae). There is little information on tropical subtidal hermit crab PLDs. As crabs typically go through 4–6 larval stages lasting approximately a week (Lang & Young 1977), we estimated the mean PLD to be 50 days. The log transformation minimizes the effects of imprecise large values, and this one point has little leverage on the linear fit. Depth range (in m), maximum total length (body or colony size in cm) and estimates of generation time (in years) were available for all species and used on a continuous scale and  $\log_{10}$  transformed. Species were divided into habitat specialists and generalists. Generalists utilize sand, rubble or reef, whereas specialists are restricted to, or limited by, specific habitat features which may have smaller total area and/or distinct spatial arrays of habitat that differ greatly from the array of shallow habitat across the archipelago (e.g. damselfish requiring nesting sites, hermit crabs sheltering in certain corals, limpets limited to intertidal basalt which is patchy or absent at islands, corallivores requiring live coral). Five basic trophic categories were designated: corallivore, detritivore/sediment, invertivore, piscivore, algivore and planktivore, but analyses collapsed these into a binary categorization (algivore and planktivore vs. others) given the sample size of the data set. Other binary categorizations were examined: predator (invertivore and piscivore vs. others) and benthic feeders (corallivores, detritivore/sediments, invertivores and algivores vs.

**Table 2** Taxonomic and life history traits of 37 species used in the study

Genus species	Taxon	PLD	Range	Depth range	Max. length	Gen. Time	Eggs	Habitat	Trophic group
1. <i>Abudefduf abdominalis</i>	Pomacentridae	24	End	49	30	2.5	Att	S: rubble	Planktivore
2. <i>Abudefduf vaigiensis</i> *	Pomacentridae	20	IP	49	20	1	Att	G: reef	Planktivore
3. <i>Acanthurus nigrofuscus</i>	Acanthuridae	31	IP	25	20	2.5	Free	G: reef	Herbivore
4. <i>Acanthurus nigroris</i>	Acanthuridae	58	End	89	25	1	Free	G: reef	Herbivore
5. <i>Acanthurus olivaceus</i>	Acanthuridae	60	Pac	43	35	2.5	Free	G: reef	Herbivore
6. <i>Acanthaster planci</i>	Asterpodea	14	Pac	3	30	2	Free	G: reef	Corallivore
7. <i>Calcinus hazletti</i>	Decapoda	50	Pac	15	1	4	Att	S: coral	Detritivore
8. <i>Cellana exarata</i>	Gastropoda	6	End	2	7	3	Free	S: intertidal	Herbivore
9. <i>Cephalopholis argus</i>	Serranidae	28	IP	39	60	2.5	Free	S: high relief	Predator
10. <i>Chaetodon fremblii</i>	Chaetodontidae	40	End	61	13	1	Free	S: coral	Invertivore
11. <i>Chaetodon lunulatus</i>	Chaetodontidae	40	Pac	17	14	2.5	Free	S: coral	Corallivore
12. <i>Chaetodon miliaris</i>	Chaetodontidae	60	End	250	13	1	Free	G: reef	Omnivore
13. <i>Chaetodon multicinctus</i>	Chaetodontidae	40	End	109	12	1	Free	G: coral and rubble	Corallivore
14. <i>Ctenochaetus strigosus</i>	Acanthuridae	58	End	112	14	1	Free	G: coral and rubble	Herbivore
15. <i>Hyporhamphus quernus</i>	Serranidae	40	End	360	122	15	Free	S: high relief	Predator
16. <i>Etelis coruscans</i>	Lutjanidae	40	IP	157	120	8	Free	G: deep reef	Predator
17. <i>Etelis marshi</i>	Lutjanidae	40	IP	128	127	8	Free	G: deep reef	Predator
18. <i>Halichoeres ornatissimus</i>	Labridae	40	End	11	18	2.5	Free	S: coral	Invertivore
19. <i>Heterocentrotus mammillatus</i>	Echinoidea	8	IP	49	8	1	Free	G: reef	Herbivore
20. <i>Holothuria atra</i>	Holothuroidea	15	IP	30	60	1	Free	G: sand	Sediments
21. <i>Holothuria whitmaei</i>	Holothuroidea	15	IP	20	30	1	Free	G: sand	Sediments
22. <i>Lutjanus kasmira</i> *	Lutjanidae	31	IP	262	40	2.5	Free	G: reef	Predator
23. <i>Monitpora capitata</i>	Scleractinia	3	Pac	17	200	10	Free	G: reef	Planktivore
24. <i>Mulloidichthys flavolineatus</i>	Mullidae	60	IP	75	43	2.5	Free	G: sand and reef	Invertivore
25. <i>Mulloidichthys vanicolensis</i>	Mullidae	36	IP	112	38	2.5	Free	G: sand and reef	Invertivore
26. <i>Myripristis berndti</i>	Holocentridae	55	IP	12	30	1	Free	G: high relief	Planktivore
27. <i>Ophiocoma erinaceus</i>	Ophiuroidea	50	IP	27	20	1	Free	G: sand and reef	Sediments
28. <i>Ophiocoma pica</i>	Ophiuroidea	50	IP	27	10	1	Free	G: sand and reef	Sediments
29. <i>Panulirus marginatus</i>	Decapoda	365	End	142	40	4	Att	S: high relief	Invertivore
30. <i>Panulirus penicillatus</i>	Decapoda	270	IP	15	40	4	Att	S: reef and rock	Invertivore
31. <i>Parupeneus multifasciatus</i>	Mullidae	44	Pac	158	35	2.5	Free	G: sand	Invertivore
32. <i>Porites lobata</i>	Scleractinia	3	IP	23	200	8	Free	G: reef	Planktivore
33. <i>Pristipomoides filamentosus</i>	Lutjanidae	45	IP	360	100	5	Free	G: deep reef	Predator
34. <i>Stegastes fasciolatus</i>	Pomacentridae	30	End	29	16	1	Att	S: reef and rock	Herbivore
35. <i>Stenella longirostris</i>	Cetacea	0	IP	250	200	13	Int	S: all	Predator
36. <i>Triaenodon obesus</i>	Chondrichthyes	0	IP	32	200	8	Int	S: all	Predator
37. <i>Zebrasoma flavescens</i>	Acanthuridae	54	Pac	43	20	1	Free	S: reef	Herbivore

PLD, estimates of mean pelagic larval duration in days; End, endemic to Hawai'i; IP, Indo-Pacific wide; Pac, Pacific wide (including eastern Indian Ocean); depth given in metres; maximum length refers to body or colony size in cm; Gen. time, generation time or minimum doubling time in years; Att, eggs attached to substrate or body; Free, eggs spawned into the water column; Int, direct development; Habitat, habitat association; S, specialist; G, generalist categories.

\*Indicates recent arrivals to Hawai'i (<60 years); taxon lists family (fishes) or higher order (invertebrates, dolphin).

piscivores and planktivores) but provided no further insights into the analyses. Remaining life history categorizations were: endemic to Hawai'i vs. nonendemic and free-floating eggs vs. attached to body or substrate. Higher taxonomic affiliation was also used as a categorical variable (fish vs. invertebrate and dolphin) representing fundamental but unspecified characteristics that may tend to be shared across these highly diverse species, such as adult mobility (for which reliable data are lacking for key species). A PCA using the four continu-

ous variables and the five binary variables allowed us to assess the variation in the life history traits across species and visualize colinearities between life history traits, which were then confirmed with univariate linear regression or *t*-tests.

