

The role of recombination and hybridisation in adaptive evolution

A thesis submitted to the School of Environmental Sciences of
the University of East Anglia in partial fulfillment of the
requirements for the degree of Doctor of Philosophy

By Ben J. Ward

Registration number: 1000135331

March 2017

© This copy of the thesis has been supplied on condition that anyone who consults it is understood to recognize that its copyright rests with the author and that no quotation from the thesis, nor any information derived therefrom, may be published without the author's prior, written consent.

© Copyright 2017

by

Ben J. Ward

Abstract

As one of the five evolutionary forces, recombination fulfills both a cleansing role, as well as a role in generating genetic diversity. Recombination cleanses by separating deleterious mutations from their genomic background, increasing the efficacy of purifying selection and curtailing the continuous accumulation of deleterious mutations. Recombination also plays a fundamental role in the repair of damaged DNA, and it can be a creative force, resulting in the formation of novel genotypes, haplotypes and alleles, thereby playing a key role in adaptive evolution. By uniting beneficial mutations that exist at different loci in separate lineages, meiotic recombination during sex accelerates adaptive evolution. Although recombination leaves a distinct signature or footprint in the genome of organisms, identifying this force can be difficult; subsequent recombination events tend to wipe out their past genomic footprints. This thesis presents the development of a novel software package called HybridCheck, for the detection of genomic regions affected by recombination in Next Generation Sequence data, and the rapid molecular dating of recombination events. Hybrid-Check was used to analyze recombination signal in different races of the plant pathogen *Albugo candida*, a generalist obligate biotroph that infects Brassica plants. I show that recombination facilitated occasional introgression and gene flow between host-specialized races. This may have accelerated the rate of adaptive evolution, and possibly broadened

the pathogen's host-range. Finally, the genome of the polar diatom *Fragilariopsis cylindrus* contains diverged alleles that are differentially expressed in different environmental conditions. The hypothesis that ancient asexuality explains how the diverged alleles evolved is challenged, but not rejected, based on evidence of recombination presented in this thesis. An alternative hypothesis is proposed: allelic divergence might have evolved despite the homogenizing effect of meiotic recombination as a result of very large effective population sizes and strong diversifying selection on *F. cylindrus* in the polar environment.

Acknowledgments

I would like to thank my friends and family for their support over the past 4 years. I would also like to thank Professor van Oosterhout for being both a friend and supervisor, for having confidence in my abilities, and for guiding my development into a young scientist, as well as comments and corrections on this manuscript. I thank Professor Kamoun for his guidance and advice about the direction of my PhD, as well as for comments and corrections on this manuscript. I thank Mark McMullan for every comment and discussion ever made over many informal conversations which helped me develop my own understanding. Finally, I thank Dr. Matt Clark and his Plant and Microbial Genomics Team at The Genome Analysis Center for welcoming me and providing the best environment I could ask for, in which to work and write up my thesis in my final year.

Contents

Abstract	v
Acknowledgments	vii
List of tables	xiii
List of figures	xv
1 General Introduction	1
1.1 The five forces of evolutionary change	2
1.1.1 Selection	2
1.1.2 Genetic Drift and finite population sizes	11
1.1.3 Mutation	17
1.1.4 Population structure and gene flow	25
1.1.5 Recombination and linkage	34
1.1.6 Hybrid zones, introgression, and hybrid speciation .	43
1.2 The role of Bioinformatics in population genetics	48
1.2.1 Sequence Alignment	49
1.2.2 Variant Calling	51
1.2.3 Haplotype phasing	52
2 HybridCheck	55
2.1 Introduction	56
2.2 Implementation	64

2.2.1	Four Taxon Tests	64
2.2.2	Sequence triplet scans for recombination signal . . .	67
2.2.3	Estimating the age of recombinant regions	68
2.2.4	Performance Testing	69
2.3	Results	70
2.4	Discussion	72
2.4.1	Performance of detecting recombinant blocks	73
2.4.2	Performance of estimating the age of recombination events	77
3	The role of introgression in the adaptive evolution of the generalist plant pathogen, <i>Albugo candida</i>	80
3.1	Introduction	81
3.2	Methods	88
3.2.1	Isolation and cultivation of races used in the study .	88
3.2.2	Genome assemblies of isolates	89
3.2.3	Detection of recombination events	90
3.2.4	Dating identified recombination events	91
3.3	Results	92
3.3.1	Distribution of polymorphisms across races	92
3.3.2	Recombination blocks identified using RDP	94
3.3.3	Estimated ages of recombination events	95
3.4	Discussion	96
3.4.1	Hybridisation and clonal reproduction of <i>A. candida</i>	97
3.4.2	Biology of genetic introgression and hybridisation . .	99
3.4.3	Introgression and evolution of <i>Albugo candida</i> in the wider context	103
4	Allelic divergence in the polar diatom <i>Fragilariopsis cylindrus</i>	107
4.1	Introduction	108

4.1.1	Sexual reproduction and recombination	108
4.1.2	<i>Fragilariopsis cylindrus</i> and Diatoms	111
4.1.3	The <i>Fragilariopsis cylindrus</i> genome project	115
4.2	Materials and Methods	119
4.2.1	Materials	119
4.2.2	Methods	120
4.3	Results	126
4.3.1	Estimating coalescence times of alleles	126
4.3.2	Testing for recombination in the PCR amplified alleles with the PHI-test	127
4.4	Discussion	133
4.4.1	Sex and the diatom reproductive cycle	135
5	General Conclusion	144
5.0.1	Summary and Conclusions	144
5.0.2	Impact and potential future directions	148
	Appendices	196
A	FALCON assembly haplotype divergence	197

List of tables

1.1	Fitness values for different fitness relationships, adapted from Hedrick 2010.	5
1.2	Expected frequencies for different gametes in a two-allele, two-locus system, adapted from Hedrick 2010.	37
4.1	PHI-Score and Neighbor Similarity Scores of the PCR amplified sequences for three different window sizes.	129

List of figures

2.1	The mean(± 5 - 95%CI) false positive rate (α) of HybridCheck as a function of the ancestral divergence time μt (i.e. the amount of time of the sequences diverged before recombination). As sequences become more diverged, the false positive rate decreases.	71
2.2	The mean(± 5 95%CI) statistical power ($1 - \beta$) of HybridCheck as a function of the divergence time of the sequences after recombination (expressed in μt) for sequences with ancestral divergence times $\mu t = 0.2, 0.4, 0.6$ and 0.8 generations. Recombination between moderately diverged sequences can be detected in $>95\%$ of the cases, as long as the recombination event was relatively recent.	72
2.3	The mean (\pm SEM) estimated age (expressed in μt) of recombinant blocks calculated using the dating algorithm with a JC correction in HybridCheck, versus their actual age. In most of the scenarios, HybridCheck returns an unbiased estimate of the divergence time. However, the age is underestimated in cases of ancient recombination between populations that have ancestral divergence of 0.2	73
3.1	Nucleotide identity amongst the homologous genomic regions of Ac2V, AcBoT and AcNc2. The mean identity was calculated for the sliding window of 20 Kb.	93

3.2	<p>Extensive variation in sequence similarity between <i>Albugo candida</i> races. A) An sequence alignment between base positions 158,779 and 167,382 within contig 1 of <i>A. candida</i> races AcNc2, AcBoT and Ac2V. Two recombination blocks coloured blue and green are visible, displaying high sequence similarity between races. B) The sequence similarity across the length of contig 1, amongst three <i>A. candida</i> races. Similarity is visualised using the colours of a RGB colour triangle in the software HybridCheck. Areas where two contigs have the same colour (yellow, purple or turquoise) are indicative of two races sharing the same polymorphisms. The linear plot of the proportion of SNPs shared between the three pairwise comparisons between the races. Shown on the X-axis is the actual base position.</p>	94
3.3	<p>A) Age of the 675 recombination blocks, identified across the whole genome, estimated using the HybridCheck binomial mass function, assuming a substitution rate of $\mu = 10^6$; B) A box plot of the median (plus first nation blocks and third quartile) log-age of recombination events in contigs. Only contigs with eight or more events are shown. There is no significant difference in age of events between contigs (GLM: $F^{22, 233} = 1.06, p = 0.387$).</p>	96
4.1	<p>Smoothed density plot of the maximum coalescence times (in generations) calculated for allelic pairs of the ABC Iron Transporter (red), Large Ribosomal Subunit (green) and allelic pairs from the genome (blue).</p>	126

4.2	Incompatibility score matrices computed for A). The ABC iron Transporter and B). The Large Ribosomal Subunit. Yellow boxes indicate two informative sites are compatible, and darker boxes indicate the two sites are incompatible. The presence of incompatible sites in the alignments is suggestive of recombination.	127
4.3	Network of simulated allelic pairs, evolved under an asexual reproduction scheme. The first copies of each allelic pair form a clade, and the second copies of each allelic pair form a clade. This is because there is no recombination during gamete formation, as with clonal reproduction, offspring are clones of their parent.	128
4.4	Split Networks of the ABC Iron Transporter and Ribosomal Subunit sequences have average branch lengths close to 10^{-2} and contain 225 splits.	130
4.5	Quantifying the branch lengths and number of splits in networks produced from simulations with varying levels of recombination and values of θ . Larger values of θ cause longer branches (a), and higher recombination rates result in more splits (b).	131
4.6	Networks computed from simulations with three different values of θ . Larger values of θ result in longer outer branches of networks.	132
4.7	Networks computed from simulations with three different levels of recombination, relative to the mutation rate μ . Larger values of R result in more splits in networks.	134

A.1 Sequence similarity calculated with sliding windows across each haplotype 'bubble' in chromosome 000002F, from the *F. cylindrus* FALCON genome assembly. Regions of divergence and indels are apparent. 198

1 CHAPTER 1

2 General Introduction

3 This thesis presents work investigating the role that recombination plays in
4 the adaptive evolution of two eukaryotic microorganisms, *Albugo candida*
5 and *Fragilariopsis cylindrus*. Both of these organisms exist in environments
6 that may be considered very dynamic.

7 In addition, methodological work was also conducted which imple-
8 mented and tested software dedicated to making it easier to detect re-
9 combination in Next Generation Sequencing data. The software was also
10 designed to help solve current methodological issues with distinguishing
11 mosaic regions that are the result of hybridisation, and those that are the
12 result of incomplete lineage sorting.

13 These works are presented in chapters 2, 3, and 4. Each has a more
14 detailed and focused introduction to the concepts specific to them. It is
15 the purpose of this chapter to provide an overview of the key concepts
16 of population genetics that are relevant to this work and provide a wider
17 context for the next three chapters.

18 In order to understand adaptive evolution, it is necessary to understand
19 the five forces of population genetics and how they drive adaptive evolution.
20 What follows is an overview of the five fundamental forces of evolutionary
21 change. Afterwards, an overview of hybrid zones, and an overview of

22 current common Bioinformatics procedures and how they are used in
23 population genetics analyses are presented.

24 **1.1 The five forces of evolutionary change**

25 **1.1.1 Selection**

26 Selection is the non-random, differential survival and reproduction of or-
27 ganisms as a result of their different phenotypes. A population contains
28 many individuals, and these individuals vary in their genetic makeup; the
29 population has genetic variance. This genetic variation, in combination
30 with some environmental effects, is the cause of the phenotypic variation in
31 a population (Ridley 2004). This phenotypic variation results in variation
32 in survival, fecundity, and mating ability, and this ultimately determines
33 whether an individual contributes any alleles to the next generation of that
34 population: Individuals may be better or worse at surviving, or may not be
35 chosen by the opposite sex to mate (Hedrick 2010). This can be expressed
36 in terms of relative fitness. Relative fitness can be defined as the relative
37 ability of different genotypes to pass on their alleles to future generations
38 Charlesworth and Charlesworth 2010. Individuals with genotypes that have
39 a higher relative fitness are expected to survive and pass their alleles on to
40 the next generation, and so over several generations, those genotypes will
41 increase in frequency in the population.

