# MOLECULAR ECOLOGY

# Local inter-species introgression is the main cause of outlying levels of intra-specific differentiation in mussels

Journal:	Molecular Ecology
Manuscript ID:	MEC-15-0353
Manuscript Type:	Original Article
Date Submitted by the Author:	31-Mar-2015
Complete List of Authors:	Fraisse, Christelle; CNRS-Université de Montpellier 2, ; University of Cambridge, Genetics Belkhir, Khalid; CNRS-Université de Montpellier 2, Welch, John; University of Cambridge, Genetics Bierne, Nicolas; CNRS-Université de Montpellier 2,
Keywords:	Genome scan, Adaptation, Selection, Outlier loci, Target enrichment sequencing, BAC assembly



1

5

**Original Articles** 

# Local inter-species introgression is the main cause of outlying levels of intra-specific differentiation in mussels

- Christelle Fraïsse<sup>1,2,3</sup>, Khalid Belkhir<sup>1</sup>, John J. Welch<sup>3</sup>, Nicolas Bierne<sup>1,2</sup>
- 6 1. Université Montpellier 2, Institut des Sciences de l'Évolution, UMR 5554, Montpellier Cedex 05,
- 7 France.
- <sup>8</sup> 2. CNRS, Institut des Sciences de l'Évolution, UMR 5554, Station Méditerranéenne de
- <sup>9</sup> l'Environnement Littoral, CNRS, Sète, France.
- <sup>10</sup> 3. Department of Genetics, University of Cambridge, Downing Street, Cambridge, UK.

# 11 Keywords:

<sup>12</sup> Genome scan, adaptation, selection, outlier loci, target enrichment sequencing, BAC assembly

# **13** Author for correspondence:

14 Christelle Fraïsse, e-mail: Christelle.Fraisse@univ-montp2.fr

# **15 Running title:**

<sup>16</sup> Local introgression in mussels.

#### 17

# 18 Abstract

Structured populations, and replicated zones of contact between species, are an ideal opportunity 19 to study regions of the genome with unusual levels of differentiation; and these can illuminate the 20 genomic architecture of species isolation, and the spread of adaptive alleles across species ranges. 21 Here, we investigated the effects of gene flow on divergence and adaptation in the Mytilus complex 22 of species, including replicated parental populations in quite distant geographical locations. We used 23 target enrichment sequencing of 1269 contigs of a few Kb each, including some genes of known func-24 tion, to infer gene genealogies at a small chromosomal scale. We show that geography is an important 25 determinant of the genome-wide patterns of introgression in *Mytilus*, and that gene flow between dif-26 ferent species, with contiguous ranges, explained up to half of the intra-specific outliers. This suggests 27 that local introgression is both widespread and tends to affect larger chromosomal regions than purely 28 intraspecific processes. We argue that this situation might be common, and this implies that genome 29 scans should always consider the possibility of introgression from sister species, unsampled differ-30 entiated backgrounds, or even extinct relatives, e.g. Neanderthals in humans. The hypothesis that 31 reticulate evolution over long periods of time contributes widely to adaptation, and to the spatial and 32 genomic reorganisation of genetic backgrounds, needs to be more widely considered in order to make 33 better sense of genome scans. 34

35

# 36 Introduction

The literature now contains many genome-wide surveys of differentiation, in a wide variety of systems (e.g. *Anopheles*, Turner *et al.* 2005; *Ficedula*, Ellegren *et al.* 2012; *Heliconius*, Martin *et al.* 2013; *Helianthus*, Renaut *et al.* 2013; *Corvus*, Poelstra *et al.* 2014; *Timena*, Soria-Carrasco *et al.* 2014; *Gasterosteus*, Jones *et al.* 2012). One of the most striking and consistent results is the heterogeneity of differentiation across the genome, including highly differentiated regions, sometimes called "genomic islands of differentiation" (Turner *et al.* 2005).

Several theories have been proposed to explain this pattern. The most prominent involves specia-43 tion with gene flow, driven by local adaptation (Via & West 2008; Nosil et al. 2009), but simulation 44 models suggest that genomic islands arise only in restricted biological conditions (Feder & Nosil 2010; 45 Feder et al. 2012; Flaxman et al. 2013). Other theories invoke background selection and hitchhiking 46 in closely related species (Noor & Bennett 2009; Roesti et al. 2012; Renaut et al. 2013; Yeaman 2013; 47 Cruickshank & Hahn 2014), the sorting of shared ancestral variation (Nielsen & Wakeley 2001), or 48 variable persistence after secondary contact of differences accumulated in allopatry (e.g. Fraïsse et 49 al. 2014a). Distinguishing between these scenarios is difficult, because the origins of semipermeable 50 genetic barriers to gene flow are intrinsically difficult to trace (Endler 1977; Barton & Hewitt 1985; 51 Harrison 1986). 52

Interpreting of regions of increased differentiation is also linked to questions regarding the origins of adaptive genotypes (Bierne *et al.* 2013; Roesti *et al.* 2014; Welch & Jiggins 2014). Barrier loci will delay the introgression of neutral alleles in proportion to their linkage (Barton & Bengtsson 1986; Charlesworth *et al.* 1997), but universally advantageous alleles can usually cross species barriers without much delay (Piálek & Barton 1997). As such, hybridization could lead to adaptive introgression of complex co-adapted haplotypes, and candidates have been reported in plants (Arnold 2004) and animals (Hedrick 2013), including humans (Mendez *et al.* 2012, 2013; Huerta-Sanchez *et al.* 2014).

In all of these cases, inferences are strengthened by the fact that the adaptation was not specieswide, allowing the researchers to focus on local introgression, and the abnormal differentiation of single populations (e.g. Europeans *vs* Africans, or Tibetans *vs* Hans in the study of human introgression from extinct relatives, Huerta-Sanchez *et al.* 2014). Nevertheless, even in such cases it is usually necessary to complement genome scans with surveys of genetic differentiation along small-scale chro-

mosomal regions (e.g. in mussels, Bierne 2010), reconstructions of the historical divergence of alleles
at candidate genes (e.g. in mice, Domingues *et al.* 2012) and ultimately experimental tests (e.g. rodent
poison resistance in mice, Song *et al.* 2011).

Here, we performed target enrichment sequencing in the *Mytilus* complex of species, focussing 68 on long anonymous regions obtained from BAC sequencing and cDNAs from databases and a tran-69 scriptome survey (Romiguier et al. 2014). The Mytilus complex includes three incompletely isolated 70 species of marine mussels, Mytilus edulis, Mytilus galloprovincialis and Mytilus trossulus. These 71 species have experienced a complex history of divergence with periods of gene exchange during the 72 Ouaternary (Roux et al. 2014). They lie along a gradient of genetic divergence: M. edulis and M. 73 galloprovincialis diverged 2.5 million years ago (Roux et al. 2014), while their divergence from M. 74 trossulus is estimated at 3.5 million years (Rawson & Hilbish 1995). Today, these species are in con-75 tact in several places in the northern hemisphere (Figure 1), and replicated parental populations are 76 found in quite distant geographical locations. 77

We took advantage of the original spatial genetic structure of the *Mytilus* mussels to explore the 78 consequences of local inter-species introgression on patterns of differentiation within and between 79 species. In Europe, a mosaic hybrid zone between M. edulis and M. galloprovincialis extends from the 80 Mediterranean Sea to the North Sea (Figure 1A, Bierne et al. 2003; Hilbish et al. 2012; Ouesada et al. 81 1995b), enabling the investigation of varying levels of inter-species introgression along a geographical 82 gradient. Equally important, natural replication of contact zones between M. edulis and M. trossulus 83 (one in Europe, Figure 1B, and one in North America, Figure 1C, Riginos & Cunningham 2005) 84 gave us the opportunity to study how the genomic architecture of species isolation varies in space. 85 Finally, we gain further insights into the history of adaptation within species by studying regions of the 86 genome with abnormal levels of differentiation. Careful analyses of genetic differentiation and gene 87 genealogies at a small chromosomal scale revealed that local introgression contributes significantly 88 to intra-specific outliers in mussels. Such outliers would have been misinterpreted had the analysis 89 considered only a single species. As it is common to address the question of the origin of adaptations 90 by scanning genomes, we argue that introgression from sister species or unsampled differentiated 91 backgrounds should be considered whenever possible. 92

# **Materials and Methods**

# 94 Sampling

The *Mytilus edulis* species complex comprises three species that hybridize at several places in the 95 northern hemisphere. We sampled individuals from eleven locations from both sides of the Atlantic 96 Ocean (Figure 1 and TableS1). Sampling took place outside hybrid zones in recognized patches of 97 panmictic populations. To investigate how patterns of genetic differentiation evolved along a gradient 98 of genetic divergence, we made use of the original genetic structure of the mussel complex of species. 99 Within M. galloprovincialis, previous studies have reported a genetic break either side of the Almeria-100 Oran front between Atlantic and Mediterranean populations (Quesada et al. 1995a,b). Samples from 101 two localities along the eastern Atlantic coast were obtained (the Iberian Coast and Brittany) as well 102 as samples either side of the Sicilo-Tunisian strait in the Mediterranean Sea (in Sete, France, for the 103 western basin and Crete for the eastern basin). The two closest species of the mussel complex, M. 104 edulis and M. galloprovincialis, meet along the French coast in a well-studied mosaic hybrid zone 105 characterized by three successive transition zones (Bierne et al. 2003; Hilbish et al. 2012). To inves-106 tigate gradients of introgression between them, we sampled populations enclosed within the mosaic 107 zone (the Bay of Biscay for M. edulis and Brittany for M. galloprovincialis) together with external 108 populations (the North Sea for M. edulis and the Iberian Coast for M. galloprovincialis). At the ex-109 treme of the gradient of divergence, M. edulis and M. trossulus met on two independent occasions in 110 the northern Atlantic (in Europe and in North America, Riginos & Cunningham 2005), giving us the 111 opportunity to study the outcomes of replicated contacts. In Europe, the two species meet in a clinal 112 hybrid zone at the entrance of the Baltic Sea (Väinölä & Hvilsom 1991). Individuals of M. trossulus 113 were sampled at the bottom of the Baltic Sea, in the gulf of Finland. In North America, the zone of 114 contact between the two species extends from Maine to Nova Scotia (Koehn et al. 1984). M. trossulus 115 mussels were obtained from the Saint-Lawrence river in Canada. M. edulis mussels were sampled in 116 Rhode Island (USA). In total, eight individuals per sample were examined, except for the American 117 *M. edulis* sample which comprised 11 individuals (see TableS1), the *M. trossulus* sample from Saint-118 Lawrence river in which a *M. edulis* individual was found (see TableS3), and the *M. galloprovincialis* 119 sample from the Iberian Coast for which two individuals had to be removed for technical reasons (low 120 coverage, see below and TableS3). Genomic DNA was extracted from adults using the DNeasy Blood 121

and Tissue Kit (Qiagen) following the manufacturer's protocol.