#### Redundancy analysis

Canonical analysis was used to assess how much the suite of life history traits explains the variation in

genetic traits as a whole and to visualize which traits most closely associate (Legendre & Legendre 2012). Redundancy analysis (RDA) is an ordination with regression; we used the package *VEGAN* in R for calculations (R code is available from Dryad doi:10.5061/dryad.1n246). The genetic metrics ( $Y$ ) are first transformed to  $Y'$  by fitting the values to a linear regression of each life history trait ( $X$ ). A PCA is then carried out on the  $Y'$  values. Colinearity of life history traits was examined before proceeding, leading to the elimination of generation time, which was correlated with maximum length but generally measured with much less precision (ordinary least squares  $r = 0.71$ ). The genetic summary statistics used were the same as described above for PCA (Table S1, Supporting information).  $F_{ST}$  was used in place of  $\phi_{ST}$  for the four nuclear marker data sets as missing data are not allowed in the analysis.  $F_{ST}$  and  $\phi_{ST}$  were correlated (OLS  $r = 0.71$ ), but both were included in the analysis to reveal differences in their responses to species traits, as this was our primary goal for the RDA instead of statistical hypothesis testing. All other genetic traits showed low colinearity. Two data sets with outlier  $F_{ST}$  values ( $F_{ST} > 0.2$ ) were removed (*Cellana exarata* and *Chaetodon lunulatus*) because outliers have disproportional influence on ordinations. We examined sampling factors as covariates in the analysis by performing a partial RDA with all factors (alleles, marker type, number of sites sampled, recent arrival species), but no effects were found. Adjusted  $R^2$  was calculated following the Ezekiel method (Legendre & Legendre 2012). The RDA triplot provided guidance on where to concentrate tests of particular associations of life history and genetic traits (i.e. it showed which traits have the strongest associations)

to avoid a large ratio of alternative models to sample size (Burnham & Anderson 2002).

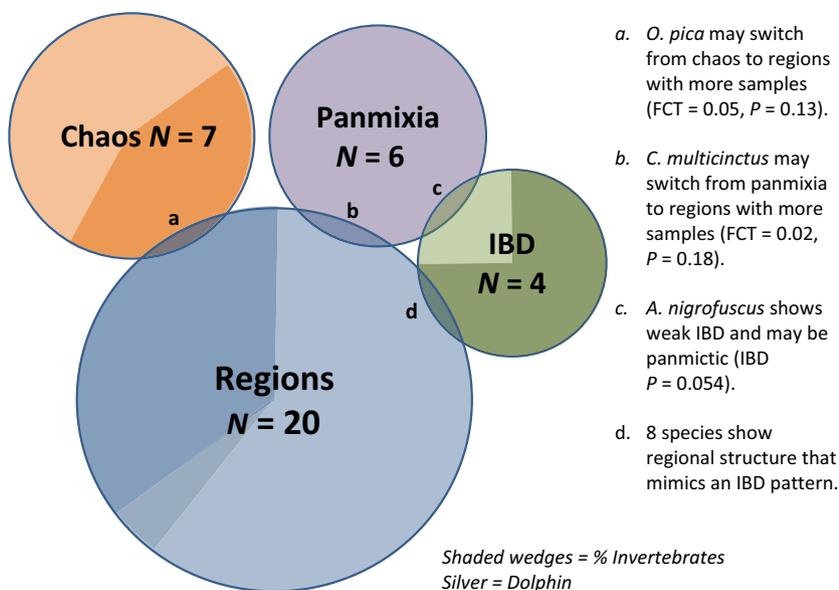
### Linear models of genetic differentiation

Based on the RDA results, correlations of several genetic and species traits were examined. Univariate correlations of continuous variables were made using ordinary least squares (OLS) regression, and  $t$ -tests were used for assessing significant association of genetic traits with categorical variables. Multivariate explanatory models combining categorical and continuous variables were made with generalized linear models (GLM) using normal distribution and identity link function in the software program *JMP*. Akaike's information criterion ( $AIC_C$ ) was used to select the most parsimonious models. The same two data sets were removed (*C. exarata* and *Ch. lunulatus*) because they were extreme outliers (i.e. their  $F_{ST}$  values were more than twice the value of the next highest values).

## Results

### Genetic categorizations

Data sets were divided into all four possible categories of genetic structuring: regional, IBD, chaotic and panmictic (Fig. 2, Table S1, Supporting information). Among species sampled with two marker types, five of the nine showed congruent categorization of both data sets. The remaining four species all had one data set with a low number of alleles that showed panmixia, and the other data set had a high number of alleles that showed structuring. The association of low polymor-



**Fig. 2** Venn diagram showing the categorization of 37 species into four models of population genetic structuring. Size of circles corresponds to number of species, shaded wedges correspond to proportion of invertebrates, and shaded sliver in largest circle indicates the dolphin data set. Overlapping edges of circles indicate grey areas where categorization of one or more data sets was borderline between the two models; each overlap is lettered and explained on right side.

phism with panmixia was the most pronounced sampling factor associated with data set categorization (Table S2, Supporting information). To avoid double counting species, we selected the data set with the more polymorphic marker(s) because of greater statistical power. For the congruent pairs, the mtDNA data sets were preferentially selected to increase consistency in marker type across data sets.

Regional grouping was the most common type of spatial genetic structure, observed in 20 of the 37 species. Eight species showed support for two spatial regions, eight for three regions, three for four regions and one for five regions (Table S1, Supporting information). In some cases, 'regions' comprised only one sample separated from others by a significant genetic break. In all but three of these species, hierarchical AMOVA showed no evidence of significant finer scale structuring within regions (i.e. significant  $F_{SC}$  values; exceptions were *Stenella longirostris* and *Montipora capitata*). Although 10 of the 20 regionally structured species showed significant overall IBD results (uncorrected  $P < 0.05$ ), none showed a significant stratified Mantel test, which would indicate IBD within regions. Thus, these IBD signals are likely an artefact of the regional structuring (the increased mean pairwise  $F_{ST}$  across regions compared to within regions), although for a minority, the stratified Mantel test may have lacked power to detect a true within-region IBD signal. Interestingly, 12 of the 20 regionally structured species showed global  $F_{ST}$  values not significantly different from 0. Regional data sets showed slightly less sampling coverage than other categories, and IBD data sets showed slightly more, especially in the MHI, although the mean difference in number of sites between the two is small and nonsignificant (Table S2, Supporting information).

Only four species were categorized as IBD, because they showed significant IBD after examining regional structuring. One of these, *Acanthurus nigrofuscus*, had a very weak signal ( $P = 0.05$ ) and most pairwise  $F_{ST} < 0$ . The other three species classified as IBD data sets were shallow invertebrates: a sea star (*Acanthaster planci*), a coral (*Porites lobata*) and a brittlestar (*Ophiocoma pica*).

Seven species categorized as chaotic showed highly significant global differentiation among sample sites and many significant pairwise  $F_{ST}$  values, but with no obvious spatial organization. However, one of these species, the brittlestar *Ophiocoma erinaceus*, showed a nearly significant IBD test ( $r = 0.42$ ,  $P = 0.07$ ) and nearly significant test for two regions ( $F_{CT} = 0.049$ ,  $P = 0.15$ ) that might have gained significance with more specimens. Two of the chaotic data sets yielded  $F_{ST}$  with  $P > 0.05$ , but  $D_{EST}$  and/or  $\phi_{ST}$  were highly significant.