42 **1.1.1.1 The basic diploid model**

43 The basic diploid model of selection models how selection operates for
44 a single diploid locus, with two alleles. The model assumes that there is
45 random mating among individuals in a population, and that selection is
46 operating identically for both sexes. In this model, selection occurs through
47 differences in viability and it is constant through space and time i.e. it acts

48 on every individual in every generation, regardless of location. Generations
49 are discrete and non-overlapping and no mutation is occurring. No gene
50 flow or inbreeding occurs and the size of the population is infinite so there
51 is no genetic drift (Charlesworth and Charlesworth 2010; Hedrick 2010).
52 Despite these assumptions it is still a very useful model to explore and
53 describe how selection operates.

54 Assume there are two alleles of a single locus, denoted as A_1 , and
55 A_2 . With these two alleles, three possible diploid genotypes are possible.
56 Two of them are heterozygous: A_1A_1 , and A_2A_2 , and the third, A_1A_2 is
57 heterozygous. The relative fitnesses of A_1A_1 , A_1A_2 , and A_2A_2 are denoted
58 as w_{11} , w_{12} , and w_{22} respectively (Wright 1937). The contribution of each
59 genotype to the next generation can be calculated as the product of its
60 relative fitness and its frequency prior to selection. The contributions of
61 A_1A_1 , A_1A_2 , and A_2A_2 are $p_0^2w_{11}$, $2p_0q_0w_{12}$, and $q_0^2w_{22}$, where p is defined as
62 the frequency of A_1 and q is defined as the frequency of A_2 (Charlesworth
63 and Charlesworth 2010; Hedrick 2010). Assuming Hardy-Weinberg allele
64 proportions before selection, the mean fitness of the population is:

$$\bar{w} = p_0^2w_{11} + 2p_0q_0w_{12} + q_0^2w_{22} \quad (1.1)$$

65 The frequency of a genotype after selection can be calculated by dividing
66 its contribution by the mean fitness, for example, for A_1A_1 this is $p_0^2w_{11}/\bar{w}$.
67 The frequency of the alleles A_1 and A_2 after selection (p_1 and q_1) can be
68 obtained by noting that the frequency of any of the two alleles is the sum
69 of the frequency of the homozygous genotype and half the frequency of
70 the heterozygous genotype (Charlesworth and Charlesworth 2010; Hedrick
71 2010).

$$p_1 = \frac{p(pw_{11} + qw_{12})}{\bar{w}} \quad (1.2a)$$

$$q_1 = \frac{q(pw_{12} + qw_{22})}{\bar{w}} \quad (1.2b)$$

72 The change in q over one round of selection can be defined as $\Delta q =$
 73 $q_1 - q_0$. Substituting q_1 and simplifying the formula gives equation 1.3
 74 (Charlesworth and Charlesworth 2010; Hedrick 2010).

$$\Delta q = \frac{pq(w_2 - w_1)}{\bar{w}} \quad (1.3)$$

75 If p or q are 0, then there can be no change in frequencies of that allele,
 76 as it is not present in the population (Charlesworth and Charlesworth 2010;
 77 Hedrick 2010).

78 1.1.1.2 Different fitness relationships

79 The formulas and quantities just described can be used to explore the ef-
 80 fects of selection for different fitness relationships. Different relative fitness
 81 values of w_{11} , w_{12} , and w_{22} can be generated for different fitness relation-
 82 ships through the combination of two other coefficients: s is the selection
 83 coefficient which measures the amount of selection against a homozygote,
 84 and h is the level of dominance (Charlesworth and Charlesworth 2010;
 85 Hedrick 2010). When h is multiplied by s , this measures the amount of
 86 selection against a heterozygote (Charlesworth and Charlesworth 2010;
 87 Hedrick 2010). These different fitness relationships are displayed in table
 88 1.1.

89 A recessive lethal allele describes an allele which has a detrimental
 90 effect on the individual that is so severe it leads to death of the individual.
 91 Examples of alleles with such effects include those that cause Tay-Sachs
 92 disease in humans (Myerowitz 1997). Relative fitnesses for this situation

Table 1.1: Fitness values for different fitness relationships, adapted from Hedrick 2010.

Fitness Relationship	A_1A_1	A_1A_2	A_2A_2
Recessive lethal	1	1	0
Recessive detrimental	1	1	$1 - s$
Additive detrimental	1	$1 - (s/2)$	$1 - s$
Purifying Selection	1	$1 - hs$	$1 - s$
Positive Selection	$1 + s$	$1 + hs$	1
Overdominance	$1 - s_1$	1	$1 - s_2$
Underdominance	$1 + s_1$	1	$1 + s_2$

93 are given in row one of table 1.1. Using these values in the formulas 1.1 and
 94 1.3 it can be demonstrated that the mean fitness of a population reaches
 95 1 when there is no A_2 allele in the population ($q = 0$). Furthermore, Δq is
 96 largest when q is large, and is smaller when q approaches 0 (Hedrick 2010).
 97 Therefore, when the frequency of a recessive lethal is high it is purged by
 98 selection very quickly from the population. The reason lethal recessive
 99 alleles are not purged as quickly when they are at low frequency is that they
 100 are present in heterozygotes, therefore the deleterious recessive alleles
 101 are not subject to differential selection (Hedrick 2010).

102 Some recessive alleles are not lethal, but they are detrimental to the
 103 fitness of an individual (Charlesworth and Willis 2009; Charlesworth and
 104 Charlesworth 2010). This type of fitness relationship is called a recessive
 105 deleterious relationship. Fitness values for this scenario are given in row 2
 106 of table 1.1. The selection coefficient (s) reflects how detrimental allele A_2
 107 is. If $s = 1$, then A_2 would be a recessive lethal allele and selection would
 108 act as previously described. Mean fitness is maximized when $q = 0$ and Δq
 109 is greatest when $q_0 = 2/3$, and lower for smaller values of q (Hedrick 2010).
 110 Again this is because A_2 mostly occurs in individuals with a heterozygote
 111 genotype for low q .

112 Heterozygous individuals may have phenotypes that are intermediate to
 113 those of the two homozygotes. If the phenotype of a heterozygote is exactly

114 halfway between that of the homozygotes this is referred to as additivity.
115 Fitness values for additivity are shown on line 3 of table 1.1. In this scenario
116 Δq is larger when both alleles are equally frequent in the population. Δq
117 is greater at low value of q than in the previous scenarios. For low q , A_2 is
118 mostly in heterozygotes, but the deleterious effects of A_2 are not masked in
119 the heterozygotes when the fitness relationship is additive (Charlesworth
120 and Charlesworth 2010; Hedrick 2010).

121 Alleles with additive and recessive effects have been discussed, but
122 every possible level of dominance can be represented in the model with the
123 h coefficient. Fitness relationships modeling different levels of dominance
124 with h are shown on lines 4 and 5 of table 1.1. These fitness arrays
125 describe purifying and positive selection. Purifying selection acts to reduce
126 the frequency of a detrimental allele in a population (Hedrick 2010). In
127 contrast, positive selection acts to increase the frequency of an alleles with
128 effects that are beneficial in the current environment of a population. In
129 reality, selection acts in both positive and purifying roles simultaneously.
130 In both the models if $h = 0$ then the allele is recessive, if $h = 0.5$ it is
131 additive, and if $h = 1$ it is dominant (Charlesworth and Charlesworth 2010).
132 For positive selection, the fastest increase in p occurs when the allele is
133 dominant. When the allele is additive, then p still increases quickly. However,
134 it takes longer for p to increase when A_1 is recessive. At low frequencies,
135 the beneficial A_1 allele typically occurs in heterozygotes, and as a recessive
136 allele, selection does not act on it (Hedrick 2010).

137 In the scenarios previously described selection is a force acting to
138 reduce genetic variation as an allele either increases or decreases in
139 frequency in a population. However, circumstances can cause selection
140 to maintain allelic diversity in the population. This is possible when the
141 heterozygote individuals have a higher fitness than individuals of either of
142 the two homozygote genotypes. The phenomenon is called overdominance.

143 The fitness values for overdominance are listed on row 6 of table 1.1. For
144 selection to maintain both alleles in a population, Δq must be equal to
145 0 for some initial q_0 between 0 and 1 (Charlesworth and Charlesworth
146 2010). This is called the equilibrium frequency of q , and it is a function of
147 both the selection coefficients for the two homozygotes. When q is below
148 this equilibrium frequency, Δq is positive. When q is above the equilibrium
149 frequency, Δq is negative. Thus, as q is perturbed away from this equilibrium
150 Δq shifts such that q will return to this equilibrium (Hedrick 2010). Therefore,
151 both alleles are maintained in the population at a certain ratio.

152 Warfarin resistance in Rats is an example of heterozygote advantage.
153 Resistance was conferred to the rats by a dominant allele (R) at the
154 VKORC1 locus. Individuals with one copy of R were resistant to War-
155 farin, but homozygous individuals had a much greater requirement for
156 Vitamin K (Greaves et al. 1977). Heterozygote advantage has also been
157 invoked to explain polymorphism at loci in the major histocompatibility com-
158 plex (MHC) (Spurgin and Richardson 2010). Overdominance is also an
159 explanation of hybrid vigour (heterosis) (Baranwal et al. 2012) and so this is
160 of particular relevance to chapter 3, where the plausibility of of a generalist
161 plant pathogen evolving through repeated hybridisation is discussed.

162 Underdominance describes the situation where heterozygous individu-
163 als have a lower fitness than homozygous individuals. Fitness values for this
164 relationship are shown on the last line of table 1.1. As with overdominance,
165 there is an equilibrium frequency of q for which $\Delta q = 0$. However, unlike
166 overdominance, with underdominance, Δq is positive above the equilibrium
167 point and negative below it (Hedrick 2010). Therefore the equilibrium is
168 unstable, and allele frequencies move away from it, rather than towards it.

169 1.1.1.3 Selection and dynamic environments

170 The basic model of selection described effectively demonstrates the key
171 concepts of when considering how selection acts. However there are
172 extensions to the model, for example, the model has been extended to
173 account for more than two alleles. Selection is the mechanism that causes
174 adaptive evolution and directional selection and molecular evidence of past
175 positive selection is abundant (Hoekstra and Coyne 2007). Most of the
176 phenotypic characteristics we associate with species are thought to be the
177 end result of selection, even if the adaptive function is not obvious.

178 However, the efficiency of selection can be reduced: Muller introduced
179 the concept of Genetic Load. This is defined as the reduction in fitness
180 from the maximum possible in a population (Davis and Columbia 2011).
181 The principal factors causing genetic load are thought to be the presence
182 of deleterious recessive mutations, maintained by a mutation-selection
183 balance (see section 1.1.3), and the segregation of homozygotes when
184 there is heterozygote advantage (Davis and Columbia 2011). Small isolated
185 populations may suffer from genetic load because they can become fixed
186 for detrimental alleles (see section 1.1.2).