# **BAC sequencing and assembly**

## 124 BAC sequencing

A BAC library was constructed by Rx Biosciences (Rockville, MD, USA) from whole genomic DNA 125 of three *M. edulis* individuals; and using pCC1BAC BamHI as vectors. BAC sequencing was per-126 formed in three different experiments. First, 32 clones were prepared for Roche 454 pyrosequencing. 127 Indexed libraries were pooled in equimolar proportions and sequenced on a Roche (Branford, CO, 128 USA) GS FLX instrument that generated single reads of 600 bp on average. The sequencing of eight 129 libraries failed, and they were sequenced again on an Illumina MiSeq platform (San Diego, CA, USA) 130 that produced paired-end reads of 250 bp. Because the BAC inserts proved to be shorter (25Kb on 131 average) than expected (100 Kb) and because of assemby difficulties with the AT-rich mussel genome, 132 we conducted a third sequencing experiment, in which a single pool of 192 anonymous clones was 133 sequenced on a single lane of an Illumina HiSeq2000 instrument that generated paired-end reads of 134 101 bp. Reads were trimmed for index sequences and low-quality terminal bases. They were deposed 135 in the NCBI Short Read Archive [[XXX]]. 136

#### 137 *de novo* BAC assembly

In the absence of a sequenced *Mytilus* genome at the time of the experiments, we conducted *de novo* 138 assemblies of BAC sequences in two steps for each sequencing run. We first assembled reads into 139 contigs with different programs depending on the sequencing technology. 454 reads were assembled 140 with Newbler v1.0.1 (Margulies et al. 2005), a de novo DNA assembler designed for pyrosequencing, 141 with the following parameters: seed length of 6 bp, minimal read length of 40 bp, minimal overlap 142 length of 40 bp, minimal overlap identity of 90%. We also used the program MIRA v3.4 (Chevreux 143 2005) with 454 and accurate settings. MiSeq and HiSeq reads were assembled using ABySS v1.2.1 144 (Simpson et al. 2009) and SOAP denovo v2.0.4 (Luo et al. 2012), two de novo short-reads assemblers, 145 using a k-mer size of 61 bp and an insert size of 500 bp for the MiSeq paired-reads and of 200 146 bp for the HiSeq paired-reads. Secondly, contigs generated in the first step were assembled into 147 longer fragments with the program CAP3 (Huang 1999). The resulting genomic assemblies for each 148

sequencing run were compared with the *DC-MEGABLAST* algorithm. We retained the longest contig
of each assembly as well as uniquely assembled contigs. DNA contaminants (BAC vectors, bacteria,
etc.) were removed with *DC-MEGABLAST*. Contigs of length < 1 Kb were discarded, except for the</li>
HiSeq assembly in which we used a 5 Kb threshold. Finally, a comparison of the 454, MiSeq and
HiSeq genomic assemblies was performed with *DC-MEGABLAST* to produce a final contig set of 378
sequences (average length 8.5 Kb, maximal length 26.4 Kb) without duplicates.

# **155** Target enrichment sequencing

#### 156 Targets

We enriched genomic DNA for 3 Mb of target regions using a SureSelectXT Custom system (Ag-157 ilent Technologies, Santa Clara, CA) comprising  $\sim 55,000$  RNA probes of 120 bp (Mamanova et 158 al. 2010). First, probes were designed from our M. edulis BAC assembly totaling 2 Mb of filtered 159 sequences. Second, we designed probes from cDNA contigs. We used a random panel of 338 M. 160 galloprovincialis cDNA contigs of  $\sim 1.6$  Kb on average (0.5 Mb in total) previously generated for a 161 transcriptome sequencing project (Romiguier et al. 2014). BAC and cDNA reference sequences were 162 annotated with the genome annotation program MAKER v2.31 (Cantarel et al. 2008) using ab-initio 163 gene predictions as well as local alignements onto Mytilus GenBank collections. In addition, we de-164 signed an additional set of probes from publicly released expressed sequence tags (ESTs) databases by 165 focusing on genes with functions of interest. We targeted 553 ESTs averaging  $\sim 1$  Kb in length (0.5 166 Mb in total): 262 immunity genes (from Mytibase, Venier et al. 2009 and another Mytilus repertoire 167 of immune genes, Philipp et al. 2012); 133 genes involved in cytonuclear interactions (identified from 168 mitodrome, D'Elia et al. 2006, and mitores, Catalano et al. 2006, in Drosophila); 30 reproduction-169 related genes, 20 habitat-related genes (because of the known association with wave action, proteins 170 from the foot, the byssus filament and adhesive plaques are suspected to be involved with habitat spe-171 cialization) and 6 nucleoporines (following Nolte et al. 2012) directly recovered from GenBank. To 172 these genes of known function, we added a panel of control genes composed of 102 genes known to 173 be single gene orthologous for phylogenomics analysis (OrthoDB in vertebrates (Kriventseva et al. 174 2008) and genes used in molluscan phylogenomics Kocot et al. 2011). We eliminated redundancy 175 between the two coding gene sets (from RNA-seq and EST libraries) by a local alignment analysis 176

(*DC-MEGABLAST*). To maximize capture of unique sequences, we identified and masked repetitive and low-complexity regions with *WindowMasker* (Morgulis *et al.* 2006) using *Mytilus* GenBank collections and our own BAC sequences as references. We designed 120 bp probes (2X tilling) covering the final masked genomic data with *OligoTiler* (*http* : //*tiling.gersteinlab.org*). Orphan probes were duplicated and low GC-content probes (< 10 *GC*%) were quadruplicated. The production of the probe library was performed by Agilent SureDesign services.

# **183** Capture and sequencing

Illumina paired-end sequencing libraries with insert sizes of 300 - 600 bp were prepared for each 184 individual. The standard Illumina TruSeq DNA Sample Preparation (http://support.illumina.com) 185 was followed, except that we used custom paired-end adaptors incorporating a unique 6-bp index as 186 well as Indexed Blocking Reagent to perform pre-capture multiplexing (Kenny et al. 2010). A total of 187 88 libraries were produced (additional samples not discussed here were also included). Individuals of 188 the same population were pooled in equimolar proportions prior to being subjected to TruSeq Custom 189 Enrichment (*http://support.illumina.com*). Several enrichment protocols were attempted in pilot 190 runs on Miseq and GA2X instruments in order to increase capture specificity (TableS3). We sequenced 191 all libraries on a full flow-cell of HiSeq2000 producing 101 bp paired-end reads. Reads were trimmed 192 for index sequences and low-quality terminal bases; low-quality reads were discarded. They were 193 deposed in the NCBI Short Read Archive [[XXX]]. The capture step and pilot runs were subcontracted 194 to the Plateforme Génome Transcriptome (CGFB, Université Bordeaux Segalen, France) and HiSeq 195 sequencing was performed at the sequencing plateforme of the "Génomique et maladies métaboliques" 196 laboratory (CNRS UMR 8199, Lille, France). 197

# <sup>198</sup> Genomic enrichment mapping to reference contigs and SNP calling

## 199 Read mapping

We conducted mapping with *bwa-mem v0.7.5a* (Li & Durbin 2009), a fast aligner that use the Burrows-Wheeler algorithm. Two challenges had to be overcome. First, because our reference genomic assembly consisted of sequences from the three closely-related *Mytilus* species, a trade-off between stringency and specificity had to be found. Second, because masking was necessarily incomplete given the

relatively little genomic information available in mussels, sequencing depth was highly variable along 204 BAC sequences. Typically, depth in repeated regions was above 1000X whereas regions of interest 205 were often < 50X. As a consequence, we evaluated the mapping performance of different settings in 206 one individual of each species based on the proportion of mapped reads and sliding window analysis 207 of depth on BAC sequences. Ultimately, bwa-mem was run with the following non-default parame-208 ters: a clipping penality of 3, a mismatch penality of 2, a gap open penality of 3, and a minimum seed 209 length of 10. Mapping was independently performed for each sequencing run. We then generated final 210 mapping for each individual by merging the sorted alignements (either GA2X or MiSeq, with HiSeq) 211 using SAMtools v0.1.19 (Li et al. 2009). 212

# 213 SNP and genotype calling

SNP calling was performed following successive steps to obtain a dataset of high-quality SNPs across 214 the three species. We used the *mpileup* tool of *SAMtools* to pileup the merged alignments of each 215 individual. We then called variant candidates from all individuals with bcftools (Li 2011) to produce 216 an initial database of SNPs. This database was subsequently used for multisample variant calling, 217 performed separately in each population and applying various quality filters with VCF tools v0.1.12a 218 (Danecek et al. 2011). To reduce as much as possible bias in allele frequencies, we excluded calls from 219 positions with an averaged depth across individuals below 10 reads. We also applied an upper depth 220 limit to exclude unmasked repeated regions. As populations varied in terms of sequencing depth, 221 we used a maximal value set at the 98.5th percentile of the depth distribution across all positions 222 (see TableS3). To filter out paralogous regions, we excluded sites deviating from Hardy-Weinberg 223 equilibrium (*p-value* < 0.05) using an exact test implemented in VCFtools. 224

Only variants that passed filters across all populations were retained for subsequent analysis. Ac-225 cording to the algorithm implemented in *bcftools*, a variant was called when the posterior probability 226 of non-reference allele counts was above 50%, assuming a standard allele frequency spectrum. For 227 each variant detected, the maximum *a posteriori* genotype was assigned to each individual assuming 228 Hardy-Weinberg proportions in genotype prior probabilities. Only genotype calls with a quality score 229 above 10 were retained, otherwise missing data was applied. For genealogical analyses, we used a 230 higher genotype quality threshold set at 30. Any position with more than 20% of missing data were 231 discarded. For outlier gene analysis, contigs with fewer than 20 SNPs were discarded. 232

# 233 Data analysis

The final data set consisting of 1269 reference sequences and polymorphism data is available on http: //www.scbi.uma.es/mytilus/index.php. Population genetic analyses were performed using custom scripts in R (R Core Team 2012).