The remaining six species were panmictic, with nonsignificant and very low global  $F_{ST}$ ,  $\phi_{ST}$  and  $D_{EST}$  val-

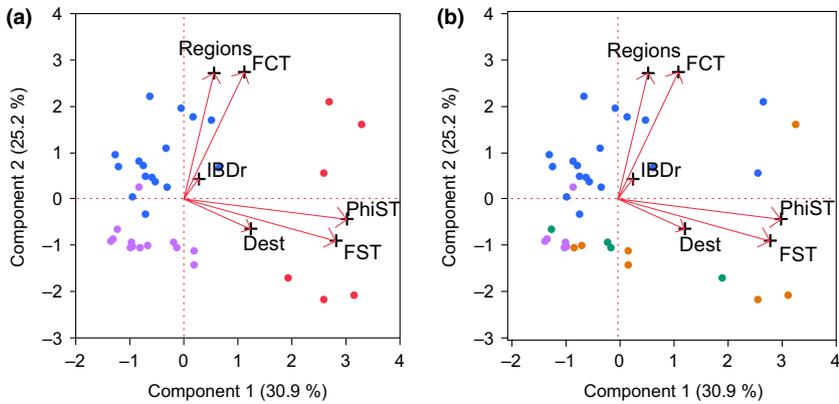
ues. Three of these had low allele counts such that their results may be considered inconclusive (surgeonfish *Acanthurus olivaceus*, slate urchin *Heterocentrotus mammilatus* and butterflyfish *Chaetodon multicinctus*).

Changing minimum sample size from 10 or more individuals per location to 20 or more affected the categorization of only one data set (butterflyfish *Chaetodon miliaris* lost four sites and switched from panmictic to regionally structured). Several other species lost enough sites to be excluded from analysis.

#### Ordinations of genetic and life history traits

The six genetic summary statistics ( $F_{ST}$ ,  $\phi_{ST}$ ,  $D_{EST}$ ,  $F_{CT}$ , number of significant regions, IBD  $r$ ) showed only moderate-to-low colinearity. The most correlated values were  $\phi_{ST}$  and  $F_{ST}$  (OLS  $r = 0.71$ ).  $D_{EST}$  was uncorrelated with  $\phi_{ST}$  and  $F_{ST}$ . A PCA using these six genetic summary statistics showed two data sets (limpet *Cellana exarata* and butterflyfish *Chaetodon lunulatus*) to be outliers to the rest because their values of  $\phi_{ST}$ ,  $F_{ST}$ ,  $D_{EST}$  and  $F_{CT}$  were much larger than the others (e.g.  $F_{ST} > 0.2$  vs.  $< 0.09$ , Table S1, Supporting information). These two data sets were removed from all further analyses, and the PCA was repeated to lessen the influence of skew on the analysis. The first four PCs showed eigenvalues  $> 1$  (Fig. 3a). The first PC, which showed high loadings for both  $\phi_{ST}$  and  $F_{ST}$ , separated out six data sets for which differentiation among sites is largest (e.g.  $F_{ST} > 0.02$ ; red markers in Fig. 3a). PC2 separated most of the data sets with one region and low  $F_{CT}$  values (purple in Fig. 3a) from those with multiple regions and high  $F_{CT}$  values (blue in Fig. 3a). PC3 was correlated with IBD  $r$  and PC4 with  $D_{EST}$ . To show how the four categories of genetic structuring map to the PCA results, the biplot is recoloured in Fig. 3b; it indicates that chaotic and IBD spatial organization are not clustered into a small range of values of  $F_{ST}$  or  $\phi_{ST}$ .

A PCA of the life history traits shows that species had diverse combinations of traits instead of a few clusters of associated types (see biplot Fig. S1, Supporting information). Pairwise linear regressions revealed two notable apparent correlations among life history traits. As previously well known, generation time and maximum length positively associate (OLS  $R^2 = 0.52$ ,  $P < 0.0001$ ). We excluded Generation Time from the RDA analysis due to colinearity. Also, fishes showed significantly broader depth ranges than invertebrates ( $R^2 = 0.30$ ,  $P < 0.0001$ ), but both were retained in the RDA as correlation was weak. The biplot shows that the genetic types contain a mixture of life history traits, but some tend to be absent from certain quadrants: for PC1 vs. PC2, chaotic data sets tend to be in the upper left (all were nonendemic and habitat generalists),



**Fig. 3** PCA biplot of six genetic summary statistics for 35 species ( $F_{ST}$ ,  $\phi_{ST}$ ,  $D_{EST}$ ,  $F_{CT}$ , number of regions, IBD fit). (a) Data sets are colour coded to show inherent clustering of data sets by genetic trait values (red = large  $F_{ST}$  values, purple = single region, blue = multiple regions). (b) Data sets are colour coded by genetic structure categorizations as in Fig. 2 (purple = panmixia, green = IBD, orange = chaos, blue = regions). PCA, principle components analysis.

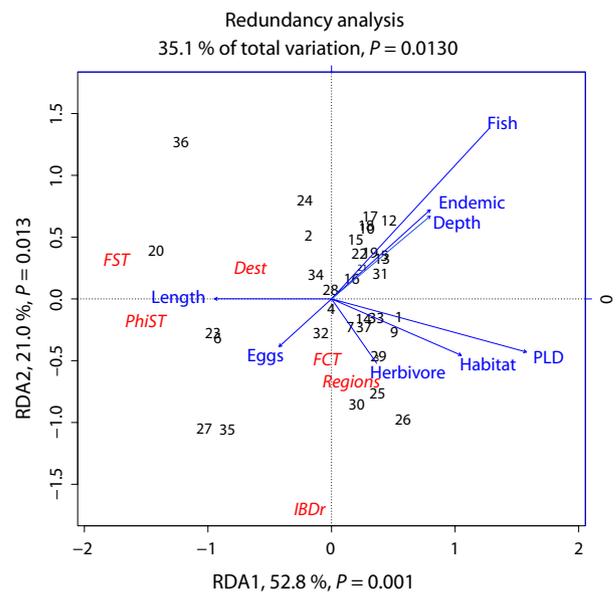
panmictic data sets in the upper half (broad and deep depth ranges and mostly fish) and IBD data sets in the lower half (shallow and invertebrate), whereas regional species are widely distributed over the plot (Fig. S1, Supporting information).

*RDA of life history and genetic traits*

The multivariate linear relationship between eight life history and six genetic traits was significant but not strong ( $R^2 = 0.35$ , adj.  $R^2 = 0.11$ ,  $P = 0.01$ ). Only the first PC showed eigenvector  $>1$ , indicating interpretability. Excluding  $\phi_{ST}$  to reduce redundancy with  $F_{ST}$  has no effect on results. The triplot is useful for visualizing which life history traits associate with genetic traits (Fig. 4). Three traits (PLD, fish, habitat) have longest vectors indicating strongest explanatory power. Three genetic traits (regions,  $D_{EST}$  and  $F_{CT}$ ) sit close to the centre of the ordination indicating that they are poorly explained by the life history traits. The vector for IBD is opposite the vectors of Fish, Endemic and Depth range, indicating negative relationships of these traits to stepping stone dispersal.  $F_{ST}$  aligns most closely with the PLD vector,  $D_{EST}$  weakly with the PLD vector, and  $\phi_{ST}$  is more influenced by Fish and Endemic than are  $F_{ST}$  and  $D_{EST}$ . As in Fig. 3a, the plot separates the six species with strong differentiation on the right side, away from the majority of other data sets. These species are a shark, a dolphin, a sea cucumber, a coral, a sea star and a brittlestar, a group encompassing all three types of spatial genetic structuring.

*Linear modelling of genetic traits*

The RDA indicates that Fish, Endemic and Depth have strong negative impacts on IBD  $r$ . A comparison of multivariate GLM models for these three traits and their interactions showed the most parsimonious model of IBD  $r$  includes Endemic and Depth only (adj.  $R^2 = 0.22$ ,

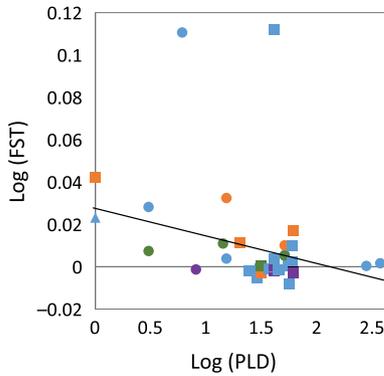


**Fig. 4** RDA triplot showing the associations of life history traits (blue font) with genetic traits (red font) for 35 species. Labels for Endemic and Depth are partially overlapping. Data points are represented by the species numbers, which correspond to names in Table 1. RDA, redundancy analysis.