187 Evidence of balancing-selection; selection that maintains polymorphism
188 like overdominance, is not as common (Bubb et al. 2006), but there are sce-
189 narios in which selection does maintain polymorphism. Selection varying in
190 time and space, frequency dependent selection, and host-pathogen evolu-
191 tion, are three such models that are particularly pertinent to the research
192 presented in this thesis as they model selection operating in a dynamic and
193 changing environments. A common aspect of these models is that they
194 violate an assumption of the basic model: constant fitness (Charlesworth
195 and Charlesworth 2010). If constant fitness is not assumed, it can be shown
196 that selection may maintain polymorphism even in absence of heterozygote

197 advantage.

198 Relative fitnesses may depend on the frequency of of the different
199 genotypes in the population. An allele may have a greater fitness when it is
200 present in the population in low numbers and less fitness when it is present
201 in larger numbers (Hedrick 2010; Charlesworth and Charlesworth 2010).
202 This is called negative frequency dependent selection. Alternatively, an
203 allele might increase in fitness as it increases in frequency (Hedrick 2010;
204 Charlesworth and Charlesworth 2010).

205 Frequency dependent selection occurs where there are host-pathogen
206 interactions. Pathogens have genes known as virulence factors and effector
207 genes, which enable them to infect a host. New mutations in a host species
208 that confer resistance to a pathogen will be at low frequencies but have
209 a high selective advantage. As a result, the allele will start to spread in
210 the host population. As the allele becomes more common, the pathogen
211 will find fewer new hosts they can infect (Charlesworth 2006; Frank 1993;
212 Seger and Antonovics 1988). Therefore, pathogen numbers decrease and
213 the advantage gained by being resistant diminishes. Indeed, if there is a
214 cost to maintaining the resistance it will even become detrimental. This
215 process also happens with the pathogens. As hosts acquire resistance to a
216 pathogen, pathogens with new mutations allowing them to infect previously
217 resistant hosts will have a strong selective advantage. The now susceptible
218 host genotype will decrease in frequency, as the pathogen increases in
219 frequency. The selective advantage of the pathogen genotype is reduced
220 and may even suffer a cost if it is less virulent than other pathogen geno-
221 types at infecting other host genotypes (Charlesworth 2006; Frank 1993;
222 Seger and Antonovics 1988). Parasite genotype frequencies may therefore
223 become balanced in a population, resulting in highly polymorphic genes
224 in pathogens, such as antigenic genes in malaria, and effector genes in
225 pathogens like *Phytophthora infestans*(Morgan and Kamoun 2007; Policy

226 and Conway 2001). This type of process, typically assuming gene-for-gene
 227 interactions between host and pathogen, leads to cycles of allele frequency
 228 changes in both the host and pathogen (May and Anderson 1983). This may
 229 be of particular importance to haploid pathogens which by definition, will
 230 not have their polymorphism maintained by heterozygote advantage, and
 231 may be subject to clonal interference which restricts levels polymorphism
 232 and the speed of adaptation (Gerrish and Lenski 1998).

233 In addition to existing in balance, polymorphisms in a host or pathogen
 234 pathogen can become fixed due to their selective advantage, which can
 235 lead to a succession of fixation events in both host and pathogen as each is
 236 under selection pressure to counter adapt each others previous adaptations.
 237 This is called an evolutionary arms race, and can lead to long term variability
 238 and rapid evolution of DNA sequences such as effector genes in plant
 239 pathogens, and R genes in plants, and accelerated molecular evolution
 240 (see chapter 3) (Brown 2003; Charlesworth 2006; Morgan and Kamoun
 241 2007; Paterson et al. 2010; Rose et al. 2004).

242 Selection may maintain variation when there is enough temporal varia-
 243 tion in relative fitnesses of different genotypes. An allele with detrimental
 244 effects in one generation may confer an advantage in subsequent genera-
 245 tions, should conditions change. This scenario is pertinent to chapter 4 as
 246 the environment of *Fragilariopsis cylindrus* is also temporally dynamic with
 247 seasonal changes such as freezing and thawing events. Models of tempo-
 248 rally changing fitnesses have shown that polymorphism is only maintained
 249 by selection under very strict conditions: The geometric mean of fitness
 250 over n generations for both homozygotes must be smaller than that of the
 251 heterozygote (equation 1.4) (Haldane and Jayakar 1963).

$$\left(\prod_{i=1}^n w_{11 \cdot i} \right)^{1/n} < 1 > \left(\prod_{i=1}^n w_{22 \cdot i} \right)^{1/n} \quad (1.4)$$

252 This can be illustrated by considering two seasons, A_1 is advantageous
253 in one season, and A_2 is advantageous in the other. Fitness values in
254 season one then are $1+s$, 1 , and $1-s$ for A_1A_1 , A_1A_2 , and A_2A_2 respectively.
255 In the second season, this is reversed and A_1A_1 has fitness $1-s$ and A_2A_2
256 has fitness $1+s$. If the same number of generations is spent in each season,
257 conditions for polymorphism are met, otherwise directional selection will
258 result instead. Such expectations from theory have been validated in
259 experimental evolution studies with bacteria, where serial transfer regimes
260 were used to emulate the effects of temporal variation (Rainey et al. 2000).
261 Therefore, it seems that there is little evidence polymorphism is maintained
262 by selection where fitnesses vary in time, without heterozygote advantage
263 or frequency dependent selection.

264 1.1.2 Genetic Drift and finite population sizes

265 Genetic drift is the chance changes in allele frequency that result from
266 the random sampling of gametes from generation to generation in a finite
267 population.

268 1.1.2.1 The effect of drift

269 Genetic drift has the same expected effect on all loci in a genome. In a
270 large population, on average only a small change in allele frequencies
271 will occur as a result of genetic drift. However, for smaller populations,
272 genetic drift can cause larger fluctuations in allele frequencies and may
273 even lead to the loss of fixation of alleles purely by chance alone (Hedrick
274 2010; Charlesworth and Charlesworth 2010). Simulations of genetic drift
275 reveal that small population sizes can cause replicate populations to drift
276 apart in allele frequency. The probability that an allele goes to fixation
277 as a result of genetic drift in a finite population is proportional to its initial
278 frequency, assuming differential selection is not occurring. $u(q) = q_0$ Over

279 replicate simulated populations, the mean allele frequency does not change
 280 as a result of drift, but the distribution of allele frequencies over replicate
 281 populations does (Hedrick 2010; Charlesworth and Charlesworth 2010).
 282 Therefore, drift is often examined by considering heterozygosity or the
 283 variance in allele frequencies of replicate populations.

284 Consider a Wright-Fisher model population with N (diploid) individuals
 285 and assume each contributes two haploid gametes to the next generation
 286 (Crow and Kimura 1970). For an offspring individual, the probability of draw-
 287 ing the same allele twice from the parents is $2N[1/(2N)]^2$. The probability
 288 that they are different is $1 - 1/(2N)$. Two alleles may also be identical by
 289 descent with probability:

$$f_{t+1} = \frac{1}{2N} + \left(1 - \frac{1}{2N}\right)f_t \quad (1.5)$$

290 This can be rewritten and the expected heterozygosity after t genera-
 291 tions derived:

$$H_{t+1} = \left(1 - \frac{1}{2N}\right)H_t \quad (1.6a)$$

$$H_t = \left(1 - \frac{1}{2N}\right)^t H_0 \quad (1.6b)$$

292 This demonstrates that each generation, heterozygosity decreases at a
 293 rate that is an inverse function of the population size, and it is possible to
 294 calculate the expected heterozygosity after t generations (Hedrick 2010;
 295 Charlesworth and Charlesworth 2010). In addition, it is possible to relate
 296 observed, heterozygosity to the difference in expected heterozygosity and
 297 the variance in allele frequency. Taking account of this into the above
 298 equations and rearranging produces a formula for for the variation in allele

299 frequencies at time t . The formula shows that as the number of generations
300 increases, the variance approaches a maximum value of p_0q_0 . This Wright-
301 Fisher model assumes parents produce many gametes and zygotes, and of
302 those N are chosen to form the next generation. It is implicit that individuals
303 are hermaphrodites and there is a small probability of self-fertilization.
304 The mean time until fixation of an allele due to drift depends on initial
305 frequencies of the allele and the initial frequency of the allele (Hedrick
306 2010; Charlesworth and Charlesworth 2010). As population size increases,
307 the effect of drift becomes smaller as it takes more consecutive chance
308 increases of an allele to fix it in the population. For any given population
309 size, the lower the initial allele frequency is, the longer it is for that allele
310 to become fixed by drift. With new neutral mutants, the expected time to
311 fixation is four times the population size.

312 Explanations of drift often mention the population size N . However, in
313 many situations the relevant value is the number of breeding individuals.
314 This may be very different from the census population size. The concept of
315 an effective population size makes it possible to consider an ideal population
316 of size N in which all parents have an equal expectation of being a parent
317 of any individual progeny. i.e the Wright-Fisher model. Effective population
318 size can be measured by different methods: inbreeding, variance, and
319 eigenvalue. When a population remains the same size these measures are
320 similar, however they may differ when populations are growing or shrinking
321 (Kimura and Crow 1963; Waples 2002). The effective population size can
322 be influenced by the frequency of different sexes in a population, variance in
323 reproduction, and varying numbers of individuals over several generations.

324 Bottlenecks and founder events are two specific cases where a popula-
325 tion changes size significantly, influencing the effective population size. A
326 bottleneck describes a situation in which something occurs to drastically

327 reduce the number of individuals which survive in a population, or other-
328 wise get to contribute to the next generation of the population. Typically,
329 these are events such as natural disasters, overwintering, or epidemics. A
330 founder event describes a situation in which a population is started from a
331 low number of individuals, for example individuals being carried to a new
332 island or location. In both cases, these events can cause large random
333 changes in allele frequencies, resulting in lower heterozygosity and fewer
334 alleles than the ancestral population. The changes in allele frequencies
335 resulting from bottlenecks and founder events generate genetic distance
336 between two populations, equation 1.7 gives the standard genetic distance
337 (Nei 1987) after a bottleneck or founder event, where t is the the number of
338 generations the event lasted (Chakraborty and Nei 1977).

$$D_t = -\frac{1}{2} \ln \left(\frac{1 - H_0}{1 - H_t} \right) \quad (1.7)$$

339 1.1.2.2 Drift and selection

340 In a finite population, when there is no differential selection at a locus, an
341 allele may become fixed or lost as a result of genetic drift.

342 In a population of infinite size, by definition there is no genetic drift,
343 and selectively favored alleles increase in frequency and asymptotically
344 approach fixation. Detrimental alleles always reduce in frequency and
345 approach loss. In finite populations however, because of the effects of
346 genetic drift, alleles may not always be fixed when they are favorable, and
347 detrimental alleles may be fixed despite their detriment. The probability of a
348 favorable allele in a finite population is a function of the initial frequency of
349 the allele, the extent to which selection favours that allele, and the size of
350 the population. Kimura 1962 developed an equation that takes these factors
351 to compute the probability of fixation of A_1 (Kimura and Ohta 1971). The
352 probability of fixation of an allele is a function of its initial frequency, the level

353 of dominance, the effective population size, and its selective advantage.
354 The probability of fixation of an allele increases with increasing initial allele
355 frequency and with increasing Ns (the product of population size and
356 selection coefficient). When $Ns \ll 1$, this indicates that $s \ll 1/N$ and
357 that the selective advantage of an allele is very low. In this case, changes
358 in allele frequency are determined by drift. When $Ns \gg 1$, then s is higher
359 than $1/N$ and changes in allele frequency depend more on selection than
360 on drift. The effect where alleles with low selection coefficients (and hence
361 only slightly deleterious effects), may act as if they were neutral in small
362 populations was first identified by Wright 1931, and described in terms of
363 molecular evolution by Ohta 1973, who called it the nearly neutral model.