#### 237 Patterns of genetic differentiation

We initially explored the data using a Principal Component Analysis implemented in the R package 238 ade4 (Dray & Dufour 2007) based on genotypes and excluding any position with missing data and sin-239 gletons. FST values (Weir & Cockerham 1984) were calculated using the R package hierfstat (Goudet 240 2005) for each SNP between pairs of populations. Data were analysed at the level of the contig to limit 241 pseudoreplication of closely-linked SNPs. Because we were interested in genomic regions that were 242 highly differentiated between populations, we calculated the 90th percentile of the FST distribution 243 of each contig  $(FST_{90})$ , as well as its maximal FST value  $(FST_{max})$ . Joint distributions of interspecific 244 FST<sub>90</sub> values were analysed by Standardised Major Axis regressions, slope and elevation were esti-245 mated and tested using the R package smatr (Warton et al. 2012). For outlier identification, we fitted 246 a null empirical  $FST_{90}$  and  $FST_{max}$  distributions across all contigs. Contigs in the upper 2.5% of the 247 empirical  $FST_{90}$  and/or  $FST_{max}$  distributions were categorized as outliers. They were further analysed 248 by estimating the allele frequency variation along the sequence. Open reading frames were predicted 249 with ORF finder (http://www.ncbi.nlm.nih.gov/projects/gorf) and non-synonymous changes were 250 identified using *BioEdit* v7.2.5 (Hall 1999). We also evaluated the proportion of exclusively shared 25 SNPs in a given set of populations, after removing singletons. 252

#### 253 Phylogenomic analysis

Genotype data were phased with *BEAGLE v3.3.2* (Browning & Browning 2007) using genotype likelihoods provided by *bcftools*. All individuals were included in the analysis to maximize linkage disequilibrium and 20 haplotype pairs were sampled for each individual during each iteration of the phasing algorithm to increase accuracy. FASTA haplotype sequences were then generated using a custom perl script. A phylogenetic network analysis was conducted with *SplitsTree4 v4.12.6* (Huson & Bryant 2006) to get insight into the population relationships across the three hybridizing species. All hap-

lotype sequences were compiled to create an artificial chromosome of 51,878 variable positions and analysed using the neighbour-net method. Haplotype sequences of each candidate locus were also individually analysed with the neighbour-net method. Finally, we quantified the degree of exclusive ancestry between populations by computing a Genealogical Sorting Index (Cummings *et al.* 2008) for each locus based on allelic genealogies inferred with the R package *ape* (Paradis *et al.* 2004) using a neighbour-joining algorithm with F84 distances (Felsenstein & Churchill 1996).

# 266 **Results**

# <sup>267</sup> BAC assembly, target enrichment performance and variant calling

Screening and assembly results for each BAC clone are summarized in TableS2. The final BAC assem-268 bly contained 378 contigs, with an average length of 8.5 Kb, including 5% above 18 Kb (maximum 269 26.4 Kb), and a total assembly length of 3.2 Mb. Performance of DNA target enrichment sequencing 270 and mapping are reported in TableS3. A total of 1269 contigs (378 BAC and 891 cDNA contigs) were 271 captured in 75 individuals spanning the three species (see Material and Methods for details). Reads 272 were aligned against a single reference made of sequences from the three species. On average, 55% 273 of the reads mapped to the reference and all three species had a similar proportion of reads aligned. 274 The performance of target enrichment mainly resulted from a combined effect of library quality and 275 capture protocol (the second alternative protocol led to higher capture specificity). 85% of targeted 276 sequences (1079) were successfully captured and assembled with mean read depth of 35X. After qual-277 ity filtering (see Material and Methods for details), variant calling produced a total of 122,144 SNPs 278 across all individuals. Among populations, the proportion of variant sites (13,827 on average) and 279 observed heterozygosity ( $h_o = 0.032$  on average) were roughly similar. However, M. trossulus popu-280 lations tended to be more polymorphic ( $h_{o\_trossulus} = 0.042$ ) than populations from the other species 281  $(h_{o-galloprovincialis} = 0.030 \text{ and } h_{o-edulis} = 0.029)$ , mainly due to private SNPs (Figure 2). 282

# 283 Genome-wide species relationships

A genome-wide network of genetic relationships (Figure 2A) was built from a subset of 51,878 highquality SNPs genotyped in 72 individuals (3 individuals were discarded due to misidentification or sequencing failure; Table S3). We observed that each species formed a distinct cluster suggesting that

on average a high proportion of SNPs supports the "species tree" topology. Given their more recent 287 divergence, M. edulis and M. galloprovincialis were less differentiated from each other than from M. 288 trossulus. This is also apparent in the multivariate analysis of genotypes (Figure 2B) in which M. 289 trossulus individuals were clearly separated from individuals of the other species in the first axis. The 290 second axis differentiated *M. edulis* from *M. galloprovincialis* individuals. In the following axes (not 29: shown), the American and European populations of *M. edulis* were separated as well as the Atlantic 292 and Mediterranean populations in M. galloprovincialis; then the east and west Mediterranean Sea and 293 finally the enclosed patch from the peripheral parental population in *M. galloprovincialis* (Brittany vs 294 Iberian Coast) and M. edulis (Bay of Biscay vs North Sea). A last line of evidence comes from Figure 295 2C showing that the majority of exclusively shared SNPs stands between populations of the same 296 species ( $S_{trossulus} = 8245$ ,  $S_{galloprovincialis} = 1673$  and  $S_{edulis} = 417$ ) or between genetic clusters within 297 species ( $S_{edulisEU} = 555$  and  $S_{edulisAM} = 358$ ;  $S_{galloprovincialisATL} = 387$  and  $S_{galloprovincialisMED} = 1011$ ). 298 This is even clearer when considering the subset of shared and fixed SNPs (Table1) which were nearly 299 all species-specific (i.e. polymorphic or fixed in one species but absent in others). We noted in 300 Figure 2A that a *M. trossulus* individual from America ranked at the bottom of its cluster suggesting 301 significant levels of recent introgression, which is consistent with the close geographical proximity of 302 *M. edulis* populations there. This is further supported by Figure 2B in which this individual appears 303 to group outside of its population, shifted toward the *M. edulis* clusters. Together with the *M. edulis* 304 individual found in this population (excluded from analysis), this reveals that admixed groups coexist 305 in the Saint-Lawrence river. 306

# 307 Genetic differentiation: geography

Alongside the "species tree" topologies, discrepant gene histories were also clearly identified. These 308 may be due to shared ancestral polymorphism, or gene flow experienced by populations during their 309 history. We characterized patterns of differentiation between populations at different levels of di-310 vergence ("intraspecific" vs "interspecific"), and geographical isolation ("parapatric" vs "allopatric") 311 using the upper decile of the FST distribution for each contig ( $FST_{90}$ , Figure 3). The intraspecific 312  $FST_{90}$  distribution was L-shaped, with most loci undifferentiated between populations and a few loci 313 highly differentiated. With increasing geographical isolation, the genome-wide average FST<sub>90</sub> value 314 increased from 0.076 to 0.131 in *M. edulis* and from 0.086 to 0.107 in *M. galloprovincialis* (Table 1). 315

The two *M. trossulus* populations were by far the most differentiated, with an average  $FST_{90}$  value of 0.362 (also visible in Figure 2). The genome-wide variance in  $FST_{90}$  increases with increasing levels of divergence and tends to become bimodal in interspecific comparisons, including both a higher proportion of highly and lowly differentiated loci.

M. edulis and M. galloprovincialis populations range along a gradient of geographical distances 320 from pure allopatry with American M. edulis populations, to different degrees of parapatry in Eu-321 rope (Figure 1). As expected based on geography, average FST<sub>90</sub> values between M. edulis and M. 322 galloprovincialis were higher in allopatric populations ( $FST_{90 galloprovincialis-edulisAM} = 0.41$ ) than in 323 parapatric populations ( $FST_{90}$  <sub>galloprovincialisMED-edulisEU</sub> = 0.39 and  $FST_{90}$  <sub>galloprovincialisATL-edulisEU</sub> = 324 0.31). The joint distribution of interspecific  $FST_{90}$  in allopatric and parapatric populations further 325 confirmed the effects of gene flow in reducing differentiation (Standardised Major Axis regression, 326 *elevation* = 0.05; p < 0.0001, Figure 4A). Similarly, the joint distribution of interspecific  $FST_{90}$  in 327 the Atlantic (closer to M. edulis) and Mediterranean (further from M. edulis) populations confirmed the 328 genome-wide difference in introgression rates (SMA regression: *elevation* = 0.03; p < 0.0001, Fig-329 ure 4B). However this also revealed that some outlier loci showed the opposite pattern of introgression 330 (i.e. highly differentiated in the Atlantic but lowly differentiated in the Mediterranean Sea), suggesting 331 a more complex history between M. edulis and M. galloprovincialis. Interspecific  $FST_{90}$  correlations 332 between Mediterranean populations showed no differences in introgression rates between the Eastern 333 and Western basins (non-significant elevation, Figure 4C); whereas the internal Atlantic population 334 was significantly more introgressed than its external counterpart (*elevation* = 0.02; p = 0.001, Figure 335 4D). Overall, this is consistent with the genome-wide gradient of increasing *M. edulis* introgression 336 observed in Figure 2, from the enclosed patch in Brittany to the Iberian Coasts, the West and the East 337 of the Mediterranean Sea. With regard to M. edulis in Europe, the joint distribution of interspecific 338 FST<sub>90</sub> showed that the enclosed patch in the Bay of Biscay was not more introgressed than its external 339 reference (non-significant *elevation*, Figure 4E). 340