$P = 0.006$ ); shallow, nonendemic invertebrates show stronger IBD patterns. Maximum depth, not minimum depth, drives the correlation with depth range.

A linear fit of PLD vs.  $F_{ST}$  was highly significant (OLS  $R^2 = 0.41$ ,  $P < 0.0001$ , Fig. 5), but when the two species (shark and dolphin) that lack larval development are excluded, the fit drops ( $R^2 = 0.16$ ,  $P = 0.019$ ). The fit strengthens slightly ( $R^2 = 0.19$ ) without the six panmictic data sets, for which  $F_{ST}$  is less informative because it is measured with larger error and is likely an artefact of low marker polymorphism for many of the data sets (Table S2, Supporting information).

For this subset of 27 nonpanmictic species with  $PLD > 0$ , a comparison of multivariate GLM models



**Fig. 5** Plot of PLD vs.  $F_{ST}$  for all species. Line excludes two outliers at top edge (*Cellana exarata*, *Chaetodon lunulatus*). Species are coded by their taxon (fish = squares, invertebrate = circles, dolphin = triangle) and type of spatial structuring (purple = panmixia, green = IBD, orange = chaos, blue = regions). PLD, pelagic larval duration.

based on  $AIC_c$  showed that a model with Fish and PLD is more parsimonious than a model of PLD alone (adj.  $R^2 = 0.32$ ,  $P = 0.002$ ,  $\Delta AIC_c = 4.1$ , Table 3). Adding Fish improved the fit because invertebrates have a higher intercept than fishes due to their generally higher  $F_{ST}$  values. Adding in the other eight species traits as additional factors does not improve the model (only individual additions were tested to minimize number of models compared and size of models given the small sample size). However, adding an indicator of whether the data set shows spatially organized structure (regional or IBD) or disorganized structure (chaotic) as a covariate improves the model significantly (adj.  $R^2 = 0.51$ ,  $P < 0.0001$ ,  $\Delta AIC_c = 6.6$ , Table 3). The model improve-

ment occurs because chaotic data sets have higher  $F_{ST}$  values on average and thus a higher intercept for PLD vs.  $F_{ST}$ . Categorization of data sets as regional or chaotic was made based on the  $P$ -value of  $F_{CT}$  in a hierarchical AMOVA, which shows no correlation with  $F_{ST}$  ( $R^2 = 0$ ,  $P = 1.0$ ). With the two direct developers (PLD = 0) included, the fit of this model is boosted (adj.  $R^2 = 0.63$ ). Interactions were not significant and thus excluded from the best fit model.

The 10 endemics show no relationship between PLD and  $F_{ST}$ , although all endemic PLD values are fairly large (PLD > 23 days). No species traits significantly explain  $F_{ST}$  for the endemics. Excluding endemics from the nonpanmictic, PLD > 0 group, results in the same best fit model, but with higher explanatory power (adj.  $R^2 = 0.58$ , Table 3).

We examined how GLM models explain variation in  $\phi_{ST}$  for the mtDNA data sets compared to those for  $F_{ST}$ . Fish and structure type (i.e. spatially organized vs. disorganized) without PLD best explained  $\phi_{ST}$  for the PLD > 0 nonpanmictic set of species (adj.  $R^2 = 0.55$ ,  $P < 0.002$ ).  $\phi_{ST}$  also shows no relationship to PLD for endemics, but shows a highly significant positive relationship to both maximum length and herbivory for endemics (adj.  $R^2 = 0.66$ ,  $P = 0.009$ ).  $D_{EST}$ ,  $F_{CT}$  and regions showed no significant linear relationships to the species traits collected for all species combined, as indicated by the RDA.

## Discussion

It is well known that marine species exhibit extensive variation in their genetic patterns that is poorly pre-

**Table 3** Upper: alternative model comparison for linear modelling of  $F_{ST}$  using 27 species (left side; two outliers excluded) and 20 nonendemic species (right side). Lower: detailed results of final GLM model.  $\Delta AIC_c$  = delta  $AIC_c$ , the difference in  $AIC_c$  value between the model and the top model

Parameter	Non-panmictic species with PLD > 0 ( $N = 27$ )				Endemics omitted ( $N = 20$ )			
	$P$ -value	$R^2$	Adj. $R^2$	$\Delta AIC_c$	$P$ -value	$R^2$	Adj. $R^2$	$\Delta AIC_c$
Fish + PLD + Structure	<0.0001	0.562	0.505	0	<0.0001	0.650	0.584	0.0
PLD + Structure	0.002	0.363	0.310	7.1	0.003	0.441	0.376	5.7
Fish + PLD	0.002	0.374	0.321	6.6	0.006	0.399	0.328	7.2
Fish	0.002	0.238	0.207	9.1	0.012	0.268	0.227	8.0
PLD	0.016	0.193	0.161	10.7	0.017	0.248	0.206	8.5
Structure	0.019	0.185	0.152	11.0	0.055	0.168	0.122	10.5
Term	Coefficient estimate			Std. error	$\chi^2$ $P$ -value			
Intercept	0.020146			0.00438	<0.0001			
Fish	0.004262			0.001218	0.0015			
PLD	-0.00751			0.002743	0.0101			
Structure	0.004766			0.001399	0.0019			

PLD, pelagic larval duration.

dicted by ecological or species traits. By focusing on a single isolated geographic region with a simplistic habitat array and calculating genetic metrics in a standardized way from raw data of many species, we investigated the extent to which variation in genetic patterns is constrained across species using two complementary approaches. The dual approaches help illuminate the extent to which our assessment is sensitive to the chosen metrics, categories and statistical framework. First, the PCA analysis focused on the strength of genetic differentiation and the spatial scale of structuring (i.e. number of regions) and revealed three clusters within the species set: single region/low differentiation, multiregional/moderate differentiation and a small number of high differentiation species with a mix of single and multiple regions. Second, the categorization of the data sets into four a priori types focused primarily on the spatial organization of the structuring and not the strength of differentiation (i.e. regional, IBD, chaotic and panmictic). Comparison of these two approaches revealed that species with strongest structuring show diverse spatial organization of structuring and likely a diversity of causes for that high structuring. Life history analyses revealed that chaotic species were all nonendemic and habitat generalists, IBD was most common for shallow invertebrate habitat generalists, regionally structured species showed a variety of life history associations, and panmixia was mostly limited to fishes with broad and deep depth ranges or associated with low allele number indicating low statistical power. Polymorphism creates precision much the way increasing the number of samples would (Kalinowski 2002).