364 In a neutral situation in a finite population, the loss of heterozygosity
365 is $1 - 1/(2N)$. For any given balancing selection regime, the decay in
366 heterozygosity can be defined as $H_{t+1} = (1 - d)H_t$, where d is the loss
367 from unfixed allele frequency states and the gain for the absorbing states.
368 With no selection, d is $1/(2N)$ i.e. the expression reduces to the neutral
369 model of heterozygosity loss as a result of drift already described. The
370 ratio of decay for a neutral locus over one undergoing selection is called
371 a retardation factor (Robertson 1962). This factor is one when there is
372 neutrality, but when d is less than 1, then selection can slow the rate of
373 fixation, or when $d > 1$, then selection is increasing the rate of fixation.
374 Even though selection may be balancing in an infinite population, in a finite
375 population, less genetic variation may be retained than in a population
376 with no selection. Populations with heterozygote advantage, and unequal
377 homozygote fitness values genetic variation is eliminated faster than in
378 populations with neutrality.

379 1.1.2.3 Impact of genetic drift

380 Genetic drift needs to be considered when studying plant pathogens and
381 organisms in very dynamic environments, as those populations may ex-
382 perience periodic population expansions or contractions. Analysis of Q_{ST}
383 values of eight traits, and F_{ST} values of eight neutral loci of the pathogenic
384 fungus *Rhynchosporium commune* revealed that the majority of the traits
385 analysed were evolving according to stabilizing selection, although a trait
386 for growth at 22 degrees centigrade was subject to diversifying selection
387 and local adaptation (Stefansson, McDonald, and Willi 2014). This was
388 proposed to be due to the fact the pathogen exists in large rather homoge-
389 neous environments (i.e. homogeneous monoculture systems) where they
390 mostly experience one host genotype, and therefore stabilizing selection
391 plays a greater role than does drift or directional selection. Furthermore, the
392 cycles of frequency dependent selection and maintenance of diversity previ-
393 ously described would only be expected to occur if there were some allelic
394 diversity - rare advantageous alleles - in the host. Other plant pathogens
395 have been significantly affected by changes in their population size. For
396 example, the global pandemic of *Phytophthora infestans* was initiated by
397 a single clone, which escaped to North America, and then to Europe, and
398 then to the rest of the world (Goodwin, Cohen, and Fry 1994). Analyses
399 of RFLP loci of the pathogen *Mycosphaerella graminicola* isolated from
400 different locations, indicated that Mexican and Australian populations have
401 low gene diversity (Zhan, Pettway, and McDonald 2003), consistent with
402 founder events and genetic drift. Steele et al. 2001 found that in Australia,
403 *Puccinia striiformis* originates from a single founder event, the founding race
404 identified corresponded to a race previously identified in Europe.

405 1.1.3 Mutation

406 Mutation is the alteration of the nucleotide sequence of the genome of
407 an organism. Mutations may be caused by errors in the DNA replication
408 process, the insertions of a transposable element, chromosome breakage,
409 and errors in meiosis. Mutations may be caused by chemicals or
410 radiation, and these mutagens cause certain kinds of mutation, for example,
411 ultraviolet light (Kozmin et al. 2005).

412 Many spontaneous mutations may have detrimental effects as they affect
413 the normal functioning of a gene. However, many mutations have neutral
414 or almost neutral effects, as they do not result in changes to proteins or
415 otherwise change DNA only slightly (Grauer and Li 2000). A few mutations
416 will confer beneficial effects and change proteins in a way that enhances
417 the fitness of organism with the allele. Of course whether or not a mutant is
418 beneficial, deleterious, or neutral also depends on the environment (Grauer
419 and Li 2000).

420 Typically, the term mutation is often used to describe the smaller scale
421 mutations which give rise to a new allele or sequence, larger alterations
422 are often referred to as copy number variations, structural variations, or
423 chromosomal abnormalities (Grauer and Li 2000; Hedrick 2010). A mutation
424 may involve a change in one nucleotide base, or it may involve changes in
425 several nucleotides. Short mutations where a few nucleotides are removed
426 or inserted into the DNA sequence are called indels, which may cause
427 a frame-shift mutation if the number of bases inserted or deleted is not
428 a multiple of three. The change affects the grouping of nucleotides into
429 codons, affecting the reading frame or possibly introducing a stop codon.
430 Both base mutations and indels can cause a change in the protein produced
431 transcription and translation of the gene (Grauer and Li 2000). Transposable
432 elements are portions of DNA that can replicate themselves and move

433 location within the genome of an organism (Grauer and Li 2000; Wicker et al.
434 2007). 60% of the maize genome and 15% of the *Drosophila melanogaster*
435 genome consists of transposable elements (Biémont and Vieira 2006).
436 Transposable elements have been characterized as junk, neutral, and
437 agents of mutation and adaptation. Their behavior ranges from that of
438 an extreme parasite, to that of a mutualist depending on the transposable
439 element, the organism, and the area of the genome affected by one (Grauer
440 and Li 2000).

441 To understand genome evolution, mutation by gene duplication, deletion,
442 and gene conversion are important. Many genes such as globins, histones,
443 enzymes, and MHC genes are members of multigene families. Such
444 families are composed of several homologous genes, with similar function,
445 and are often situated close together on a chromosome i.e. they are
446 closely linked (Hedrick 2010). Such multigene families are thought to
447 have evolved through serial duplication of an ancestral gene. Duplicate
448 genes may cause dosage effects, or they may diverge, resulting in new
449 functionality (neofunctionalisation), or they may retain only a subset of their
450 original functionality (subfunctionalisation). Further duplication and deletion
451 of genes may occur through unequal crossing over or gene conversion
452 (Grauer and Li 2000). Gene conversion is a process by which the nucleotide
453 sequence of one allele or allele segment is replaced by a homologous
454 sequence from another allele. Voordeckers et al. 2012 demonstrated
455 how the MALS family of genes, which code for proteins specialised to act
456 on disaccharides, were likely to have evolved through duplication of an
457 ancestral gene. By reconstructing the ancestral genes, and testing their
458 activity on different substrates, they found the ancestor was mostly active
459 on maltose like substrates, but had some function on isomaltose like sugars.
460 Duplication and mutation resulted in a series of enzymes specialised for
461 different substrates. Many species of plant pathogens have genomes rich

462 in both repeats and transposable elements (Raffaele and Kamoun 2012;
463 Kemen and Jones 2012) and it is therefore suspected they play a role in
464 the evolution of effector repertoires and can influence the expression of
465 effectors (Whisson et al. 2012).

466 Mutations may occur anywhere across the genome stochastically, ac-
467 cording to a mutation rate, however there are hotspots in the genome which
468 experience mutations more often than other regions. Research into *E.coli*
469 by Shee, Gibson, and Rosenberg 2012 has indicated such hotspots can
470 be caused by double strand breaks in DNA which then lead to stress in-
471 duced mutagenesis. In the plant pathogen *Neurospora crassa* duplicate
472 sequences in DNA are detected and mutated during its sexual phase. The
473 mechanism could cause linked duplicated genes to diverge further than
474 unlinked ones (Cambareri, Singer, and Selker 1991).

475 It is often assumed that likelihood of mutation occurring is unaffected
476 by selection, however there are exceptions. In microorganisms it is known
477 that mutator phenotypes can arise (Barrick et al. 2009). These increase
478 the number of mutations occurring in the population, and facilitate the
479 adaptation of large asexual populations to new conditions, even when the
480 frequency of the mutators is low. Such hyper-mutation can be genetically
481 inherited, or can be transient. Clinical isolates of many pathogens such as *E.*
482 *coli*, *Streptococci spp.*, and *Staphylococci spp.* have been found to contain
483 high proportions of hypermutators (Jayaraman 2011). Localization of the
484 hyper-mutation to contingency genes or specific regions of the genome
485 limit the risk of accumulating too many detrimental mutations through hyper-
486 mutation (Jayaraman 2011). In the case of an inheritable hyper-mutator
487 allele, it may increase in frequency in a population through hitchhiking; it
488 is physically linked to a selectively beneficial mutation it caused to occur
489 (Giraud et al. 2001). Several models demonstrating how hypermutators
490 persist and succeed exist (Taddei and Radman 1997; Tenailon et al. 1999),

491 and Hyper-mutation is particularly beneficial strategy for microorganisms
492 that are exposed to frequent and possibly unpredictable stresses (like
493 pathogens) (Visser 2002; Tanaka, Bergstrom, and Levin 2003).

494 Mutation is an important evolutionary force that generates the variation
495 the other forces act on. Several mechanisms in microbes and pathogens
496 have been described through which such variation is generated, in addition
497 to ways in which an organism might increase the rate at which this variation
498 is generated during times of stress for for certain alleles. Next the effects
499 mutation has on populations and how it exists in balance with previously
500 described forces is presented.

501 1.1.3.1 Effect of mutations on populations

502 The effect of mutation on population allele frequencies can be evaluated
503 by assuming a forward-backward model of mutation (Hedrick 2010). In
504 this model, two types of allele are possible, a wild type allele (A_1) and a
505 detrimental mutant (A_2). In addition, mutation is reversible and may change
506 wild type alleles to the mutant alleles (forward mutation), and the mutant
507 alleles may mutate back to the wild type (backward mutation). It is assumed
508 forward mutations are more common than backward mutations. This is
509 because forward mutations are mutations that resulting in gene malfunction.
510 It is assumed only a limited number of possible mutations could compensate
511 for such forward mutations and result in a backward mutation. Mutation
512 from A_1 to A_2 occurs at a rate u , and mutation from A_2 to A_1 occurs at rate
513 v . The change in frequency of A_2 due to only mutation is $\Delta q = up - vq$. This
514 expression is linearly related to the allele frequency, but as u and v are small
515 - mutation rates are typically low - mutation does not significantly affect the
516 proportion of alleles in the population (Hedrick 2010). An equilibrium is
517 achieved if the forward and backward mutation rates are equal, and if u is
518 higher than v then it is expected that the frequency of detrimental alleles

519 would be higher than the wild type alleles (Hedrick 2010). However this
520 expectation is not realistic as it does not consider selection.

521 When mutations occur, they are the only copy in the entire population. All
522 the individuals in the population immediately after mutation are homozygous
523 for the wild type allele (A_1A_1), and the mutant is heterozygous (A_1A_2). This
524 one heterozygous individual must mate with a homozygous individual. The
525 new mutant may be lost, only homozygous wild type offspring may be the
526 outcome, or some offspring may be heterozygous with the new mutant
527 allele. If mating results in only one offspring, then there is a 50% chance it
528 is A_1A_1 , and if A_1A_2 is the result, then there is still only one A_1A_2 individual
529 in the population. If mating results in two offspring, then the probability of
530 losing A_2 is halved. So the frequency of A_2 in generations following the
531 mutation event depends on how many progeny are the result of mating,
532 and what type they are (Hedrick 2010).