# 341 Genetic differentiation: replicated contacts

M. *edulis* and *M. trossulus* are currently in contact in Europe and in America (Figure 1). To evaluate the degree of genetic parallelism between contacts, we assessed the level of interspecific  $FST_{90}$  correlation between the European and American replicates (Figure 5A). Despite being lower than the correlation

in *M. edulis* and *M. galloprovincialis* comparisons (Figure 4), it was still significant ( $r_{pearson} = 0.57$ ) 345 showing that the outcomes were to some extent similar in the two contacts (highly/lowly differentiated 346 regions partially overlap). More importantly, Figure 5A showed that genome-wide differentiation be-347 tween European species was lower than their American counterparts (*elevation* = 0.17; p < 0.0001) 348 reflecting genome-wide asymmetry in *M. edulis* introgression rates between the two species; the Eu-349 ropean M. trossulus population being more introgressed by M. edulis alleles at the genomic level. 350 Outlier loci were also asymmetric with a deficit of loci both highly differentiated in Europe and lowly 351 differentiated in America. Figure 5B further confirmed the asymmetry of introgression. It shows 352 that loci characterized by a strong phylogenetic separation between American M. trossulus and other 353 species (high GSI values) were less exclusive to *M. trossulus* in Europe (*elevation* = 0.03; p < 0.01). 354 Also, this is clear across Figure 2 in which the European *M. trossulus* individuals were closer to *M.* 355 edulis individuals; and from Figure 3 in which lower genome-wide differentiation was observed in 356 comparisons involving the European *M. trossulus* population. 357

# 358 Within-species outliers: examples of the different categories

Making clear the evolutionary relationships between species of the *Mytilus* complex allowed us to 359 identify the causes of outlying levels of differentiation within species and to reconstruct the evolution-360 ary histories of outlier loci. Based on analyses of variation in differentiation at small chromosomal 361 scales, as well as allele frequencies and allele genealogies (Figure 6), candidate outliers were placed 362 in different categories depending on whether they were most plausibly due to differentiation of in-363 traspecific alleles, or introgression of heterospecifc alleles. Figure 6A and 6B illustrate representative 364 cases of candidate loci for local introgression (see Figure S1 for additional examples). Figure 6A 365 represents a complete introgression of *M. edulis* haplotypes of several Kb in length, into the European 366 *M. trossulus* population; while American *M. trossulus* remained unaffected by *M. edulis* introgression. 367 In most cases, a major part of the heterospecific genetic diversity has introgressed into the recipient 368 species. In a few cases however (3/83), a small proportion of the heterospecifc diversity has swept 369 and introgressed between species. A clear example is a female-specific transcript for which a single 370 edulis haplotype has introgressed into Atlantic populations of *M. galloprovincialis* (Figure 6B). The 371 chromosomal footprint shows a peak of FST at the 3' side of the sequence; and several radical amino 372 acid changes have occured in the protein. In addition, a distinct group of haplotypes are found in the 373

Mediterranean Sea and in the American *M. trossulus* populations. Together, these patterns are consistent with the hypothesis of adaptive introgression across the mosaic hybrid zone in Europe. Figure 6C illustrates an example of an outlying level of differentiation that did not involve introgression from another species. It is a classical sweep in the external Mediterranean population of *M. galloprovincialis*, characterized by a star-shape clade and a high frequency of a new variant at the maximal FST value (see Figure S1 for additional examples).

# 380 Within-species outliers: major causes

Table1 reports the total number of intraspecific outliers from within-species comparisons. The number of such outliers increased steadily with the average level of genomic differentiation, from  $\sim 15$ in parapatric comparisons, to 40 between the allopatric *M. trossulus* populations. Outlier tests were not possible between species because the distribution of *FST* values was overdispersed. At the contig level, some loci were repeatedly involved in outlying levels of differentiation ("shared outliers"; TableS4), but their number was not significantly different than would be expected by chance.

Overall, outlier analysis revealed that local introgression of heterospecific variants contributes to 387 a significant part of the within-species differentiation (Table1 and TableS5 for details). In M. gal-388 loprovincialis, introgression of M. edulis variants explains between 55% and 80% of the outliers. 389 Both basins of the Mediterranean Sea (East and West) presented highly introgressed variants; whereas 390 outlying levels of introgression were much more asymmetric between the two Atlantic populations, 391 corroborating the genome-wide trend. Between European and American M. edulis populations, more 392 than half of the outliers were due to local introgression. While heterospecific variants came from M. 393 trossulus in America, the European populations were introgressed by M. galloprovincialis variants. 394 Notably, outlying levels of differentiation between the European *M. edulis* populations were never 395 due to introgression, in agreement with the genome-wide picture. Moreover the number of outliers 396 was small, their level of differentiation low and restricted to a small region of the contigs. In M. 397 trossulus, local introgression from *M. edulis* explains more than 90% of the outliers. As is the case 398 genome-wide, outlying levels of introgression were highly asymmetric: European populations being 399 more permeable to *M. edulis* introgression than were their American counterparts. 400

Remarkably, the different sources of within-species differentiation often led to subtle chromosomal footprints, which would be difficult to identify in a large-scale survey of genome-wide differen-

tiation. For example, introgression candidates generally involved several haplotypes from the donor 403 species leading to a "soft sweep" patterns (Figure 6 and Figure S1). Similarly, variants that swept to 404 high frequency in a single population often segregated as ancestral polymorphism in other populations, 405 instead of arising as new mutations (TableS5). The chromosomal scale of the footprint sometimes ap-406 peared to extend beyond the scale investigated here (a few Kb). This was particularly the case for local 407 introgression candidates, because heterospecific haplotypes are initially introgressed from one species 408 to another. Nevertheless, in some cases, a single peak of a few hundred base pairs was detected (Fig-409 ure 6 and Figure S1) allowing the identification of candidate causal variants (TableS5). When direct 410 selection was suspected (i.e. when non-synonymous variants were identified in open reading frames), 411 amino acid changes were rarely ever differentially fixed between populations, suggesting a multigenic 412 architecture for the traits implicated in differentiation. We also investigated whether any functional 413 category was overrepresented among the list of EST outliers (TableS6). Across comparisons, there 414 was no category represented more often that would be expected by chance. However, when consider-415 ing EST outliers differentially fixed between species ("diagnostic loci"), the most likely candidates to 416 represent direct barrier to gene flow, we noted a slight excess of immune genes together with a slight 417 deficit of single copy orthologous genes. 418

# 419 **Discussion**

Early empirical studies of the build up of reproductive isolation, emphasized the role of gene flow 420 in shaping genome-wide patterns of differentiation. But depite this, the evolutionary importance of 421 gene flow for adaptation and speciation remains unclear (Hedrick 2013; Seehausen et al. 2014). This 422 is partly because variation in levels of genetic differentiation are also influenced by variation in lo-423 cal genomic parameters (such as mutation and recombination rates), which determine the speed of 424 divergence tuned by selective interference (including background selection, hitchhiking, and Hill-425 Robertson effects), which implies that candidate regions may not represent interspecific barriers to 426 gene flow (Cruickshank & Hahn 2014). Indeed, genome scans in regions of low recombination strug-427 gle to disentangle the relative contribution of reduced gene flow and selection at linked sites in pro-428 ducing high level of differentiation (Nachman & Payseur 2012). Moreover, even if gene flow and 429 hybridization occur during speciation, they are not expected to greatly impede divergence (Endler 430

<sup>431</sup> 1977). In contrast, the acquisition of beneficial variants by migration from another population or
<sup>432</sup> species (i.e. adaptive introgression) may be crucial in adaptation (Arnold 2004; Hedrick 2013). This
<sup>433</sup> shows the importance of considering the history of related lineages when asking how organisms adapt
<sup>434</sup> to their environment.

Here, we used targeted enrichment sequencing to perform fine-scale empirical genome scans of 435 differentiation between closely related species of mussels in different geographical contexts. Given 436 that detection of loci influenced by divergent selection using genome-scans carries a significant risk of 437 false positives (extreme values of FST might not necessarily be due to selection) and that multifarious 438 processes can lead to FST outliers, we further investigated candidate outliers by analyzing allele fre-439 quencies and gene genealogies along contigs. A phylogenetic network (Figure 2) suggested discordant 440 genealogical histories across the genome, i.e. genome regions that do not follow the "species-tree" 441 topology. Incomplete lineage sorting, in addition to introgression, can lead to interspecies shared poly-442 morphism. This is especially the case for recently-diverged lineages with large effective population 443 sizes, such as *Mytilus* mussels. We were able to control for this, by using an allopatric population 444 (e.g. *M. edulis* in America) as a reference for the level of shared ancestral polymorphism, and thereby 445 showed that interspecific differentiation between parapatric populations was lower than between their 446 allopatric counterparts (Figure 3 and Figure 4). This result matches other studies highlighting the role 447 of gene flow in eroding genetic differentiation in neutral regions (e.g. in Heliconius, Martin et al. 2013 448 and in *Helianthus*, Renaut et al. 2013). If selection against hybrids/immigrants is sufficiently strong 449 compared to the migration rate, genomic regions involved in the interspecific genetic barrier are not 450 expected to introgress between species. As such, outlier loci, highly differentiated in both allopatric 451 and parapatric comparisons, are good candidates for reproductive isolation genes (Figure 4). 452