#### *Regional boundaries across the archipelago*

The finding that regional structuring was most common, and IBD least common, was surprising given the stepping stone habitat array. Although every interisland channel along the chain was a possible boundary for at least one data set in the study, the most frequent site of a regional boundary, shared by 13 of the 20 regionally structured species, occurred at the centre of the archipelago, in the vicinity of French Frigate Shoals. This trend could lend insight into the factors enabling regional structuring for a diversity of taxa in this system. First, it might be possible to produce such a boundary in a stepping stone dispersal system with finite ends, because this would concentrate genetic differences on either end, especially when gene flow is relatively high, by elevating the importance of the increased drift at the edges (Rousset 2004). However, the stratified Mantel test results indicate that this scenario is unlikely, because IBD within regions was rare. Second, there

may be an oceanographic divergence zone at the centre of the chain. Larval dispersal might be biased away from the centre due to the eastern flowing Subtropical Countercurrent splitting as it encounters the archipelago, combined with the westerly North Hawai'i Ridge Current which may drive larvae westward (Fig. 1; Qiu *et al.* 1997; Kobayashi 2006). However, complex eddying indicated by meso-scale circulation modelling and simulated larval dispersal results suggest this scenario may also be overly simplistic and unlikely (Kobayashi 2006; Rivera *et al.* 2011). Finally, heterogeneity in demographic processes might drive departure from an IBD pattern, perhaps due to differences in habitat area or effective population size ( $N_e$ ) among locations. The sampling gaps in the data sets are a key consideration in this context. Despite a bias towards selecting study species that are abundant and easy to sample, many of the sampling gaps were caused by absence or very low density of organisms at sites. Thus, uneven abundance or density distributions across islands may lead to hierarchical structuring despite stepping stone dispersal for some species. Almost all endemics and all habitat specialists were regionally structured (except a few cases of panmixia). Both groups are more likely to have variable abundance across the chain due to spatially varying microhabitats, supporting this cause for regional structuring. This phenomenon begs for marine population genetic studies to carefully consider sampling design and run simulations to test effects of sampling factors on results. Additional factors and analysis approaches might add insight into the current results. For example, if ocean currents are important drivers of gene flow, the seasonal timing of larval dispersal, and other larval traits for which we were unable to find data, may help generate variation in genetic structuring across these species. We will explore the relative roles of history, oceanography, sampling and habitat factors in generating the observed variation in genetic patterns across species in future studies.

#### *Statistical considerations of genetic structure analysis*

Our data set proved to be a good example of the 'trouble with isolation by distance', described recently by Meirmans (2012). Nine data sets showed significant test results for IBD that on closer examination were driven only by regional structuring, evident both by examining site membership of data points on the IBD plot and by a stratified Mantel test. Our algorithm for categorizing a data set by its spatial genetic structuring was inspired by this study, and at least in the case of marine species, it appears that understanding whether the data set is spatially auto-correlated, regionally structured or chaotically structured is an important first step to interpreting

population genetic analyses. Meirmans (2012) found that 70% of a sample of studies testing for IBD found it. These were mostly terrestrial or aquatic studies. It is already known that marine species show much lower rates of IBD, and in this study, despite an uncommonly clear-cut stepping stone habitat array, only 10% of species showed IBD. It is unlikely sampling gaps influenced the categorization of data sets as IBD vs. regionally structured, because the average number of samples, number of sites in the MHI, NWHI and whole chain, and the size of the largest sampling gap were nearly identical for the categories (Table S2, Supporting information).

Many species showed global estimates of  $F_{ST}$  and  $\phi_{ST}$  near 0 despite strong regional structuring. Despite the fact that island-scale differentiation (global  $F_{ST}$ ) correlated with PLD, many species in the highest PLD category showed two or three genetically distinct regions, perhaps indicating that regional boundaries are not caused by dispersal related processes and instead may be a product of historical events (Marko 2004) and/or local adaptation. The  $K$ -means clustering approach to guide hierarchical AMOVA has not been widely used, but is more sensitive than a priori designation of groups. While it has the potential to uncover large-scale structure that is missed by other approaches (e.g. Kelly & Eernisse 2007; Díaz-Ferguson *et al.* 2010), it is also possible that the approach has inflated type 1 error.

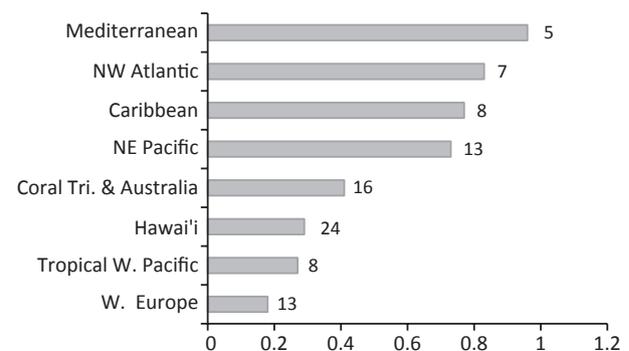
#### Variation in correlation of PLD and $F_{ST}$

We found that accounting for whether a data set is spatially organized improves insight into the relationship between genetic structure and species traits. It is interesting that the chaotic data sets showed a significant correlation of  $F_{ST}$  and PLD, but with higher mean  $F_{ST}$  values than species with spatially organized structure. The pattern suggests that these data sets are not chaotic simply because they are out of drift-migration equilibrium, but rather that they have an additional factor inflating differentiation. Consistent with this idea, recent simulation studies indicate that chaotic genetic patchiness can arise via small local effective population size and mildly aggregated dispersal of kin (Broquet *et al.* 2012), which may occur even in species with extremely long pelagic developmental periods (Iacchei *et al.* 2013).

The correlation of PLD and  $F_{ST}$  ( $R^2 = 0.22$  for species with pelagic larvae) was lower than the value for a global sample of studies ( $R^2 = 0.30$ ) derived from a variety of spatial scales, habitat configurations, regions and environmental settings (Selkoe & Toonen 2011). Furthermore, our sample of  $\phi_{ST}$  showed no significant relationship to PLD, instead correlating just with taxon,

consistent with  $\phi_{ST}$  having higher sensitivity to demographic history and mutation than  $F_{ST}$  (Bird *et al.* 2011a, b; Meirmans & Hedrick 2011). For any isolated marine habitat, retention strategies are crucial to persistence, but PLD may be less indicative of realized dispersal distance in this extremely isolated ecosystem compared to other places (but see Schultz & Cowen 1994; Robertson 2001). This possibility is supported by two additional insights. First, endemics, which may be more dependent on larval retention for persistence than widespread species, showed no relationship of  $F_{ST}$  to PLD. However, all endemics in the study have PLD > 20 days, after which the linear relationship saturates. Thus, retention of larvae may not be highly related to the mean PLD in this setting. Second, comparison of the correlation strength for Hawai'i to that found in other regions shows that the PLD vs.  $F_{ST}$  correlation is relatively weak in Hawai'i (Fig. 6).

Aside from poor correlation of PLD and dispersal distance, there are many other factors that can decouple PLD and  $F_{ST}$  values. One recent focus has been the influence of coalescent time on  $F_{ST}$ , because holding coalescent time constant should improve PLD vs.  $F_{ST}$  correlation (Dawson 2014). An analysis using hierarchical approximate Bayesian computation (Hickerson & Meyer 2008; Beaumont 2010) of our data set indicates nearly uniform timing and rate of expansion following the last glacial maximum (Y. Chan, K. Selkoe, H. Corpus & R. Toonen, in preparation). Therefore, we conclude that different coalescent histories are unlikely to



**Fig. 6** Log pelagic larval duration (PLD) vs. log  $F_{ST}$  correlation (OLS  $R^2$ ) for regional subsets of genetic studies sampled from the literature, which calculated  $F_{ST}$  using at least five sites per study. Data for all regions except Hawai'i taken from Selkoe and Toonen (2011).  $R^2$  value for Hawai'i used an analogous calculation and sample filtering of the present data set (i.e.  $F_{ST} < 0.001$  excluded, PLD = 0 included, and outliers *Cellana exarata* and *Chaetodon lumulatus* included). Numbers of studies used to calculate the  $R^2$  values of each region are indicated above columns. OLS, ordinary least squares.

drive variation in genetic patterns of these species. Contrary to claims that comparisons among synchronously diverging codistributed (SDC) species 'consistently evince higher gene flow in species with higher dispersal potential' (Dawson 2014), results from SDC taxa in Hawai'i (and previous global analyses; Selkoe & Toonen 2011) mandate a more nuanced treatment of the many forces impacting the population genetics of marine species.