533 The way in which purifying selection keeps detrimental alleles from
534 increasing in frequency has previously been described. The entire genome
535 is subject to the opposite effects of mutation and selection, and the joint
536 effects of mutation and selection is called the mutation-selection balance.
537 Assume that A_2 is deleterious and recessive, selection will act to reduce the
538 frequency of A_2 as previously described. Equation 1.8 rewrites 1.3 using
539 the fitness values for a recessive deleterious allele from table 1.1 (Hedrick
540 2010).

$$\Delta q_s = \frac{sq^2p}{1 - sq^2} \quad (1.8)$$

541 The increase in q due to mutation then is $\Delta q_{mu} = up$, and assuming
542 back mutation occurs at a low rate compared to u , as these forces have
543 opposite effects, there is a point where they are at equilibrium (equation 1.9)
544 and the total change in allele frequency is $\Delta q = \Delta q_{mu} + \Delta q_s = 0$ (Hedrick

545 2010).

$$up = \frac{sq^2p}{1 - sq^2} \quad (1.9)$$

546 If it is assumed that q^2 is small then equation 1.9 can be solved for
 547 the equilibrium genotype frequency ($q_e^2 = u/s$), and the equilibrium allele
 548 frequency ($q_e = \sqrt{u/s}$). This frequency is increased as a result of either
 549 higher mutation rate or lower selective disadvantage. If the deleterious
 550 mutant were not completely recessive, the level of dominance h can affect q_e .
 551 If h is much larger than 0 and q_e is small, then equilibrium allele frequency
 552 is approximately u/hs , and assuming p is almost 1, the frequency of the
 553 mutant phenotype at equilibrium is $2u/s$. As a general rule, as the level
 554 of dominance increases, the equilibrium allele frequency rapidly reduces
 555 (Hedrick 2010).

556 Mutations will contribute to the genetic load of a population, reducing its
 557 fitness from the maximum possible. For a deleterious recessive mutation
 558 the load is $L = sq^2$ and at equilibrium $u = sq^2$, load is roughly equal
 559 to the mutation rate. If the deleterious mutant is dominant, then load
 560 becomes $L = 2u$ which shows that depending on the level of dominance,
 561 the mutation load can be between the mutation rate and twice the mutation
 562 rate. If independence of fitness between loci is assumed, the fitness at
 563 locus i may be defined as \bar{w}_i , and the overall fitness of the population is
 564 defined as $\bar{w} = \bar{w}_i^n$. The overall load is $L = 1 - \bar{w}$. Crow and Kimura 1970
 565 gave a formula for approximating the total load caused by mutation:

$$L \approx C \sum u_i \quad (1.10)$$

566 Where C is a constant between 1 and 2 and u_i is the mutation rate of
 567 the locus i .

568 Joint consideration of mutation and drift forms the basis of the neutral

theory. The initial frequency of a new mutant A_1 in a population of A_2 alleles has an initial frequency of $p_0 = \frac{1}{2N}$. The two alleles are neutral respective to each other, thus the probability of this mutant being fixed in the population is equal to its initial frequency as described in section 1.1.2, and the probability of losing the mutant from the population is $u(q) = 1 - \frac{1}{2N}$. Unless a population is very small, a new neutral mutation is likely to be lost from the population by drift alone (section 1.1.2). Loss of a mutant due to drift occurs more quickly than fixation. This is because the change in frequency necessary to lose a new mutant is much smaller than that necessary to fix the new mutant. Kimura and Ohta 1971 formulated the average time to fixation and loss of a new mutant due to drift alone:

$$T_1(p) = 4N_e \quad (1.11a)$$

$$T_0(p) = 2 \left(\frac{N_e}{N} \right) \ln(2N) \quad (1.11b)$$

Assuming $N = N_e$ then the time to loss reduces to $2N/[\ln(2N)]$. As a result, polymorphism is often transient. Mutation acts to increase the number of alleles, whereas drift acts to reduce the number of alleles. The properties of this equilibrium for the infinite alleles model were explored by Kimura and Crow 1964 using the inbreeding coefficient. Recall that equation 1.5 gives the expected inbreeding coefficient. This may be modified by the probability both alleles did not mutate:

$$f_t = \left[\frac{1}{2N_e} + \left(1 - \frac{1}{2N_e} \right) f_{t-1} \right] (1 - u)^2 \quad (1.12)$$

Setting $f_0 = 1$ (heterozygosity $H_0 = 0$) and $u = 10^{-5}$ and examining the change in heterozygosity over many generations for various values of N_e it can be shown that it takes many generations, but eventually heterozygosity rises to approach an asymptotic value. Furthermore, the asymptotic level of heterozygosity is greater when N_e is greater. As a consequence, when

593 population size is small, the rise to the smaller asymptotic value occurs
594 more quickly as genetic drift has a greater impact on the genetic variation
595 change than does mutation (Kimura and Crow 1964; Hedrick 2010). If an
596 equilibrium between mutation adding variation and drift eliminating variation
597 from a population is assumed $f_t = f_{t-1} = f_e$, formula 1.12 reduces to:

$$f_e \approx \frac{1}{4N_e u + 1} \quad (1.13)$$

598 Because $H = 1 - f$, equilibrium heterozygosity for the infinite allele
599 neutral model can be obtained, where $\Theta = 4N_e u$:

$$H_e = \frac{\Theta}{\Theta + 1} \quad (1.14)$$

600 This equilibrium is different to equilibrium previously described, as the
601 allele frequencies are constantly changing, but the distribution of alleles
602 remains mostly constant. The above equation demonstrates that when
603 $\Theta \approx 1$, then $H_e \approx 0.5$. When $\Theta \gg 1$ then mutation primarily affects
604 heterozygosity rather than drift and so H_e is quite high. The opposite is
605 true, when $\Theta \ll 1$ then drift is the major determinant of heterozygosity and
606 H_e is low (Kimura and Crow 1964; Hedrick 2010).

607 To examine the effect of a population bottleneck, assume a population
608 starts at mutation-drift equilibrium. The population goes through a bot-
609 tleneck and grows large once again (Nei 2005). The expected genetic
610 variation after the bottleneck depends on heterozygosity prior to the bottle-
611 neck, the size of the bottleneck, and the rate of increase after the bottleneck
612 (Nei, Maruyama, and Chakraborty 1975). The size of the bottleneck has a
613 large effect on the number of alleles in a population, but average heterozy-
614 gosity is mostly affected by the rate of growth after the bottleneck. This is
615 because whilst heterozygosity is reduced by the decrease in population size,
616 when growth of the population after the bottleneck is slow, heterozygosity

617 is lost each generation until it is large enough. Faster population growth
618 rates allow populations to rebound as loss of heterozygosity only occurs
619 during the first few generations following the bottleneck (Nei, Maruyama,
620 and Chakraborty 1975).

621 Mutations can have selective effects. When s is less than $1/(2N)$ genetic
622 drift is the stronger factor affecting allele frequency than selection and the
623 mutant behaves neutrally, and deleterious mutants may become fixed
624 as if they were neutral in small populations (Kimura 1983; Lynch and
625 Gabriel 1990; Lande 1994). Over time, fitness declines which can lead to
626 further reductions in population size, and hence mutations of increasingly
627 detrimental effect behave as if they are neutral, and are more likely to be
628 fixed. Such a feedback is called mutation meltdown, and in theory could
629 make small populations go extinct, (Lynch, Conery, and Burger 1995).

630 **1.1.4 Population structure and gene flow**

631 Populations may be split into subpopulations due to geographical, eco-
632 logical, or behavioral factors. When a population is divided or there is
633 more than one population, the amount of genetic exchange, or gene flow,
634 between the subpopulations may differ between the different populations
635 or subpopulation. When gene flow is high between two populations or
636 subpopulations, they are highly connected genetically and the amount of
637 genetic variation between them is homogenized. Conversely, when the
638 amount of gene flow is low between populations or subpopulations, then
639 genetic drift, selection, and mutation in the populations and subpopulations
640 may lead to genetic differentiation (Charlesworth and Charlesworth 2010;
641 Hedrick 2010).

642 Some types of movement of individuals like migrations will not actually
643 result in gene flow, especially if the individual is only transiently passing
644 through a population and does not breed with members of the population

645 (Hedrick 2010). Gene flow may be distinguished from simple migration as
646 movement between groups that results in genetic exchange (Endler 1977).

647 When considering population subdivision it is often assumed that the
648 subpopulations are always present. Another view assumes they can die
649 out, but they are repopulated from neighboring subpopulations, this is
650 termed a metapopulation (Hanski 1998), and the dynamics of extinction
651 and re-population make metapopulations differ from the basic concept of a
652 subdivided population. What follows is a basic description of how gene flow
653 effects populations using a simple genetic model, before the joint effects of
654 gene flow and drift, and gene flow and selection are considered.

655 The continent-island model models a situation in which a large continent
656 population is connected to a smaller island population (Charlesworth and
657 Charlesworth 2010). The smaller island population receives migrants from
658 a larger continent population. The larger continent population is assumed
659 to be large enough to render the effect of genetic drift negligible compared
660 to the effect of gene flow. Gene flow is assumed to have negligible effect
661 on the source population. In this model, the proportion of migrants moving
662 to the island is m , and the proportion of residents in the island population is
663 $1m$. The proportion of A_2 in the migrants coming from the continent is q_m
664 and the frequency of A_2 on the island before the gene flow is q_0 (Hedrick
665 2010).

666 Frequency of A_2 on the island after gene flow is calculated as:

$$q_1 = (1 - m)q_0 + mq_m \quad (1.15)$$

667 Formula 1.15 can be reduced to $q_0 - m(q_0 - q_m)$.

668 The change in frequency of q is then defined as:

$$\Delta q = q_1 - q_0 \quad (1.16)$$

669 Formula 1.16 reduces to $-m(q_0 - q_m)$.

670 q_m and m are assumed to be constant (Hedrick 2010). From these
671 equations it is clear that $m = 0$ then there is not migration from the continent
672 to the island and so there is no change in allele frequency. If $q_0 < q_m$ then
673 the frequency of q increases on the island. If $q_0 > q_m$ the frequency
674 decreases. This indicates that there is a stable equilibrium freq of A_2 at
675 $q_m = q_0$.

676 A general formula to calculate the frequency of A_2 for any generation t
677 has been derived as:

$$q_t = (1 - m)^t q_0 + [1 - (1 - m)^t] q_m \quad (1.17)$$

678 In this formula, as t increases the first term approaches 0, and the
679 second term approaches q_m (Hedrick 2010). Therefore eventually the
680 frequency of A_2 in the island population converges to the frequency of A_2
681 in the continent population. This is because gene flow is unidirectional,
682 and therefore eventually all in the island population are descended from
683 migrants. Thus, the allele frequencies approach that of the continent i.e.
684 the source of the migrants (Charlesworth and Charlesworth 2010). In this
685 model, allele frequency changes at a maximum rate initially, and as the
686 equilibrium is approached, it decreases.

687 A more general model assumes gene flow can occur among all parts of
688 a structured population. The model assumes there is k different subpopula-
689 tions, and that the proportion of individuals migrating from a subpopulation
690 i to another subpopulation j is m_{ij} (Hedrick 2010). The values of m_{ij} then
691 can form a matrix called a backward migration matrix (Bodmer and Cavalli-
692 Sforza 1968). In this matrix, the proportion of residents (i.e. not migrants)
693 in each subpopulation i are given by the diagonal values of the matrix (i.e.
694 m_{ii}). Each row of the matrix sums to 1, because it describes the proportion

695 of migrants coming into a population i from the other j populations. For this
696 model, the amount of allele A_2 in any subpopulation i after gene flow is:

$$q'_i = \sum_{j=1}^k m_{ij} q_j \quad (1.18)$$

697 To process of allele frequency change over time can be described with
698 matrix notation, where M is the migration matrix, and Q_t is the vector of
699 allele frequencies in each population at generation t :

$$Q_{t+1} = MQ_t \quad (1.19)$$

700 The above can be generalized for any t

$$Q_t = M^t Q_0 \quad (1.20)$$

701 (Hedrick 2010)

702 In this model, as with the continent-island model previously described,
703 after a period of time, allele frequencies in the subpopulations converge and
704 approach an asymptotic value. This value can be calculated with equation
705 1.18 using a migration matrix raised to a power of t large enough that all
706 elements have reached their asymptotic values. This demonstrates the
707 homogenizing effect gene flow has on populations when it is sustained for
708 a period of time (Charlesworth and Charlesworth 2010; Hedrick 2010).