In Europe, a genome-wide gradient of introgression, in agreement with current range distributions, 453 was observed in *M. galloprovincialis* populations from the Mediterranean Sea to the internal Atlantic 454 population (Figure 2). The high level of introgression detected in the internal Atlantic population 455 (Figure 4) is well explained by its enclosed position within the mosaic hybrid zone, and concords with 456 a previous study using few markers (Bierne et al. 2003). While this pattern is consistent with on-457 going neutral introgression between geographically proximate populations, it may also be explained 458 by migration of heterospecific variant favoured in the internal M. galloprovincialis population due to 459 similarities in selective pressures. Such candidates should be weakly differentiated from M. edulis 460

in the internal population, but highly differentiated elsewhere (Figure 4). Another line of evidence 461 supporting cases of non-neutral introgression is that some of the outliers of differentiation within M. 462 galloprovincialis retained the signature of ancient introgression into the Mediterranean Sea (Figure 463 4 and Table1). In addition to candidates previously reported (Gosset & Bierne 2012; Fraïsse et al. 464 2014b), many new candidates for adaptive introgression have been found in this study (TableS5). 465 However, demonstrating that introgression was initially selectively-driven is tricky, because the pat-466 tern left by adaptation from recurrent migration may be quite different from the classical hitchhiking 467 signature. Specifically, if adaptive introgression involves several haplotypes of the same beneficial 468 mutation, most of the heterospecific diversity is expected to introgress into the foreign genetic back-469 ground leading to a "soft sweep" pattern (Pennings & Hermisson 2006). Such a pattern was often 470 found among within-species outliers when local introgression was observed (Figure S1). This em-471 phasizes a neglected hypothesis in which genetic hitchhiking is not involved: a local change in the 472 permeability of a barrier to gene flow (Gagnaire et al. 2013; Fraïsse et al. 2014b), so that in some 473 genomic regions, some populations are more permeable to introgression than others. 474

In contrast to other comparisons, the two European M. edulis populations shared similar level of 475 differentiation with *M. galloprovincialis* even though the internal population is enclosed within the 476 mosaic hybrid zone (Figure 2 and Figure 4). This suggests that the genome of *M. edulis* resists cur-477 rent introgression from M. galloprovincialis and that the two European populations share a similar 478 genetic architecture with respect to interspecifc gene flow. Accordingly, all cases of outlying differen-479 tiation involved homospecific processes (Table1). For example, a candidate "global sweep", identified 480 previously (Faure et al. 2008) was characterized by fixed haplotypes in the external population and 481 non-fixed haplotypes in the internal population (Figure S1-B4). In fact two domes of differentiation 482 are expected on either side of the benefical mutation because recombination breaks down its associ-483 ation with the haplotype initially entering the internal population (Bierne 2010). Unfortunately, we 484 failed to assemble the targeted BAC to the region directly under selection and so we cannot confirm 485 this pattern. Other outliers showed a classical signature of a selective sweep, with a sharp single peak 486 of FST between the European and American *M. edulis* (Figure S1-B2). Even when non-synonymous 487 changes were identified - suggesting direct selection, variants were generally not differentially fixed 488 between them (TableS5). This may reveal pervasive polygenic adaptation which implies slight allele 489 frequency differences at a large number of loci underlying the selected traits (Pritchard & Di Rienzo 490

<sup>491</sup> 2010). Also swept variants often segregated as low-frequency polymorphisms in other populations
<sup>492</sup> (TableS5), suggesting that selection from standing genetic variation may be widespread, although it
<sup>493</sup> leaves a less pronouced signature at linked sites than does a sweep from a single new mutation (e.g.
<sup>494</sup> in mice, Domingues *et al.* 2012; in sticklebacks, Hohenlohe *et al.* 2010; Roesti *et al.* 2014; and in
<sup>495</sup> mussels, Gosset *et al.* 2014).

While most of the genome is resistant to current introgression in European M. edulis, identification 496 of within-species outliers between American and European *M. edulis* populations showed a putative 497 case of adaptive introgression in Europe. A gene expressed solely in females (H5, Anantharaman 498 & Craft 2012) has swept through M. edulis and M. galloprovincialis across the hybrid zone, while 499 being structured between American M. edulis and Mediterranean M. galloprovincialis (Figure 6B). 500 In addition, the different haplotypes carry several non-synonymous changes suggesting that the gene 501 is rapidly evolving (TableS5). Similar patterns have been documented in *Mytilus* mussels in a gene 502 expressed in male gametes (M7 lysin) which dissolves the egg vitelline coat (Riginos et al. 2006; 503 Springer & Crespi 2007). This surprising observation may suggest that a locus *a priori* implicated in 504 reproductive isolation may have swamped the species barrier between specific genetic backgrounds 505 while remaining impermeable between others. Generally speaking, recent reviews of the molecular 506 basis of species formation have not shown any specific biochemical pathway preferentially involved in 507 incipient reproductive isolation (Presgraves 2010). Instead the main feature of isolating genes is their 508 rapid evolution, often attributable to evolutionary conflicts (cyto-nuclear interactions, host-pathogen 509 interactions or sexual conflicts). For example, Burton & Barreto (2012) emphasize the role of the 510 mitochondrial genome in producing incompatible interactions with nuclear genes. Immune genes are 511 also known to evolve faster than others (e.g. in humans, Fumagalli et al. 2011 and in flies, Obbard 512 et al. 2009) and so they may be implicated in species barriers. However, our functional comparison 513 of diagnostic ESTs between species did not reveal any obvious enrichment for these genes, although 514 immune genes were slightly more common than expected by chance (TableS6). 515

The European and American contacts between *M. edulis* and *M. trossulus* (Figure 1) allowed us to investigate the extent to which differentiation is parallel in replicated contacts. Interestingly, we found incomplete parallelism in the genome-wide patterns of differentiation (Figure 5). The independent history of divergence between the two pairs is likely to be responsible for this pattern. Specifically, the genetic differences accumulated in allopatry may either scatter, or couple with each other when the

divergent lineages come into contact in locations stabilized by a natural barrier to dispersal (Barton & 521 Hewitt 1985) or environmental boundaries (Bierne et al. 2011). This leads to a genomically localized 522 breakdown or strengthening of the barrier to gene flow (Barton & de Cara 2009). Therefore, indepen-523 dent outcomes of the coupling process following secondary contact between the same two lineages (M. 524 edulis and M. trossulus) is expected to produce partially different genetic architectures of reproduc-525 tive isolation. Supporting this hypothesis, changes in the permeability of the species barrier have also 526 been reported in replicated pairs of whitefish species (Gagnaire et al. 2013). However, we observed 527 in mussels a striking pattern of genome-wide asymmetry in introgression level between the European 528 and American contacts. More importantly, outlying levels of differentiation were also highly asym-529 metric: almost all cases of local introgression from M. edulis were found in the European M. trossulus 530 lineage (Table1). Our study consolidates previous similar findings with nuclear DNA markers (Borsa 531 et al. 1999) and mitochondrial genomes, which have been completely replaced in the Baltic Sea 532 (Rawson & Hilbish 1998; Quesada et al. 1999). This asymmetry is not an expected outcome of the 533 coupling process, because independent incompatibilities can couple in either direction (Barton & de 534 Cara 2009) leading to overall symmetric reproductive isolation. Riginos & Cunningham (2005) and 535 Väinölä & Strelkov (2011) have proposed that the secondary contact in Europe may be older than that 536 in America: so that an extended period of gene flow could have led to the erosion of isolating barriers 537 still functioning in America. To assess whether or not this hypothesis is valid will require a rigorous 538 inference of the demographic history of the two lineages. Asymmetrical introgression may also sim-539 ply reflect differences in effective population size or in the environmental landscape in Europe, which 540 influence the direction of gene flow and the level of genetic drift. It worth emphasising that if the 541 *M. trossulus* ancestry of Baltic mussels is still clearly evident, a more permeable barrier would have 542 easily resulted in a more general swamping of the Baltic genome by the North Sea genome. In this 543 case, the reconstructed history of populations from the neutral fraction of the genome could easily 544 result in the mistaken inference that differentiation was recent (Bierne et al. 2013). This could have 545 happened in other species exhibiting a genetic break at the entrance of the Baltic Sea (Bierne et al. 546 2011) and explained the astonishing number of outliers sometimes found between Baltic and North 547 Sea populations (Lamichhaney et al. 2012). 548

In conclusion, this study has revealed the complexity of the relationships among populations of *Mytilus* species which are both influenced by current and past demogeography. Based on qualita-

tive analysis of genetic polymorphism, we highlighted the importance of introgression as a neglected
 source of adaptive variation to be considered in genome scans.

# **553** Acknowledgments

Analyses largely benefited from the Montpellier Bioinformatics Biodiversity computing cluster plat-554 form. It is a pleasure to thank Rémy Dernart for his support on the computing platform; Christophe 555 Hubert and Véronique Dhennin for their technical help in producing data; Baptiste Faure and Pierre-556 Alexandre Gagnaire for sampling. We also thank Nicolas Galtier, Pierre-Alexandre Gagnaire, Patrice 557 David and François Bonhomme for discussions. Molecular data were produced through the ISEM 558 platform Génomique marine at the Station Méditerranenne de l'Environnement Littoral (OSU OREME 559 (Observatoire de Recherche Méditerranéen de l'Environnement)) and the platform Génomique Envi-560 ronnementale of the LabEx CeMEB (Laboratoire d'Excellence Centre Méditerranéen de l'Environnement 561 et de la Biodiversité). This work was funded by the Agence Nationale de la Recherche (HYSEA 562 project, ANR-12-BSV7-0011) and the project Aquagenet (SUDOE, INTERREG IV B). This is arti-563 Jntpe cle 2015-XXX of Institut des Sciences de l'Evolution de Montpellier. 564