We found that fishes show less structure and less organized structure than invertebrates. This point has not been previously highlighted, but the pattern is evident in other marine data sets (Carpenter *et al.* 2011; Selkoe & Toonen 2011; Toonen *et al.* 2011). Compared to invertebrates, fishes generally are more capable of behaviours that promote dispersal, and adult and juvenile migration is possible (Eble *et al.* 2011; Poortvliet *et al.* 2013). In addition, we found that species with deeper depth ranges tend to show less structuring than shallow species (also see Etter *et al.* 2005; Kelly & Palumbi 2010; Gaither *et al.* 2012; Andrews *et al.* 2014), perhaps because shallow habitat is smaller in total area, harder for larvae to intercept, and subject to more frequent disturbance, contributing to higher rates of genetic drift. The reason that only shallow invertebrates showed IBD may be due to the double constraint of limited dispersal ability and smaller habitat patch sizes.

In our analyses, egg type, body size and trophic group showed little influence on genetic traits. However, egg type and body size correlated with  $F_{ST}$  in other synthesis studies of marine species (Bradbury *et al.* 2008; Riginos *et al.* 2011). The great variation in findings for ecological correlates with  $F_{ST}$  and  $\phi_{ST}$  indicates that such syntheses are sensitive to the species composition and/or genetic markers in the data set, as well as the statistical approach. A shift of focus away from analysing trends in global  $F_{ST}$ , which is a relatively uninformative metric, towards a deeper understanding of what drives variation in spatial patterns of genetic differentiation will bring new insights into this line of inquiry (Lowe & Allendorf 2010; Marko & Hart 2011a,b).

#### *Multivariate analysis of species and genetic traits*

Canonical analysis of how eight species traits associated with genetic traits revealed that PLD, taxonomy (fish vs. invertebrate) and habitat specialization had the strongest influences on  $F_{ST}$ ,  $\phi_{ST}$  and IBD  $r$ , but overall explanatory power was low. Our use of RDA to uncover associations of life history and genetic traits followed a similar study of 27 codistributed high alpine plants of the European Alps (Meirmans *et al.* 2011). That study used AFLP data to generate eight genetic summary statistics describing spatial genetic diversity paired with six species traits

related to dispersal and habitat preference. The six species traits together explained a very similar fraction of variation in genetic traits relative to our finding ( $R^2 = 0.30$ , adj.  $R^2 = 0.17$ ). Considering the diversity of ecological, organismal and historical factors that can impact the distribution of genetic diversity, the authors interpreted this as a large fraction. Our data set included a wider diversity of species in terms of life history and taxonomy. Interestingly, Meirmans *et al.* (2011) analysis showed the same qualitative main results we report here. First,  $F_{ST}$  was the most strongly predicted trait and was driven by dispersal factors. Plants with multiple dispersal modes showed higher gene flow, similar to our finding that fishes, which can disperse both as adults and larvae, show higher gene flow than invertebrates. Second, IBD  $r$  was the only other strongly predicted genetic metric aside from  $F_{ST}$  in both studies, and rather than associating with dispersal factors as would be expected, IBD  $r$  was best predicted by habitat factors (soil type for plants, depth range for coral reef species). Historical factors (i.e. size and distribution of refugia) may drive both the depth and soil type effects, and retrospective analyses using coalescent models are needed to distinguish ancient connectedness from contemporary gene flow. A final similarity to Meirmans *et al.* (2011) is that Jost's  $D_{EST}$  showed no correlation with life history or other genetic traits, thus providing little insight in either context.

Despite the study design to limit sources of natural variance, the species included in this study showed great variation in genetic structure, and species traits were not highly predictive of that variation. Two of the species showed extremely high spatial structuring relative to all others, one regionally structured and one chaotically structured. Their exclusion from the analysis serves only to weaken the link between genetic variation and species traits. In sum, the question of what maintains the extreme diversity in spatial genetic patterns across marine species remains largely unanswered by this study, but is narrowed by the finding that it persists despite controlling for sampling design, marker type, habitat array, major environmental and oceanographic gradients and recent history to a greater extent than possible in meta-analyses of published works. Even within a single reef community, life history of marine species is extremely diverse and likely drives high diversity of demographic and genetic patterns.

Genetic diversity is a crucial foundation for biodiversity, with demonstrated influence on fitness, persistence, species diversity and ecosystem functioning (reviewed in Hughes *et al.* 2008; Taberlet *et al.* 2012). There is great interest in integrating population genetics into community ecology to understand the forces controlling community assembly and species interactions

(Avice 2000; Wares 2002; Cavender-Bares *et al.* 2009). Continuing to characterize the forces controlling spatial genetic structure in more marine and terrestrial communities and regions is an important first step.

## Acknowledgements

Funding was provided to K. Selkoe and R. Toonen by the National Science Foundation (BioOCE Award Number 1260169) and the National Oceanic and Atmospheric Administration (NMSP MOA#2005-008/66882). O. Gaggiotti was supported by the Marine Alliance for Science and technology for Scotland (MASTS). We are grateful to Heidi Lischer for modifying the software PGDSpider on our request, to the following for contributing ideas and feedback: Eric Crandall, Mary Donovan and Alan Friedlander; and to the following for contributing raw data: Giacomo Bernardi, Marina Ramon and Nicholas Whitney. Special thanks to Randall Kosaki for enabling critical field collections.