709 1.1.4.1 Gene flow - drift balance

710 Gene flow acts to homogenize populations as described above. However
711 populations are finite in size and so genetic drift will cause differences
712 between the populations through the random fixation and loss of alleles.
713 The joint effects of gene flow and drift can be examined using a simple
714 model of replicate island populations (Wright 1940). Each island has N

715 individuals and receives a proportion of migrants each generation m , from
716 a continent population.

717 When the gene flow between islands, and the population size of the
718 islands are large the allele frequencies on the islands behave as previously
719 described: they will converge to the frequencies of the continent. However
720 if population sizes are small, and the amount of gene flow is low, then
721 the allele frequencies of the islands may differ from each other (Hedrick
722 2010). So genetic drift causes allele frequencies in subpopulations to drift
723 apart, whilst gene flow acts to homogenise the allele frequencies: Take N
724 to be equal to N_e , the probability two alleles coalesce in generation $t - 1$ is
725 $1/(2N)$ and the probability that they do not is $1 - 1/(2N)$ (Hedrick 2010).
726 The expected homozygosity in generation t can be given as:

$$f_t = \frac{1}{2N} + \left(1 - \frac{1}{2N}\right) f_{t-1} \quad (1.21)$$

727 This expression can be modified by the probability that both alleles are
728 not migrants:

$$f_t = \left[\frac{1}{2N} + \left(1 - \frac{1}{2N}\right) f_{t-1} \right] (1 - m)^2 \quad (1.22)$$

729 Assuming there is an equilibrium between gene flow homogenizing
730 variation, and drift generating variation, then $f = f_t = f_{t-1}$ and $f = F_{ST}$,
731 then

$$F_{ST} = \frac{(1 - m)^2}{2N - (2N - 1)(1 - m)^2} \quad (1.23)$$

732 (Hedrick 2010)

733 F_{ST} is the fixation index, a measure of genetic differentiation over sub-
734 populations. When $m = 0$ then $F_{ST} = 1$, and when $m = 1$, then $F_{ST} = 0$. In
735 other words when levels of gene flow are high, the genetic differentiation
736 over subpopulations is low. Ignoring the powers of two, and reducing the

737 formula, F_{ST} can be approximated:

$$F_{ST} \approx \frac{1}{4Nm + 1} \quad (1.24)$$

738 Assuming k subpopulations, the differentiation between populations can
739 be given as

$$G_{ST} = \frac{1}{4Nm \left(\frac{k}{k-1} \right)^2 + 1} \quad (1.25)$$

740 (Slatkin 1995). In both equations, Nm means the absolute number of
741 migrants entering a population every generation.

742 F_{ST} for any generation t has been derived when $m = 0$

$$F_{ST(t)} = 1 - e^{t/2N} \quad (1.26)$$

743 (Wright 1943).

744 The above expression is 0, when subpopulations are not very separated
745 in early generations, and reaches a maximum of 1 as subpopulations
746 are separated by drift. The smaller the population size, the faster the
747 subpopulations diverge due to drift. The increase in F_{ST} is fastest for the
748 first $2N$ generations, after which time it approaches the maximum of 1.

749 Iterating over formula 1.22 allows examination of the rate of approach to
750 equilibrium for different values of N and m . When population size is large
751 and the amount of gene flow is large, then approach to equilibrium is fast,
752 but when populations are large and gene flow is small, then the approach
753 to equilibrium is slow (Hedrick 2010).

754 Population subdivision also affects the N_e of populations. For the island
755 model:

$$N_e = \frac{kN}{1 - F_{ST}} \quad (1.27)$$

756 If F_{ST} is low, then $N_e \approx kN$, but if gene flow is low then N_e might be
757 larger than kN (Wright 1943).

758 Wright 1940 gave an explicit method of estimating allele frequencies
759 incorporating the effects of gene flow and drift for the island model. Assum-
760 ing the frequency of A_2 in migrants (q_m) is constant, when observing a large
761 number of islands, their average allele frequency will be q_m , but depend-
762 ing on drift and gene flow, the distribution over the islands will vary. The
763 shape of the distribution depends on $4Nm q_m$ and $4Nm(1 - q_m)$. With large
764 amounts of gene flow and large population sizes, the allele frequencies over
765 the islands will not depart far from the mean (Hedrick 2010). However, with
766 lower $4Nm q_m$ and $4Nm(1 - q_m)$, and if $q_m = 0.5$, then the distribution takes
767 on a U shape: Drift plays a greater role in determining allele frequencies
768 as alleles enter the islands by gene flow, and islands become temporarily
769 fixed for either A_2 , or instead for A_1 .

770 Other models add an extra consideration by assuming different popula-
771 tions occupy positions in space, and that gene flow is restricted to certain
772 routes or directions. For example, the stepping stone model arranges popu-
773 lations in a one dimensional structure, and restricts gene flow to occurring
774 only between populations that are adjacent in that one dimensional space
775 (Hedrick 2010). The effective population size of such a linearly divided popu-
776 lation can be approximated as $N_e \approx kN$ (Maruyama 1970). If populations
777 are distributed across a landscape according to available habitat, then there
778 may be distance-dependent gene flow between the populations. In such
779 case, expected patterns of genetic variation may be similar to the stepping
780 stone models (Wright 1943). It has been suggested that the amount of
781 genetic divergence as estimated with Nm or $F_{ST}/(1 - F_{ST})$ should change
782 as an inverse linear function of geographic distance (Nm), or as a linear
783 function of geographic distance ($F_{ST}/(1 - F_{ST})$) (Rousset 1997).

784 In metapopulations (Levins 1969), the dynamics of recolonization and

785 extinction greatly influence N_e , the genetic variation present in the metapop-
 786 ulation, and the distribution of genetic variation over the subpopulations
 787 (Slatkin 1977; Hedrick and Gilpin 1997; Whitlock and Barton 1997; Nunney
 788 1999). Many parameters can influence the rate at which genetic variation is
 789 lost, for example, the source of individuals recolonizing a previously extinct
 790 path might be from a single path, or a group of individuals from all other
 791 non-extinct patches. A metapopulation with 20 patches, an infinite popu-
 792 lation size in each patch, and no gene flow except during recolonization,
 793 will have an effective size of 150 when recolonization of a patch is from a
 794 single female from another patch. This low N_e is due to the low number of
 795 founders in each recolonization (Hedrick and Gilpin 1997; Hedrick 2010).

796 1.1.4.2 Gene flow - selection balance

797 Gene flow and selection are often both important forces driving allele
 798 frequencies in a population. Both forces are diverse in their effects on allele
 799 frequencies and so the interaction of the two forces can lead to complex
 800 results (Lenormand 2002). Therefore, only a simple scenarios of selection
 801 and gene flow is introduced here.

802 Consider again the continent-island model, if the change in allele fre-
 803 quency due selection is Δq_s , and the change in allele frequency due to
 804 gene flow is Δq_m , then the change in allele frequency due to the joint effect
 805 of the two forces is $\Delta q = \Delta q_m + \Delta q_s$ (Hedrick 2010). Assuming the fitness
 806 values of A_1A_1 , A_1A_2 , and A_2A_2 are 1, $1 - s$, and $1 - 2s$ respectively, then
 807 Δq can be expressed as $\Delta q = sq^2 - (m + s)q + mq_m$ (Li 1976). When $\Delta q = 0$,
 808 there is equilibrium, and the equilibrium frequency is found by solving the
 809 quadratic equation.

$$q_e = \frac{1}{2s} \{ (m + s) \pm [(m + s)^2 - 4msq_m]^{1/2} \} \quad (1.28)$$

810 A_1 is favored if s is positive, otherwise A_2 is favored (Hedrick 2010).
811 There are three main scenarios to consider, one where gene flow is much
812 less, greater than, or equal to the absolute value of selection ($|s|$). As m
813 increases with respect to $|s|$, genetic differentiation does not occur. This is
814 intuitive, as gene flow has a homogenizing effect as previously described,
815 and with increasing m , its effects become more influential than the effects
816 of selection, and the island's equilibrium frequency approaches that of the
817 migrants coming from the continent (Li 1976).

818 Generally, the equilibrium frequency of an island depends on the se-
819 lective advantage, the level of dominance on the island, and the amount
820 of gene flow. With high amounts of gene flow, even a favorable variant
821 can be lost from an island, no matter its level of dominance. This is called
822 patch disappearance (Haldane 1948). Thus gene flow is a force which
823 limits selection and local adaptation (Lenormand 2002).

824 1.1.4.3 Importance of gene flow

825 Gene flow and genetic structuring significantly influence plant pathogen
826 and marine plankton populations. Gene flow is the force which introduces
827 new virulence alleles into a new agricultural field, far from the source of
828 original mutation. Plant pathogen populations are often made up of one or
829 a few clonal lineages which differentiate themselves from other populations
830 (in chapter 3 these are called 'races') (Koenig, Ploetz, and Kistler 1997).
831 In such populations, it may help instead to think of genotype flow rather
832 than gene flow because of the high degree of linkage. Genotype flow
833 refers to the movement of entire genotypes between distinct populations.
834 Since many plant pathogens have an asexual stage and a sexual stage,
835 both genotype flow and gene flow can occur. An existing example of gene
836 flow between plant pathogen populations is provided by Zhan, Pettway,
837 and McDonald 2003, who demonstrated that *Mycosphaerella graminicola*

838 populations shared RFLP alleles, but no two populations had completely
839 identical fingerprints, indicating that gene flow, but not genotype flow, was
840 occurring. An example of genotype flow is the global movement of a single
841 clone of *Phytophthora infestans*, out of Mexico in the 1840's as previously
842 described. Only one mating type escaped and spread globally, and as the
843 organism has two mating types, sexuality was not possible until the other
844 mating type escaped in the 1970's (Goodwin et al. 1992; Goodwin, Cohen,
845 and Fry 1994; Goodwin et al. 1995).

846 There is substantial evidence of genetic structuring in marine plankton
847 populations despite the high dispersal capacity of those organisms that
848 might usually lead one to expect high levels of gene flow (Sildever et al.
849 2016). Oceanographic features like currents and eddies will create habitat
850 heterogeneity which in turn leads to genetic population structuring (White et
851 al. 2010; Sanford and Kelly 2011; Casabianca et al. 2012), as do chemical
852 and biotic properties of the oceans such as pH levels, temperature, salinity,
853 and the presence or absence of predators and parasites (Cousyn et al.
854 2001; Decaestecker et al. 2007; Weisse et al. 2007; Yampolsky, Schaer,
855 and Ebert 2014; Defaveri and Meril 2014). All these factors may cause
856 local adaptation resulting in population structuring.