21

# 565 **References**

566 567	• Anantharaman S, Craft JA. (2012) Annual variation in the levels of transcripts of sex-specific genes in the mantle of the common mussel, <i>Mytilus edulis</i> . <i>PLoS One</i> , <b>7</b> .
568 569	• Arnold ML. (2004) Transfer and origin of adaptations through natural hybridization: were Anderson and Stebbins right? <i>Plant Cell</i> , <b>16</b> , 562–570.
570 571	• Barton N, Bengtsson BO. (1986) The barrier to genetic exchange between hybridising populations. <i>Heredity</i> , <b>57</b> , 357–376.
572 573	• Barton NH, de Cara MAR. (2009) The evolution of strong reproductive isolation. <i>Evolution</i> , <b>63</b> , 1171–1190.
574	• Barton NH, Hewitt GM. (1985) Analysis of hybrid zones. Annu. Rev. Ecol. Syst., 16, 113–148.
575 576	• Bierne N. (2010) The distinctive footprints of local hitchhiking in a varied environment and global hitchhiking in a subdivided population. <i>Evolution</i> , <b>64</b> , 3254–3272.
577 578 579	• Bierne N, Borsa P, Daguin C, Jollivet D, Viard F, Bonhomme F, David P. (2003) Introgression patterns in the mosaic hybrid zone between <i>Mytilus edulis</i> and <i>M. galloprovincialis. Mol. Ecol.</i> , <b>12</b> , 447–461.
580 581	• Bierne N, Gagnaire PA, David P. (2013) The geography of introgression in a patchy environment and the thorn in the side of ecological speciation. <i>Curr. Zool.</i> , <b>59</b> , 72–86.
582 583	• Bierne N, Welch J, Loire E, Bonhomme F, David P. (2011) The coupling hypothesis: why genome scans may fail to map local adaptation genes. <i>Mol. Ecol.</i> , <b>20</b> , 2044–2072.
584 585 586	• Borsa P, Daguin C, Ramos Caetano S, Bonhomme F. (1999) Nuclear-DNA evidence that north- eastern atlantic Mytilus trossulus mussels carry M. edulis genes. J. Mollus. Stud., 65, 504– 507.
587 588 589	• Browning SR, Browning BL. (2007) Rapid and accurate haplotype phasing and missing-data inference for whole-genome association studies by use of localized haplotype clustering. <i>Am. J. Hum. Genet.</i> , <b>81</b> , 1084–1097.
590 591	• Burton RS, Barreto FS. (2012) A disproportionate role for mtDNA in Dobzhansky-Muller in- compatibilities? <i>Mol. Ecol.</i> , <b>21</b> , 4942–4957.
592 593 594	• Cantarel BL, Korf I, Robb SM, Parra G, Ross E, Moore B <i>et al.</i> (2008) MAKER: An easy-to-use annotation pipeline designed for emerging model organism genomes. <i>Genome Res.</i> , <b>18</b> , 188–196.
595 596 597	• Catalano D, Licciulli F, Turi A, Grillo G, Saccone C, D'Elia D. (2006) MitoRes: a resource of nuclear-encoded mitochondrial genes and their products in Metazoa. <i>BMC Bioinformatics</i> , <b>7</b> , 36.
598 599 600	• 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. <i>Genet. Res.</i> , <b>70</b> , 155–174.

601 602	• Chevreux B. (2005) MIRA: An Automated Genome and EST Assembler. Duisburg: The Ruprecht-Karls-University.
603 604	• Cruickshank TE, Hahn MW. (2014) Reanalysis suggests that genomic islands of speciation are due to reduced diversity, not reduced gene flow. <i>Mol. Ecol.</i> , <b>23</b> , 3133–3157.
605 606	• Cummings MP, Neel MC, Shaw KL. (2008) A genealogical approach to quantifying lineage divergence. <i>Evolution</i> , <b>62</b> , 2411–2422.
607 608	• Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA <i>et al.</i> (2011) The variant call format and VCFtools. <i>Bioinformatics</i> , <b>27</b> , 2156–2158.
609 610 611	• Domingues VS, Poh YP, Peterson BK, Pennings PS, Jensen JD, Hoekstra HE. (2012) Evidence of adaptation from ancestral variation in young populations of beach mice. <i>Evolution</i> , <b>66</b> , 3209–3223.
612 613	• Dray S., Dufour A. (2007) The ade4 package: implementing the duality diagram for ecologists. <i>J. stat. softw.</i> , <b>22</b> , 1–20.
614 615 616	• D'Elia D, Catalano D, Licciulli F, Turi A, Tripoli G, Porcelli D <i>et al.</i> (2006) The MitoDrome database annotates and compares the OXPHOS nuclear genes of <i>Drosophila melanogaster</i> , <i>Drosophila pseudoobscura</i> and <i>Anopheles gambiae</i> . <i>Mitochondrion</i> , <b>6</b> , 252–257.
617 618	• Ellegren H, Smeds L, Burri R, Olason PI, Backström N, Kawakami T <i>et al.</i> (2012) The genomic landscape of species divergence in <i>Ficedula</i> flycatchers. <i>Nature</i> , <b>491</b> , 756–760.
619 620	• Endler JA. (1977) Geographic Variation, Speciation and Clines. New Jersey: Princeton University Press.
621 622	• Faure MF, David P, Bonhomme F, Bierne N. (2008) Genetic hitchhiking in a subdivided population of Mytilus edulis. BMC Evol. Biol., <b>8</b> , 164.
623 624	• Feder JL, Gejji R, Yeaman S, Nosil P. (2012) Establishment of new mutations under divergence and genome hitchhiking. <i>Phil. Trans. R. Soc. B</i> , <b>367</b> , 461–474.
625 626	• Feder JL, Nosil P. (2010) The efficacy of divergence hitchhiking in generating genomic islands during ecological speciation. <i>Evolution</i> , <b>64</b> , 1729–1747.
627 628	• Felsenstein J, Churchill GA. (1996) A Hidden Markov Model approach to variation among sites in rate of evolution. <i>Mol. Biol. Evol.</i> , <b>13</b> , 93–104.
629 630	• Flaxman SM, Feder JL, Nosil P. (2013) Genetic hitchhiking and the dynamic buildup of genomic divergence during speciation with gene flow. <i>Evolution</i> , <b>67</b> , 2577–2591.
631 632	• Fraïsse C, Elderfield JD, Welch JJ. (2014a) The genetics of speciation: are complex incompatibilities easier to evolve? <i>J. Evol. Biol.</i> , <b>27</b> , 688–699.
633 634	• Fraïsse C, Roux C, Welch JJ, Bierne N. (2014b) gene flow in a mosaic hybrid zone: Is local introgression adaptive? <i>Genetics</i> , <b>197</b> , 939–951.
635 636 637	• Fumagalli M, Sironi M, Pozzoli U, Ferrer-Admettla A, Pattini L, Nielsen R. (2011) Signatures of environmental genetic adaptation pinpoint pathogens as the main selective pressure through human evolution. <i>PLoS Genet.</i> , <b>7</b> , e1002355.

638 639 640	• Gagnaire PA, Pavey SA, Normandeau E, Bernatchez L. (2013) The genetic architecture of reproductive isolation during speciation-with-gene flow in lake whitefish species pairs assessed by RAD sequencing. <i>Evolution</i> , <b>67</b> , 2483–2497.
641 642	• Gosset CC, Bierne N. (2012) Differential introgression from a sister species explains high FST outlier loci within a mussel species. <i>J. Evol. Biol.</i> , <b>26</b> , 14–26.
643 644 645	<ul> <li>Gosset CC, Do Nascimento J, Augé MT, Bierne N. (2014) Evidence for adaptation from standing genetic variation on an antimicrobial peptide gene in the mussel <i>Mytilus edulis</i>. <i>Mol. Ecol.</i>, 23, 3000–3012.</li> </ul>
646 647	• Goudet J. (2005) Hierfstat, a package for R to compute and test hierarchical F-statistics. <i>Mol. Ecol. Notes</i> , <b>5</b> , 184–186.
648 649	• Hall T. (1999) BioEdit: a user-friendly biological sequence alignment editor and analysis pro- gram for Windows 95/98/NT. <i>Nucleic Acids Symp. Ser.</i> , <b>41</b> , 95–98.
650	• Harrison RG. (1986) Pattern and process in a narrow hybrid zone. <i>Heredity</i> , <b>56</b> , 337–349.
651 652	• Hedrick PW. (2013) Adaptive introgression in animals: examples and comparison to new mutation and standing variation as sources of adaptive variation. <i>Mol. Ecol.</i> , <b>22</b> , 4606–4618.
653 654	• Hilbish TJ, Lima FP, Brannock PM, Fly EK, Rognstad RL, Wethey DS. (2012) Change and stasis in marine hybrid zones in response to climate warming. <i>J. Biogeogr.</i> , <b>39</b> , 676–687.
655 656 657	• Hohenlohe PA, Bassham S, Etter PD, Stiffler N, Johnson EA, Cresko WA. (2010) Population genomics of parallel adaptation in threespine stickleback using sequenced RAD tags. <i>PLoS Genet.</i> , <b>6</b> , e1000862.
658	• Huang X. (1999) CAP3: A DNA sequence assembly program. Genome Res., 9, 868–877.
659 660	• Huerta-Sanchez E, Jin X, Asan, Bianba Z, Peter BM, Vinckenbosch N <i>et al.</i> (2014) Altitude adaptation in Tibetans caused by introgression of Denisovan-like DNA. <i>Nature</i> , <b>512</b> , 194–197.
661 662	• Huson DH, Bryant D. (2006) Application of phylogenetic networks in evolutionary studies. <i>Mol. Biol. Evol.</i> , <b>23</b> , 254–267.
663 664	• Jones FC, Grabherr MG, Chan YF, Russell P, Mauceli E, Johnson J <i>et al.</i> (2012) The genomic basis of adaptive evolution in threespine sticklebacks. <i>Nature</i> , <b>484</b> , 55–61.
665 666 667	• Kenny EM, Cormican P, Gilks WP, Gates AS, O'Dushlaine CT, Pinto C <i>et al.</i> (2010) Multiplex target enrichment using DNA indexing for ultra-high throughput SNP detection. <i>DNA Res.</i> (pp. dsq029).
668 669	• Kocot KM, Cannon JT, Todt C, Citarella MR, Kohn AB, Meyer A <i>et al.</i> (2011) Phylogenomics reveals deep molluscan relationships. <i>Nature</i> , <b>477</b> , 452–456.
670 671	• Koehn RK, Hall JG, Innes DJ, Zera AJ. (1984) Genetic differentiation of <i>Mytilus edulis</i> in eastern North America. <i>Mar. Biol.</i> , <b>79</b> , 117–126.
672 673	• Kriventseva EV, Rahman N, Espinosa O, Zdobnov EM. (2008) OrthoDB: the hierarchical catalog of Eukaryotic orthologs. <i>Nucleic Acids Res.</i> , <b>36</b> , D271–D275.