## References

- Andrews KR, Moriwake V, Wilcox C, Grau EG, Pyle RL, Bowen BW (2014) Phylogeographic analyses of submesophotic snappers *Etelis coruscans* and *Etelis "marshi"* (Family Lutjanidae) reveal concordant genetic structure across the Hawaiian Archipelago. *PLoS ONE*, **9**, e91665. doi:10.1371/journal.pone.0091665.
- Arnaud-Haond S, Vonau V, Rouxel C *et al.* (2008) Genetic structure at different spatial scales in the pearl oyster (*Pinctada margaritifera cumingii*) in French Polynesian lagoons: beware of sampling strategy and genetic patchiness. *Marine Biology*, **155**, 147–157.
- Avice JC (2000) *Phylogeography: The History and Formation of Species*. Harvard University Press, Cambridge, Massachusetts.
- Baums ILB, Miller M, Hellberg ME (2006) Geographic variation in clonal structure in a reef-building Caribbean coral, *Acropora palmata*. *Ecological Monographs*, **76**, 503–519.
- Beaumont MA (2010) Approximate Bayesian computation in evolution and ecology. *Annual Review of Ecology, Evolution, and Systematics*, **41**, 379–406.
- Belkhir K, Borsa P, Chikhi L, Raufaste N, Bonhomme F (2002) GENETIX 4.04, logiciel sous Windows™ pour la génétique des populations. Laboratoire Genome, Populations, Interactions. CNRS UMR 500, Université de Montpellier II, Montpellier, France.
- Bird CE, Smouse PE, Karl SA, Toonen RJ (2011a) Detecting and measuring genetic differentiation. In: *Crustacean Issues: Phylogeography and Population Genetics in Crustacea* (eds Koene-mann S, Held C, Schubart C), pp. 31–55. CRC Press, Boca Raton, Florida.
- Bird CE, Holland BS, Bowen BW, Toonen RJ (2011b) Diversification of sympatric broadcast-spawning limpets (*Cellana* spp.) within the Hawaiian archipelago. *Molecular Ecology*, **20**, 2128–2141.
- Bradbury IR, Laurel B, Snelgrove PVR, Bentzen P, Campana SE (2008) Global patterns in marine dispersal estimates: the influence of geography, taxonomic category and life history. *Proceedings of the Royal Society Series B, Biological Sciences*, **275**, 1803–1809.
- Broquet T, Viard F, Yearsley JM (2012) Genetic drift and collective dispersal can result in chaotic genetic patchiness. *Evolution*, **67**, 1660–1675.
- Burnham KP, Anderson DR (2002) *Model Selection and Multi-model Inference: A Practical Information-Theoretic Approach*. Springer-Verlag, New York, New York.
- Carpenter KE, Barber PH, Crandall ED *et al.* (2011) Comparative phylogeography of the coral triangle and implications for marine management. *Journal of Marine Biology*, **2011**, 1–14.
- Cavender-Bares J, Kozak KH, Fine PVA, Kembel SW (2009) The merging of community ecology and phylogenetic biology. *Ecology Letters*, **12**, 693–715.
- Crawford N (2010) SMOGD: software for the measurement of genetic diversity. *Molecular Ecology Resources*, **10**, 556–557.
- Darriba D, Taboada GL, Doallo R, Posada D (2012) jModelTest 2: more models, new heuristics and parallel computing. *Nature Methods*, **9**, 772.
- Dawson M (2014) Natural experiments and meta-analyses in comparative phylogeography. *Journal of Biogeography*, **41**, 52–65.
- Díaz-Ferguson E, Haney RA, Haney R *et al.* (2010) Population genetics of a trochid gastropod broadens picture of Caribbean Sea connectivity. *PLoS ONE*, **5**, e12675. doi:10.1371/journal.pone.0012675.
- Eble J, Toonen R, Sorenson L, Basch L, Papastamatiou Y, Bowen B (2011) Escaping paradise: larval export from Hawaii in an Indo-Pacific reef fish, the yellow tang *Zebrasoma flavescens*. *Marine Ecology Progress Series*, **428**, 245–258.
- Etter RJ, Rex MA, Chase MR, Quattro JM (2005) Population differentiation decreases with depth in deep-sea bivalves. *Evolution; International Journal of Organic Evolution*, **59**, 1479–1491.
- Excoffier L, Lischer HE (2010) Arlequin suite ver 3.5: a new series of programs to perform population genetics analyses under Linux and Windows. *Molecular Ecology Resources*, **10**, 564–567.
- Faurby S, Barber PH (2012) Theoretical limits to the correlation between pelagic larval duration and population genetic structure. *Molecular Ecology*, **21**, 3419–3432.
- Foster NL, Paris CB, Kool JT *et al.* (2012) Connectivity of Caribbean coral populations: complementary insights from empirical and modelled gene flow. *Molecular Ecology*, **21**, 1143–1157.
- Gaither MR, Toonen RJ, Bowen BW (2012) Coming out of the starting blocks: extended lag time rearranges genetic diversity in introduced marine fishes of Hawai'i. *Proceedings of the Royal Society Series B, Biological Sciences*, **279**, 3948–3957.
- Galindo HM, Pfeiffer-Herbert AS, McManus MA, Chao Y, Chai F, Palumbi SR (2010) Seascape genetics along a steep cline: using genetic patterns to test predictions of marine larval dispersal. *Molecular Ecology*, **19**, 3692–3707.
- Hedgecock D (1994) Temporal and spatial genetic structure of marine animal populations in the California Current. *CalCOFI Reports*, **35**, 73–81.
- Hedgecock D, Pudovkin AI (2011) Sweepstakes reproductive success in highly fecund marine fish and shellfish: a review and commentary. *Bulletin of Marine Science*, **87**, 971–1002.
- Hickerson MJ, Meyer CP (2008) Testing comparative phylogeographic models of marine vicariance and dispersal using a hierarchical Bayesian approach. *BMC Evolutionary Biology*, **8**, 322.

- Hughes AR, Inouye BD, Johnson MTJ, Underwood N, Vellend M (2008) Ecological consequences of genetic diversity. *Ecology Letters*, **11**, 609–623.
- Iacchi M, Ben-Horin T, Selkoe KA, Bird CE, García-Rodríguez FJ, Toonen RJ (2013) Combined analyses of kinship and FST suggest potential drivers of chaotic genetic patchiness in high gene-flow populations. *Molecular Ecology*, **22**, 3476–3494.
- Johnson MS, Black R (1982) Chaotic genetic patchiness in an intertidal limpet, *Siphonaria* sp. *Marine Biology*, **70**, 157–164.
- Jost L (2008) GST and its relatives do not measure differentiation. *Molecular Ecology*, **17**, 4015–4026.
- Kalinowski ST (2002) How many alleles per locus should be used to estimate genetic distances? *Heredity*, **88**, 62–65.
- Kelly RP, Eernisse DJ (2007) Southern hospitality: a latitudinal gradient in gene flow in the marine environment. *Evolution*, **61**, 700–707.
- Kelly RP, Palumbi SR (2010) Genetic structure among 50 species of the northeastern Pacific rocky intertidal community. *PLoS ONE*, **5**, e8594.
- Kobayashi D (2006) Colonization of the Hawaiian Archipelago via Johnston Atoll: a characterization of oceanographic transport corridors for pelagic larvae using computer simulation. *Coral Reefs*, **25**, 407–417.
- Kokko H, López-Sepulcre A (2007) The ecogenetic link between demography and evolution: can we bridge the gap between theory and data? *Ecology Letters*, **10**, 773–782.
- Lang WH, Young AM (1977) The larval development of *Clibanarius vittatus* (Bosc) (Crustacea: decapoda: diogenidae) reared in the laboratory. *Biological Bulletin*, **152**, 84–104.
- Legendre P, Legendre L (2012) *Numerical Ecology*. Elsevier B.V, Amsterdam.
- Lischer HEL, Excoffier L (2012) PGDSpider: an automated data conversion tool for connecting population genetics and genomics programs. *Bioinformatics*, **28**, 298–299.
- Lowe WH, Allendorf FW (2010) What can genetics tell us about population connectivity? *Molecular Ecology*, **19**, 3038–3051.
- Manel S, Schwartz MK, Luikart G, Taberlet P (2003) Landscape genetics: combining landscape ecology and population genetics. *Trends in Ecology & Evolution*, **18**, 189–197.
- Marko P (2004) 'What's larvae got to do with it?' Disparate patterns of post-glacial population structure in two benthic marine gastropods with identical dispersal potential. *Molecular Ecology*, **13**, 597–611.
- Marko PB, Hart MW (2011a) Retrospective coalescent methods and the reconstruction of metapopulation histories in the sea. *Evolutionary Ecology*, **26**, 291–315.
- Marko PB, Hart MW (2011b) The complex analytical landscape of gene flow inference. *Trends in Ecology & Evolution*, **26**, 448–456.
- Meirmans PG (2012) The trouble with isolation by distance. *Molecular Ecology*, **21**, 2839–2846.
- Meirmans PG, Hedrick PW (2011) Assessing population structure: F(ST) and related measures. *Molecular Ecology Resources*, **11**, 5–18.
- Meirmans P, van Tienderen PH (2004) GENOTYPE and GENODIVE: two programs for the analysis of genetic diversity of asexual organisms. *Heredity*, **93**, 43–50.
- Meirmans PG, Goudet J, Gaggiotti OE (2011) Ecology and life history affect different aspects of the population structure of 27 high-alpine plants. *Molecular Ecology*, **20**, 3144–3155.
- Poortvliet M, Longo GC, Selkoe KA et al. (2013) Phylogeography of the California sheephead, *Semicossyphus pulcher*: the role of deep reefs as stepping stones and pathways to antitropicality. *Ecology and Evolution*, **3**, 4558–4571.
- Puritz JB, Toonen RJ (2011) Coastal pollution limits pelagic larval dispersal. *Nature Communications*, **2**, 226.
- Qiu B, Koh D, Lumpkin C, Flament P (1997) Existence and formation mechanism of the North Hawaiian Ridge Current. *Journal of Physical Oceanography*, **27**, 431–444.
- Riginos C, Douglas KE, Jin Y, Shanahan DF, Trembl EA (2011) Effects of geography and life history traits on genetic differentiation in benthic marine fishes. *Ecography*, **34**, 566–575.
- Rivera MAJ, Andrews KR, Kobayashi DR et al. (2011) Genetic analyses and simulations of larval dispersal reveal distinct populations and directional connectivity across the range of the Hawaiian Grouper (*Epinephelus quernus*). *Journal of Marine Biology*, **2011**, 1–11.
- Robertson DR (2001) Population maintenance among tropical reef fishes: inference from small island endemics. *Proceedings of the National Academy of Sciences USA*, **98**, 5667–5670.
- Rousset F (2004) *Genetic Structure and Selection in Subdivided Populations*. Princeton University Press, Princeton, New Jersey.
- Schultz ET, Cowen RK (1994) Recruitment of coral-reef fishes to Bermuda: local retention or long-distance transport? *Marine Ecology Progress Series*, **109**, 15–28.
- Selkoe KA, Toonen RJ (2011) Marine connectivity: a new look at pelagic larval duration and genetic metrics of dispersal. *Marine Ecology Progress Series*, **436**, 291–305.
- Selkoe KA, Henzler CM, Gaines SD (2008) Seascape genetics and the spatial ecology of marine populations. *Fish and Fisheries*, **9**, 363–377.
- Selkoe KA, Watson JR, White C et al. (2010) Taking the chaos out of genetic patchiness: seascape genetics reveals ecological and oceanographic drivers of genetic patterns in three temperate reef species. *Molecular Ecology*, **19**, 3708–3726.
- Storfer A, Murphy MA, Evans JS et al. (2007) Putting the "landscape" in landscape genetics. *Heredity*, **98**, 128–142.
- Taberlet P, Zimmermann NE, Englisch T et al. (2012) Genetic diversity in widespread species is not congruent with species richness in alpine plant communities. *Ecology Letters*, **15**, 1439–1448.
- Theisen TC, Bowen BW, Lanier W, Baldwin JD (2008) High connectivity on a global scale in the pelagic wahoo, *Acanthocybium solandri* (tuna family Scombridae). *Molecular Ecology*, **17**, 4233–4247.
- Toonen RJ, Grosberg RK (2011) Causes of chaos: spatial and temporal genetic heterogeneity in the intertidal anomuran crab *Petrolisthes cinctipes*. In: *Phylogeography and Population Genetics in Crustacea* (eds Koenemann S, Held C, Schubart C) Ch. 4, pp. 75–107. CRC Press Crustacean Issues Series, ISBN 1439840733.
- Toonen RJ, Andrews KR, Baums IB et al. (2011) Defining boundaries for ecosystem-based management: a multispecies case study of marine connectivity across the Hawaiian Archipelago. *Journal of Marine Biology*, **2011**, 1–13.
- Wares JP (2002) Community genetics in the Northwestern Atlantic intertidal. *Molecular Ecology*, **11**, 1131–1144.
- Weersing K, Toonen R (2009) Population genetics, larval dispersal, and connectivity in marine systems. *Marine Ecology Progress Series*, **393**, 1–12.