857 **1.1.5 Recombination and linkage**

858 In the theory introduced so far, it has been assumed that alleles at a locus
859 under consideration are transmitted independently of any alleles at any
860 other loci. This is called independent assortment (Hedrick 2010). It was
861 also assumed that the fitnesses of genotypes at any given locus were
862 independent of the fitnesses of other genotypes at other loci. However,
863 this simplification is not valid in the majority of cases. The transmission
864 of genetic variants does not occur independently of other genetic variants.
865 This is because of linkage between genetic variants; variants are distributed

866 across DNA molecules, and two variants situated on the same molecule
867 are said to be physically linked. The non-random association of alleles
868 is called linkage disequilibrium (LD) (Lewontin and Kojima 1960). The
869 amount of LD is generally an inverse function of the rate of recombination.
870 Where recombination is the rearrangement of genetic material, especially
871 by crossing over in chromosomes or by the artificial joining of segments of
872 DNA from different organisms

873 If one assumes a large randomly mating population has two alleles at
874 one locus A (A_1, A_2), and two alleles at a second locus B (B_1, B_2), then
875 four gametes or haplotypes are possible: A_1B_1, A_1B_2, A_2B_1 , and A_2B_2 .
876 The frequencies of these four haplotypes are denoted as x_{11}, x_{12}, x_{21} , and
877 x_{22} . The frequencies of each allele are $p_1 = x_{11} + x_{12}$, $p_2 = x_{21} + x_{22}$,
878 $q_1 = x_{11} + x_{21}$, and $q_2 = x_{12} + x_{22}$ for A_1, A_2, B_1 , and b_2 , respectively
879 (Lewontin and Kojima 1960).

880 Assuming random association between alleles in gametes, then the
881 frequency of each gamete is equal to the product of the frequencies of the
882 alleles it is made of. In other words $x_{11} = p_1q_1, x_{12} = p_1q_2, x_{21} = p_2q_1, x_{22} =$
883 p_2q_2 . However, when this assumption does not hold and there is nonrandom
884 association between alleles, the frequencies must be written as a function
885 of these expected frequencies, with some deviation D from the expectation.
886 Therefore, $x_{11} = p_1q_1 + D, x_{12} = p_1q_2 - D, x_{21} = p_2q_1 - D, x_{22} = p_2q_2 + D$.
887 D is the LD parameter and it is a measure of the deviation from random
888 association between alleles at different loci, $D = x_{11} - p_1q_1$ (Lewontin and
889 Kojima 1960). In other words it is the observed frequency of a gamete,
890 minus the expected frequency of the gamete. By substituting values p_1 and
891 q_1 , D may be written as:

$$D = x_{11}x_{22} - x_{12}x_{21} \quad (1.29)$$

892 The gametes can be categorized as coupling or repulsion gametes.
 893 Coupling gametes are those with alleles of the same subscript, and repul-
 894 sion gametes are those with alleles with different subscripts. D then is the
 895 product of the frequencies of the two coupling gametes, minus the product
 896 of the frequencies of the repulsion gametes (Hedrick 2010).

897 From these four gametes, 10 genotypes are possible. The genotypes
 898 and their expected proportions are listed in Table 1.2. These derivations
 899 make sense given that A_1B_1/A_1B_1 genotypes only produce A_1B_1 gametes,
 900 and that A_1B_1/A_1B_2 genotypes produce $1/2A_1B_1$ and $1/2A_1B_2$ gametes.
 901 Double heterozygotes produce gametes different from the parental gametes
 902 due to recombination, e.g. A_1B_2 and A_2B_1 gametes can be produced by
 903 recombination of A_1B_1/A_2B_2 individuals. The recombination rate is denoted
 904 as c in Table 1.2. c ranges from 0 where there is no recombination between
 905 loci A and B , to 0.5 or independent assortment. The frequency of each
 906 gamete in the next generation can be calculated the summing each of
 907 columns 3 to 6 in Table 1.2, the simplified way of working out such sums
 908 are given on the bottom line of the table, where D_0 is the initial amount of
 909 LD (Hedrick 2010).

910 The amount of D after one generation then is $D_1 = x'_{11}x'_{22} - x'_{12}x'_{21}$.
 911 After substitution and simplification this becomes $D_1 = (1 - c)D_0$, which is
 912 recursive and so can become

$$D_t = (1 - c)^t D_0 \quad (1.30)$$

913 with D_t meaning the amount of LD after t generations (Hedrick 2010).

914 With this formula we see that when there is no linkage ($c = 0.5$) most
 915 disequilibrium is lost within a few generations, and with lower recombination
 916 rate, linkage is tighter as recombination does not break up associations
 917 between alleles as frequently, and so LD does not decay as fast.

Table 1.2: Expected frequencies for different gametes in a two-allele, two-locus system, adapted from Hedrick 2010.

Genotypes	Frequencies	Gametes of offspring			
		A_1B_1	A_1B_2	A_2B_1	A_2B_2
A_1B_1/A_1B_1	x_{11}^2	x_{11}^2	—	—	—
A_1B_1/A_1B_2	$2x_{11}x_{12}$	$x_{11}x_{12}$	$x_{11}x_{12}$	—	—
A_1B_2/A_1B_2	x_{12}^2	—	x_{12}^2	—	—
A_1B_1/A_2B_1	$2x_{11}x_{21}$	$x_{11}x_{21}$	—	$x_{11}x_{21}$	—
A_1B_1/A_2B_2	$2x_{11}x_{22}$	$(1 - c)x_{11}x_{22}$	$cx_{11}x_{22}$	$cx_{11}x_{22}$	$(1 - c)x_{11}x_{22}$
A_1B_2/A_2B_1	$2x_{12}x_{21}$	$cx_{12}x_{21}$	$(1 - c)x_{12}x_{21}$	$(1 - c)x_{12}x_{21}$	$cx_{12}x_{21}$
A_1B_2/A_2B_2	$2x_{12}x_{22}$	—	$x_{12}x_{22}$	—	$x_{12}x_{22}$
A_2B_1/A_2B_1	x_{21}^2	—	—	x_{21}^2	—
A_2B_1/A_2B_2	$2x_{21}x_{22}$	—	—	$x_{21}x_{22}$	$x_{21}x_{22}$
A_2B_2/A_2B_2	x_{22}^2	—	—	—	x_{22}^2
1		$x'_{11} = x_{11} - cD_0$	$x'_{12} = x_{12} + cD_0$	$x'_{21} = x_{21} + cD_0$	$x'_{22} = x_{22} - cD_0$

918 To determine how long it will take for an initial amount of LD D_0 to decay
 919 to a given amount of LD D_t the equation 1.30 can be solved to give:

$$t = \frac{\ln(D_t/D_0)}{\ln(1 - c)} \quad (1.31)$$

920 (Hedrick 2010).

921 The measure of LD described is not the only one proposed (Hedrick
 922 1987; Lewontin 1988; Devlin and Risch 1995). To examine the extent of
 923 linkage equilibrium over chromosomes, the r^2 and D' are often used and the
 924 extent of LD measured varies with the estimated amount of recombination
 925 over chromosomes (Dawson et al. 2002).

926 The rate of recombination c is estimated as the proportion of recombi-
 927 nant gametes produced from a parent with a known gamete constitution
 928 (Hedrick 2010). The amount of recombination can vary because of a few

929 factors. Recombination can vary between the sexes, on different chromo-
930 somes, and between different regions on the chromosomes. Regions of
931 higher or lower levels of recombination than are expected are termed hot
932 spots and cold spots (Arnheim, Calabrese, and Nordborg 2003; Kauppi,
933 Jeffreys, and Keeney 2004). Patterns of LD can be used to try to putatively
934 identify such hot and cold recombination regions and estimate rates of
935 recombination (Stumpf and McVean 2003; Ptak, Voelpel, and Przeworski
936 2004; Auton and McVean 2007), and many other methods of recombination
937 detection in DNA sequences exist. In chapter 2 more methods for detecting
938 recombination are discussed along with presentation of the HybridCheck
939 software.

940 LD can be generated by multilocus selection. For example, tightly linked
941 members of a multigene family or supergene (Darlington and Mather 1950)
942 may be under selection that generates linkage disequilibrium as each gene
943 of the family is related in its adaptive function. Multigene family members
944 are created by serial gene duplication, followed by divergence through
945 mutation, drift, and differential selection. Therefore, they have historical
946 association, but interacting effects between them may cause selection to
947 maintain their association, keeping them in disequilibrium. The MHC of
948 vertebrates has properties of both supergenes and multigene families and is
949 in linkage disequilibrium (Edwards and Hedrick 1998; Beck and Trowsdale
950 2000).

951 LD can be influenced by genetic drift (Hill and Robertson 1968; Ohta and
952 Kimura 1969). The effects of drift on LD can be considered by imagining
953 the two-loci two-state model as four alleles at one locus. Drift will alter the
954 frequency of the gametes from generation to generation similar to that of a
955 single loci model. Thus, drift in small populations can lead to nonrandom
956 associations between alleles at different loci (Hedrick 2010). Recombination
957 reduces the effect of drift, reconstituting some gametes. The expected value

958 of the LD measure r^2 for a given effective population size N_e and a given
959 rate of recombination between two loci c , can be expressed as:

$$E(r^2) \approx \frac{1}{1 + 4N_e c} \quad (1.32)$$

960 (Hill and Robertson 1968; Ohta and Kimura 1969).

961 With large $N_e c$, $E(r^2)$ moves towards 0, with smaller $N_e c$ $E(r^2)$ ap-
962 proaches 1. Just as with the single locus model, founder events and
963 population bottlenecks can also influence LD. If N_e was small at some point
964 in the past, the LD caused may still be present if the LD has not decayed
965 (Hedrick 2010). With large $N_e c$, equation 1.32 is approximately

$$E(r^2) = \frac{1}{\rho} \quad (1.33)$$

966 where ρ is $4N_e c$ or the population recombination rate. This is analogous
967 to the population mutation rate $\theta = 4N_e \mu$ (Wall 2000; Stumpf and McVean
968 2003; Padhukasahasram et al. 2006), and the expected amount of LD
969 decreases as ρ increases (assuming that drift is the only thing affecting LD)
970 (Pritchard and Przeworski 2001; Hedrick 2010).

971 Mutations may also generate low levels of LD, however recurrent mu-
972 tation is unlikely to cause higher LD because as they are unlikely to occur
973 associated with the same allele repeatedly, and any buildup of LD through
974 mutation would occur more slowly than the process of recombination reduc-
975 ing LD (Hedrick 2010). However, mutation coupled with recombination and
976 gene flow are the source of new haplotypes in populations. New genetic
977 variants can increase in frequency by selection and drift, and hence all
978 these factors in concert may create additional LD (Hedrick 2010). Mutations
979 may also break up LD if the mutation rate is high enough. Assuming an
980 allele A_1 which mutates to a disease allele A_2 , creating a new gamete
981 $A_2 B_1$, if mutations from B_1 to any other B allele occur at rate μ , assuming

982 no recombination, the association between a disease allele A_2 and B_1 is
983 broken down. This effect has been found to be especially significant for
984 microsatellite loci, which are characterized by a high mutation rate relative
985 to SNP and indel mutations (Payseur, Place, and Weber 2008).