674 675 676	• Lamichhaney S, Martinez Barrio A, Rafati N, Sundström G, Rubin C-J, Gilbert ER <i>et al.</i> (2012) Population-scale sequencing reveals genetic differentiation due to local adaptation in Atlantic herring. <i>Proc. Natl. Acad. Sci. U.S.A.</i> , <b>109</b> , 19345–19350.
677 678 679	• Li H. (2011) A statistical framework for SNP calling, mutation discovery, association mapping and population genetical parameter estimation from sequencing data. <i>Bioinformatics</i> , <b>27</b> , 2987–2993.
680 681	• Li H, Durbin R. (2009) Fast and accurate short read alignment with Burrows–Wheeler transform. <i>Bioinformatics</i> , <b>25</b> , 1754–1760.
682 683	• Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N <i>et al.</i> (2009) The sequence alignment/map format and SAMtools. <i>Bioinformatics</i> , <b>25</b> , 2078–2079.
684 685	• Luo R, Liu B, Xie Y, Li Z, Huang W, Yuan J <i>et al.</i> (2012) SOAPdenovo2: an empirically improved memory-efficient short-read de novo assembler. <i>GigaScience</i> , <b>1</b> , 18.
686 687	• Mamanova L, Coffey AJ, Scott CE, Kozarewa I, Turner EH, Kumar A <i>et al.</i> (2010) Target- enrichment strategies for next-generation sequencing. <i>Nat. Methods</i> , <b>7</b> , 111–118.
688 689	• Margulies M, Egholm M, Altman WE, Attiya S, Bader JS, Bemben LA <i>et al.</i> (2005) Genome sequencing in microfabricated high-density picolitre reactors. <i>Nature</i> , <b>437</b> , 376–380.
690 691 692	• Martin SH, Dasmahapatra KK, Nadeau NJ, Salazar C, Walters JR, Simpson F <i>et al.</i> (2013) Genome-wide evidence for speciation with gene flow in <i>Heliconius</i> butterflies. <i>Genome Res.</i> , <b>23</b> , 1817-1828.
693 694 695	• Mendez FL, Watkins JC, Hammer MF. (2012) A haplotype at STAT2 introgressed from Nean- derthals and serves as a candidate of positive selection in Papua New Guinea. <i>Am. J. Hum.</i> <i>Genet.</i> , <b>91</b> , 265–274.
696 697	• Mendez FL, Watkins JC, Hammer MF. (2013) Neandertal origin of genetic variation at the cluster of OAS immunity genes. <i>Mol. Biol. Evol.</i> , <b>30</b> , 798–801.
698 699	• Morgulis A, Gertz EM, Schäffer AA, Agarwala R. (2006) WindowMasker: window-based masker for sequenced genomes. <i>Bioinformatics</i> , <b>22</b> , 134–141.
700 701	• Nachman MW, Payseur BA. (2012) Recombination rate variation and speciation: theoretical predictions and empirical results from rabbits and mice. <i>Phil. Trans. R. Soc. B</i> , <b>367</b> , 409–421.
702 703	• Nielsen R, Wakeley J. (2001). Distinguishing migration from isolation: a Markov chain Monte Carlo approach. <i>Genetics</i> , <b>158</b> , 885–896.
704 705 706	• Nolte V, Pandey RV, Kofler R, Schlötterer C. (2012) Genome-wide patterns of natural variation reveal strong selective sweeps and ongoing genomic conflict in <i>Drosophila mauritiana</i> . <i>Genome Res.</i> , <b>23</b> , 99–110.
707 708	• Noor MAF, Bennett SM. (2009) Islands of speciation or mirages in the desert? examining the role of restricted recombination in maintaining species. <i>Heredity</i> , <b>103</b> , 439–444.
709 710	• Nosil P, Funk DJ, Ortiz-Barrientos D. (2009) Divergent selection and heterogeneous genomic divergence. <i>Mol. Ecol.</i> , <b>18</b> , 375–402.

711 712	• Obbard DJ, Welch JJ, Kim KW, Jiggins FM. (2009) Quantifying adaptive evolution in the <i>Drosophila</i> immune system. <i>PLoS Genet.</i> , <b>5</b> , e1000698.
713 714	• Paradis E, Claude J, Strimmer K. (2004) APE: Analyses of phylogenetics and evolution in R language. <i>Bioinformatics</i> , <b>20</b> , 289–290.
715 716	• Pennings PS, Hermisson J. (2006) Soft sweeps II: molecular population genetics of adaptation from recurrent mutation or migration. <i>Mol. Biol. Evol.</i> , <b>23</b> , 1076–1084.
717 718 719	• Philipp EER, Kraemer L, Melzner F, Poustka AJ, Thieme S, Findeisen U <i>et al.</i> (2012) Massively parallel RNA sequencing identifies a complex immune gene repertoire in the Lophotrochozoan <i>Mytilus edulis. PLoS ONE</i> , <b>7</b> , e33091.
720 721	• Piálek J, Barton NH. (1997) The spread of an advantageous allele across a barrier: the effects of random drift and selection against heterozygotes. <i>Genetics</i> , <b>145</b> , 493–504.
722 723	• Poelstra JW, Vijay N, Bossu CM, Lantz H, Ryll B, Müller I <i>et al.</i> (2014) The genomic landscape underlying phenotypic integrity in the face of gene flow in crows. <i>Science</i> , <b>344</b> , 1410–1414.
724 725	<ul> <li>Presgraves DC. (2010) The molecular evolutionary basis of species formation. <i>Nat. Rev. Genet.</i>, 11, 175–180.</li> </ul>
726 727	• Pritchard JK, Di Rienzo A. (2010) Adaptation - not by sweeps alone. <i>Nat. Rev. Genet.</i> , <b>11</b> , 665–667.
728 729 730	• Quesada H, Beynon CM, Skibinski DO. (1995a) A mitochondrial DNA discontinuity in the mussel <i>Mytilus galloprovincialis</i> Lmk: pleistocene vicariance biogeography and secondary intergradation. <i>Mol. Biol. Evol.</i> , <b>12</b> , 521–524.
731 732 733	• Quesada H, Zapata C, Alvarez G. (1995b) A multilocus allozyme discontinuity in the mussel <i>Mytilus galloprovincialis</i> : The interaction of ecological and life-history factors. <i>Mar. Ecol. Prog. Ser.</i> , <b>116</b> , 99–115.
734 735 736	• Quesada H, Wenne R, Skibinski DO. (1999) Interspecies transfer of female mitochondrial DNA is coupled with role-reversals and departure from neutrality in the mussel <i>Mytilus trossulus</i> . <i>Mol. Biol. Evol.</i> , <b>16</b> , 655–665.
737 738 739	• Rawson PD, Hilbish TJ. (1995) Distribution of male and female mtDNA lineages in populations of blue mussels, <i>Mytilus trossulus</i> and <i>M. galloprovincialis</i> , along the Pacific coast of North America. <i>Mar. Biol.</i> , <b>124</b> , 245–250.
740 741	• Rawson PD, Hilbish TJ. (1998) Asymmetric introgression of mitochondrial DNA among European populations of blue mussels ( <i>Mytilus spp.</i> ). <i>Evolution</i> , <b>52</b> , 100–108.
742 743	• Renaut S, Grassa CJ, Yeaman S, Moyers BT, Lai Z, Kane NC <i>et al.</i> (2013) Genomic islands of divergence are not affected by geography of speciation in sunflowers. <i>Nat. Commun.</i> , <b>4</b> , 1827.
744 745	• Riginos C, Cunningham CW. (2005) INVITED REVIEW: Local adaptation and species segre- gation in two mussel ( <i>Mytilus edulis x Mytilus trossulus</i> ) hybrid zones. <i>Mol. Ecol.</i> , <b>14</b> , 381–400.
746 747	• Riginos C, Wang D, Abrams AJ. (2006) Geographic variation and positive selection on M7 lysin, an acrosomal sperm protein in mussels ( <i>Mytilus spp.</i> ). <i>Mol. Biol. Evol.</i> <b>23</b> , 1952–1965.

748 749	• Roesti M, Gavrilets S, Hendry AP, Salzburger W, Berner D. (2014) The genomic signature of parallel adaptation from shared genetic variation. <i>Mol. Ecol.</i> , <b>23</b> , 3944–3956.
750 751 752	• Roesti M, Hendry AP, Salzburger W, Berner D. (2012) Genome divergence during evolutionary diversification as revealed in replicate lake-stream stickleback population pairs. <i>Mol. Ecol.</i> , <b>21</b> , 2852–62.
753 754 755	• Romiguier J, Gayral P, Ballenghien M, Bernard A, Cahais V, Chenuil A <i>et al.</i> (2014) Compar- ative population genomics in animals uncovers the determinants of genetic diversity. <i>Nature</i> , <b>515</b> , 261–263.
756 757 758	• Roux C, Fraïsse C, Castric V, Vekemans X, Pogson GH, Bierne N. (2014) Can we continue to neglect genomic variation in introgression rates when inferring the history of speciation? a case study in a <i>Mytilus</i> hybrid zone. <i>J. Evol. Biol.</i> , <b>27</b> , 1662–75.
759 760	• Seehausen O, Butlin RK, Keller I, Wagner CE, Boughman JW, Hohenlohe PA <i>et al.</i> (2014) Genomics and the origin of species. <i>Nat Rev Genet</i> , <b>15</b> , 176–192.
761 762	• Simpson JT, Wong K, Jackman SD, Schein JE, Jones SJM, Birol I. (2009) ABySS: a parallel assembler for short read sequence data. <i>Genome Res.</i> , <b>19</b> , 1117–1123.
763 764 765	• Song Y, Endepols S, Klemann N, Richter D, Matuschka FR, Shih CH <i>et al.</i> (2011) Adaptive introgression of anticoagulant rodent poison resistance by hybridization between old world mice. <i>Curr. Biol.</i> , <b>21</b> , 1296–1301.
766 767 768	• Soria-Carrasco V, Gompert Z, Comeault A., Farkas TE, Parchman TL, Johnston JS <i>et al.</i> (2014) Stick insect genomes reveal natural selection's role in parallel speciation. <i>Science</i> , <b>344</b> , 738–742.
769 770	• Springer SA, Crespi BJ. (2007) Adaptive gamete-recognition divergence in a hybridizing <i>Mytilus</i> population. <i>Evolution</i> , <b>61</b> , 772–783.
771 772	• Turner TL, Hahn MW, Nuzhdin SV. (2005) Genomic islands of speciation in <i>Anopheles gambiae</i> . <i>PLoS Biol.</i> , <b>3</b> , e285.
773 774 775	• Venier P, De Pittà C, Bernante F, Varotto L, De Nardi B, Bovo G <i>et al.</i> (2009) MytiBase: a knowledgebase of mussel ( <i>M. galloprovincialis</i> ) transcribed sequences. <i>BMC Genomics</i> , <b>10</b> , 72.
776 777	• Via S, West J. (2008) The genetic mosaic suggests a new role for hitchhiking in ecological speciation. <i>Mol. Ecol.</i> , <b>17</b> , 4334–4345.
778 779	• Väinölä R, Hvilsom MM. (1991) Genetic divergence and a hybrid zone between Baltic and North Sea <i>Mytilus</i> populations (Mytilidae: Mollusca). <i>Biol. J. Linn. Soc.</i> , <b>43</b> , 127–148.
780	• Väinölä R, Strelkov P. (2011) Mytilus trossulus in Northern Europe. Mar. Biol., 158, 817–833.
781 782	• Warton DI, Duursma RA, Falster DS, Taskinen S. (2012) smatr 3– an R package for estimation and inference about allometric lines. <i>Methods Ecol. Evol.</i> , <b>3</b> , 257–259.
783 784	• Weir BS, Cockerham CC. (1984) Estimating F-statistics for the analysis of population structure. <i>Evolution</i> , <b>38</b> , 1358–1370.