- Weir BS, Cockerham CC (1984) Estimating *F*-statistics for the analysis of population structure. *Evolution*, **38**, 1358–1370.
- White C, Selkoe KA, Watson J, Siegel DA, Zacherl DC, Toonen RJ (2010) Ocean currents help explain population genetic structure. *Proceedings of the Royal Society Series B, Biological Sciences*, **277**, 1685–1694.
- Wright S (1943) Isolation by distance. *Genetics*, **28**, 114–138.
- Yearsley JM, Viard F, Broquet T (2013) The effect of collective dispersal on the genetic structure of a subdivided population. *Evolution*, **67**, 1649–1659.

---

K.A.S., R.J.T. and O.G. designed research, K.A.S. performed research, K.A.S. analysed data and K.A.S., R.J.T., O.G. and B.B. wrote the manuscript. Tobo Laboratory authors contributed data, analysed data and edited the manuscript.

---

### Data accessibility

Online supplemental materials (DRYAD doi:10.5061/dryad.1n246) include:

- R code for the RDA analysis ('RDA\_in\_revised.R')
- Excel file ('Selkoe\_et\_al\_Hawaii\_genetics\_dataset\_5.8.14.suppl.xls') with all other supplementary materials in the following sheets: 1 – assembled life history and genetic data with references; 2 – reference list; 3 – input table used in the RDA analysis; 4 – Table S1; 5 – Table S2; 6 – Figure S1.
- Genetic data sets will be contributed to the Indo-Pacific Research Network Data Repository located at <http://indo-pacific.wikispaces.com/>

### Appendix I

#### ToBo Laboratory authors and addresses

**Kimberly Andrews**, School of Biological Sciences, Durham University, South Road, Durham, DH1 3LE, UK; **Moisés A. Bernal**, California Academy of Sciences, 55 Music Concourse Drive, Golden Gate Park, San Francisco, CA 94118, USA; **Christopher Bird**, Marine Biology Program, Department of Life Sciences, Texas A & M University–Corpus Christi, 6300 Ocean Drive, Corpus Christi, Texas 78412, USA; **Holly Bolick**, Bishop Museum, 1525 Bernice St, Honolulu, HI, 96817, USA; **Iliana Baums**, Department of Biology, The Pennsylvania

State University, 208 Mueller Laboratory University Park, PA, 16802, USA; **Richard Coleman**, Hawai'i Institute of Marine Biology, University of Hawai'i, Kāne'ohe, HI 97644, USA; **Gregory T. Concepcion**, Pacific Biosciences, 1380 Willow Rd, Menlo Park, CA, 94025, USA; **Matthew T. Craig**, Department of Marine Science and Environmental Studies, University of San Diego, 5998 Alcalá Park, San Diego, CA 92110, USA; **Joseph D. DiBattista**, Red Sea Research Center, King Abdullah University of Science and Technology, Thuwal, Saudi Arabia; **Jeffrey Eble**, Center for Environmental Bioremediation and Diagnostics, University of West Florida, Pensacola, FL, 32561, USA; **Iria Fernandez-Silva**, California Academy of Sciences, 55 Music Concourse Drive, Golden Gate Park, San Francisco, CA 94118, USA; **Michelle R. Gaither**, School of Biological and Biomedical Sciences, Durham University, South Road, Durham DH1 3LE, UK and California Academy of Sciences, Ichthyology, 55 Music Concourse Drive, San Francisco, CA 94118, USA; **Mathew Iacchei**, Department of Oceanography, University of Hawai'i at Mānoa, 1000 Pope Rd., Honolulu, HI, 96822, USA; **Nicholas R. Polato**, Department of Ecology & Evolutionary Biology, 215 Tower Rd., Cornell University, Ithaca, NY 14853, USA; **Malia Ana J. Rivera**, Hawai'i Institute of Marine Biology, University of Hawai'i, Kāne'ohe, HI 97644, USA; **Luiz A. Rocha**, California Academy of Sciences, 55 Music Concourse Drive, Golden Gate Park, San Francisco, CA 94118, USA; **Derek Skillings**, Hawai'i Institute of Marine Biology, University of Hawai'i, Kāne'ohe, HI 97644, USA; **Molly Timmers**, Hawai'i Institute of Marine Biology, University of Hawai'i, Kāne'ohe, HI 97644, USA; **Zoltan Szabo**, Hawai'i Institute of Marine Biology, University of Hawai'i, Kāne'ohe, HI 97644, USA.

### Supporting information

Additional supporting information may be found in the online version of this article.

**Table S1** Basic genetic results for the subset of datasets used in the analyses (excludes duplicate datasets; see Dataset S1).

**Table S2** Sampling statistics by category of structure for all datasets.

**Fig. S1** PCA biplot of life history traits for all species ( $n = 37$ ).