- Welch JJ, Jiggins CD. (2014) Standing and flowing: the complex origins of adaptive variation.
   *Mol. Ecol.*, 23, 3935–3937.
- Yeaman S. (2013) Genomic rearrangements and the evolution of clusters of locally adaptive loci. *Proc. Natl. Acad. Sci. U.S.A.*, **110**, E1743–E1751.

# 789 Data Accessibility

Sequencing data were deposed in NCBI Short Read Archive, accession numbers [[XXX]] for BAC
 sequences and [[XXX]]] for target-enrichment sequences.

792

# 793 Author Contributions

C.F. and N. B. designed the research. C.F. and K.B. performed the *in silico* work. C.F., J.W. and
N.B. analysed the data. C.F. wrote the article, N.B., J.W. and K.B. revised and commented the article.





**Figure 1** Localities of *Mytilus spp.* samples. We studied three geographical areas (shaded in grey) characterized by transitional zones, denoted  $z_i$ , between genetic backgrounds. (A) Mosaic hybrid zone between *M. edulis* and *M. galloprovincialis* of the Atlantic Coast ( $z_1$ : Normandy,  $z_2$ : South of Brittany,  $z_3$ : Landes; Bierne 2003); followed by the transition with the *M. galloprovincialis* of the two Mediterranean basins ( $z_4$ : Almeria-Oran Front,  $z_5$ : Siculo-Tunisian Strait). (B) Clinal hybrid zone between *M. edulis* and *M. trossulus* ( $z_6$ : Danish Straits; Väinölä & Hvilsom 1991). (C) Mosaic hybrid zone between *M. edulis* and *M. trossulus* ( $z_7$ : Maine and Nova Scotia; Koehn 1984). *M. trossulus* samples are (1) Tvarminne (EU, light green) in the European population of the Baltic Sea and (2) Tadoussac (AM, dark green) in the American population of the Saint-Lawrence River. *M. galloprovincialis* samples are (1) Faro (ATL - external, red) in the peripheral Atlantic population of Iberian Coast, (2) Guillec (ATL - internal) in the enclosed Atlantic population of Brittany, (3) Sete (MED - west, yellow) in the Occidental Mediterranean basin and (4) Crete (MED - east, black) in the Oriental Mediterranean basin. *M. edulis* samples are (1) Holland (EU - external, light blue) in the peripheral European population of the North Sea, (2) Lupin/Fouras (EU - internal, cyan) in the enclosed European population of the Bay of Biscay and (3) Quonochontaug (AM, dark blue) in the American population of the Bay of Biscay and (3) Quonochontaug (AM, dark blue) in the American population of the Bay of Biscay and (3) Quonochontaug (AM, dark blue) in the American population of the Bay of Biscay and (3) Quonochontaug (AM, dark blue) in the American population of the Bay of Biscay and (3) Quonochontaug (AM, dark blue) in the American population of the Bay of Biscay and (3) Quonochontaug (AM, dark blue) in the American population of the Bay of Biscay and (3) Quonochontaug (AM, dark blue) in the American populatio



**Figure 2** Genome-wide relationships between populations. (A) Phylogenetic network, produced with the neighbour-net method (*SplitsTree4*) based on 51,878 high quality SNPs genotyped in 72 individuals. Haplotype sequences were statistically phased with *BEAGLE*. (B) Principal Component Analysis. Results of the first two factorial components are shown. Sites with missing data were removed. (C) Pairwise privately shared SNPs. The size of lines is proportional to the logarithm of the number of SNPs shared. Non-shared private SNPs are indicated to the side of the corresponding population. Singletons were removed. The color code matches Figure 1.



**Figure 3** Distributions of  $FST_{90}$  values between (A) intraspecific and (B) interspecific populations in (1) parapatry and (2) allopatry. Mean values are indicated with vertical stripes underneath the x-axis. Differentiation levels were measured with the 90th percentile of the FST distribution of each contig ( $FST_{90}$ ). Names match Figure 1.



**Figure 4** Joint distributions of interspecific FST<sub>90</sub> values between *M. edulis* and *M. galloprovincialis* pairs in different geographical contexts. Allopatry: (A) European *M. galloprovincialis* with "European *vs* American" *M. edulis* (slope=1.032; elevation=0.049). Parapatry in Europe: (B) *M. edulis* with "Atlantic *vs* Mediterranean" *M. galloprovincialis* (slope=1.138; elevation=0.030); (C) *M. edulis* with "West *vs* East" Mediterranean *M. galloprovincialis* (slope=0.965; elevation=0); (D) *M. edulis* with "Internal *vs* External" Atlantic *M. galloprovincialis* (slope=1.18; elevation=0.020) and (E) *M. galloprovincialis* with "Internal *vs* External" European *M. edulis* (slope=1.043; elevation=0). The Standardised Major Axis regression is indicated in dashed line. All other details match Figure 3.



**Figure 5** Parallelism between the European and American contacts of *M. edulis* and *M. trossulus*. (A) Joint distribution of interspecific  $FST_{90}$  values (slope=0.911; elevation=0.171). (B) Joint distribution of genealogical sorting index (GSI) values (slope=1.634; elevation=0.034). All other details match Figure 4.

**Figure 6** Examples of three candidate outliers in intraspecific comparisons. Top left panel: variation in FST level along the contig; the red paint shows the focal site (i.e. maximal FST value); the red rectangle marks off 1 Kb around the focal site. Top right panel: neighbor network of the contig (*SplitsTree4*) and allele frequency at the focal site in all samples. Bottom panel: neighbour-joining tree 1 Kb around the focal site (R package *ape*). (A) Introgression of *M. edulis* alleles into the European *M. trossulus* population ("Contig54420\_GA36A"). (B) Introgression of *M. galloprovincialis* alleles into the Atlantic *M. edulis* populations ("gi\_403238785\_gb\_JX297444"). (C) Sweep of *M. galloprovincialis* Mediterranean specific alleles ("H\_L1\_abyss\_Contig783"). The contig name is given in brackets (see TableS5 for details). The color code matches Figure 1.

# A. Introgression of *M. edulis* alleles into European *M. trossulus*



Pagenteogression of M. galloprovincialis alleles intro Atlantico Wayedulis



C. Mediterranean *M. galloprovincialis* specific allelasecular Ecology



#### TABLE 1. Intraspecific outliers.

					nh. s	hared	ES.	г.,	FST			nb <sub>con</sub>	s outliers			
							.5190		1 51 max		introgression					
	n a nulation				total	fined		a <b>t</b> l:aa		<b>t</b> l:	into	into	population		TOTAL	
species	population <sub>1</sub>	population <sub>2</sub>	n <sub>1</sub>	n <sub>2</sub>	totai	nxea	genomic	outliers	genomic	outliers	$population_1$	$population_2$	-specific	na	TOTAL	
															_	
M. trossulus	troAM	troEU	16	16	8245	318	0.362	0.693	0.521	1	2	35	3	0	40	
M. edulis	eduAM	eduEU	22	32	417	15	0.131	0.544	0.211	0.802	10	3	10	0	23	
M. galloprovincialis	galATL	galMED	28	32	1673	113	0.107	0.435	0.167	0.625	11	2	6	3	22	
M. galloprovincialis	galATL-external	galATL-internal	12	16	387	0	0.086	0.321	0.139	0.51	1	10	3	0	14	
M. galloprovincialis	galMED-east	galMED-west	16	16	1011	4	0.085	0.333	0.135	0.574	6	3	7	0	16	
M. edulis	eduEU-external	eduEU-internal	16	16	555	0	0.076	0.185	0.134	0.464	0	0	10	4	14	

**n**<sub>1</sub> and **n**<sub>2</sub>: the number of haplotypes sampled; **nb**<sub>SNP</sub> **shared**: the number of SNPs shared between population<sub>1</sub> and population<sub>2</sub> (without singletons); **FST**<sub>90</sub> (and **FST**<sub>max</sub>): the 90th percentile (and maximal value) of the FST distribution of each contig, averaged across all contigs ("genomic") or outlier contigs ("outliers"); **nb**<sub>contigs</sub> **outliers**: total number of outlier contigs, i.e. in the upper 2.5% of the FST<sub>90</sub> or FST<sub>max</sub> distributions. Outliers were categorized according to the cause of differentiation: (1) introgression of foreign alleles into population<sub>1</sub> or population<sub>2</sub>; (2) differentiation of population-specific alleles between population<sub>1</sub> and population<sub>2</sub>; "**na**" stands for outliers that we were unable to classify without ambiguity. In total, data include 1269 contigs; 122,144 SNPs; 471 shared and fixed SNPs.