986 Gene conversion can also affect LD, but typically only affects shorter
987 DNA segments. Assume there is gene conversion around a gene B in
988 an $A_1B_1C_1/A_1B_2C_1$ individual, gene conversion could result in a $A_1B_2C_2$
989 gamete. B_1 has been converted to B_2 . This would decrease LD between
990 A and B , and B and C . However, it would not affect LD between A and
991 C . Many close sites do not have complete association, suggesting that
992 reduction in LD is occurring through gene conversion (Ardlie et al. 2001).
993 Note however that consecutive mutations can also explain the incomplete
994 association between linked sites. For example, consider three haplotypes
995 $A_1B_1C_2$, $A_1B_2C_1$, and $A_1B_2C_2$ in a 100bp fragment in a population sample.
996 This observation is consistent with recombination (between the 1st and
997 2nd haplotype, with breakpoint between B and C , creating the 3rd haplo-
998 type). It is also consistent with gene conversion (e.g. a C_1 in an ancestral
999 2nd haplotype might have been converted by the C_2 of the 1st haplotype,
1000 thereby creating a novel 3rd haplotype). Finally, it is also consistent with
1001 mutation ($B_2 \rightarrow B_1$) in the ancestral 3rd haplotype ($A_1B_2C_2$), creating the
1002 1st haplotype, and a second mutation ($C_2 \rightarrow C_1$) in another copy of the
1003 ancestral $A_1B_2C_2$ haplotype (before or after the first mutation at any point in
1004 time) resulting in the 2nd haplotype. In other words, and in contrast to Ardlie
1005 et al. 2001, the observation that many close sites do not have complete
1006 association should not be taken as evidence for gene conversion because
1007 other evolutionary forces can explain this observation more plausibly.

1008 Gene flow can also affect LD. The amount of disequilibrium when two

1009 populations are mixed to produce a third can be expressed as

$$D = m_x m_y (p_{1 \cdot x} - p_{1 \cdot y})(q_{1 \cdot x} - q_{1 \cdot y}) \quad (1.34)$$

1010 where $p_{1 \cdot x}$ and $p_{1 \cdot y}$ are the frequencies of of the A_1 allele in the two pop-
1011 ulations being mixed (population x and population y), and $q_{1 \cdot x}$ and $q_{1 \cdot y}$
1012 are the frequencies of the B_1 allele in the two populations (Hedrick 2010).
1013 For LD to be generated, the frequencies of both loci must be different in
1014 the two populations. The greater the difference, and the more equal the
1015 contributions are from each population, the more LD is generated (Hedrick
1016 2010).

1017 Population subdivision reduces the rate of LD decay. The reduction
1018 in heterozygotes in subdivided populations due to the Wahlund effect
1019 (Wahlund 1928) reduces the opportunity to create recombinant gametes.
1020 If the amount of gene flow is small, then it can determine the rate of LD
1021 decay (Nei and Li 1973). The amount of linkage disequilibrium has been
1022 expressed as $D \approx m/c$ (Barton et al. 2007) i.e. it is a balance between the
1023 rate of gene flow creating LD, and the rate of recombination reducing LD.
1024 Since many factors including selection, drift, gene flow and mutation affect
1025 LD, it can be difficult therefore to attribute a cause of LD without historical
1026 knowledge or data.

1027 Since alleles are linked and selection occurs at one or more loci we
1028 say that alleles have a genetic background (Hedrick 2010). Multilocus
1029 phenomenon may explain some observations encountered in evolutionary
1030 genetics. Apparent heterozygous advantage at a given marker locus may
1031 actually be caused by association of alleles at a linked locus to the alleles
1032 at the marker locus (Ohta 1971). For example, Oosterhout 2009 proposed
1033 that the genetic variation at the MHC may be maintained by a linkage
1034 of the genetic load (or sheltered load) present at the peri-MHC region.

1035 Recessive deleterious mutations associated with a given haplotype prevent
1036 the fixation of that haplotype in the population because these mutations
1037 would become expressed in homozygous state, reducing the fitness of that
1038 individual. In other words, an MHC haplotype is self incompatible because it
1039 expresses its genetic load in homozygous state. Assuming that each MHC
1040 haplotype has its own sheltered load of recessive deleterious mutations,
1041 this prevents their fixation in the population, and results in a balanced
1042 polymorphism (Oosterhout 2009). Recombination between MHC alleles
1043 is further reduced by negative epistasis, with selection operating against
1044 recombination because the recombinant haplotypes are incompatible with
1045 both parental (non-recombinant) haplotypes.

1046 Furthermore, changes in allele frequencies might be the result of se-
1047 lection acting on alleles at an associated locus to one being observed.
1048 This can result in genetic hitchhiking, selective sweeps or background
1049 selection (Charlesworth and Charlesworth 2010). Genetic hitchhiking, pre-
1050 viously described as the mechanism by which hypermutator alleles can
1051 be indirectly selected for in clonal populations (section 1.1.3), is possible
1052 because of linkage. Neutral alleles can increase in frequency because
1053 of their association with a selected allele. The magnitude of hitchhiking
1054 depends on the extent of linkage, inbreeding, and the initial amount of LD
1055 (Thomson 1977; Hedrick 1980; Kaplan, Hudson, and Langley 1989). If
1056 there is no initial statistical association between the neutral and selected
1057 allele, there can be no hitchhiking, even if recombination rates are low. To
1058 fully understand the effect of hitchhiking, the rate of change in frequency of
1059 the positively selected allele must be known (Hedrick 2010). For example
1060 for a new advantageous recessive allele, initial increase in frequency due
1061 to selection will be low (see section 1.1.1), providing time for recombination
1062 to reduce initial LD, and thus reducing the amount of expected hitchhiking
1063 of neutral alleles. Hitchhiking can even create LD between two neutral loci

1064 if they are associated with a third selected locus (Thomson 1977; Hedrick
1065 1980). One of the most important effects of hitchhiking is the reduction in
1066 heterozygosity of neutral or nearly neutral variation in areas of low recomb-
1067 ination (Maynard-Smith and Haigh 1974). This is called a selective sweep
1068 and leaves a characteristic signature in genome sequences, which can be
1069 detected to provide evidence of recent selection (Hedrick 2010).

1070 The projects presented in this thesis are concerned with how recom-
1071 bination has influenced the adaptive evolution of the two species studied.
1072 Specific aspects of recombination, sex and linkage relevant to each project
1073 are introduced in detail in subsequent chapters. In the introduction to
1074 chapter 4 the advantages and disadvantages of recombination and sex are
1075 presented, to provide context to the question of why *F. cylindrus* might have
1076 abandoned sex (as is hypothesized at the start of the study). In chapter
1077 3 the evolutionary advantages and disadvantages of introgression and
1078 hybridisation is discussed in the context of results, and there multilocus
1079 concepts are important.

1080 **1.1.6 Hybrid zones, introgression, and hybrid speciation**

1081 Gene flow and recombination can result in so called hybrid zones, a physical
1082 location where hybrid offspring of two diverged taxa occur (Hewitt 1985).
1083 A hybrid zone may form where divergence is occurring between adjacent
1084 populations of a species that was previously homogenous. Parapatric and
1085 peripatric speciation is most likely to result in hybrid zones because the
1086 divergence and speciation is driven not by geographical isolation. With
1087 parapatric speciation, changes in environmental conditions between the
1088 adjacent population can result in adaptations and reproductive isolation
1089 (Mayr 1942). Founder events and random genetic drift play an important
1090 role during peripatric speciation. Before reproductive isolation has evolved,
1091 ongoing gene flow and recombination between the two adjacent populations

1092 could result in a hybrid zone. In this case, the hybrid zone is called a primary
1093 hybrid zone. Hybrid zones may also form as a result of secondary contact
1094 between two populations of diverged taxa which were previously allopatric
1095 and had diverged as a result of geographic isolation. In the latter case,
1096 partial pre-zygotic reproductive isolation has evolved, but this is broken
1097 down, for example due to changes in environmental conditions that could
1098 hinder conspecific mate choice. It is often difficult to distinguish between
1099 primary and secondary hybrid zones (Endler 1982).

1100 Such hybrid zones have a cline in the genetic composition across the
1101 zone from one of the parental forms to the other, as novel alleles from
1102 either side (that is either parental population) flow into the hybrid zone.
1103 Such clines can either be gradual or stepped, and they can be observed by
1104 recording the frequency of diagnostic alleles for the parental populations,
1105 across the transect between the two parental populations (Hewitt 1985).
1106 When quantifying the cline in this way, the frequency of diagnostic alleles
1107 is often characterized by a sigmoid curve, and the width of the cline is
1108 dependent on the ratio of hybrid survival to rate of recombination (Hewitt
1109 1985). In addition to a cline of genetic composition, hybrid zones often
1110 exhibit a higher variability in fitness within the zone. In the middle of the
1111 cline hybridzymes may also be found. Hybridzymes are rare alleles from both
1112 the parental taxa, which reach high frequencies where hybrids are formed,
1113 due to genetic hitchhiking of those alleles with alleles that contribute to
1114 hybrid fitness (Schilthuizen, Hoekstra, and Gittenberger 1999).

1115 It is possible for alleles to flow back into the distinct parental popula-
1116 tions through introgression (subsequent backcrossing of a hybrid individual
1117 breeding with a parental individual). As a result, they appear to present a
1118 problem for the biological definition of a species if it is defined as a popula-
1119 tion of (potentially) interbreeding individuals that produce fertile offspring,
1120 however if the two parental populations remain identifiably distinct then

1121 there is no problem for the alternative concept of a species as taxa that
1122 retain their identity, despite gene flow (Mayr 1942).

1123 When introgression occurs, each generation is less able to replace
1124 itself with genetically similar individuals as a result of the influx of alleles
1125 from across the hybrid zone, and this may lead to genetic assimilation and
1126 homogenization of the two parental populations (Robbins et al. 2014). How-
1127 ever, hybridisation does not always lead to the merging and homogenizing
1128 of the two populations involved. The different evolutionary outcomes of
1129 hybridization occur through different pathways in addition to introgression,
1130 like consequences of ecology such as hybrid vigour or hybrid inferiority
1131 (Edmands 1999; Johansen-Morris and Latta 2006; Rieseberg and Carney
1132 1998).

1133 Hybrid vigour can lead to a slowing of the growth rates of the two
1134 parental populations, because of the competition with the more fit hybrids
1135 (Slattery et al. 2008). But equally, if the increased hybrid fitness only
1136 applies in the hybrid zone, then a stable situation occurs in which the
1137 two parental populations are not threatened with assimilation, and instead
1138 hybrid speciation may occur, whereby hybridisation leads to hybrids which
1139 are reproductively isolated from either of the two parental populations.
1140 Some hybrid zones can persist for thousands of years (White et al. 1966).
1141 This is possible as the hybrid zones are so called tension-zones. In tension
1142 zones, there is a balance between ongoing hybridisation, dispersal of
1143 parental forms, and natural selection against hybrids (hybrid inferiority). If
1144 those forces are in equilibrium, a stable tension zone persists (Bazykin
1145 1969). Recent studies identifying the signature of admixture across the
1146 genomes of native westslope cutthroat trout, and an invasive rainbow trout,
1147 revealed genome-wide selection against the invasive alleles, and that this
1148 was consistent across environments and populations (Kovach et al. 2016).
1149 It is important to note when considering the possible paths the evolution of a

1150 hybrid zone may take, that the different outcomes are not exclusive either/or
1151 scenarios: For example, even though a hybrid zone may be maintained
1152 by negative selection acting on hybrids, and whilst some alleles from a
1153 parental population will be prevented from flowing into the other parental
1154 population as a result of negative selection, other alleles that are neutral or
1155 positively selected for may be able to flow across the hybrid zone and into
1156 the other population (Hewitt 1985). Both of these processes are occurring
1157 at once, with the outcome varying across the genome, depending on the
1158 alleles. In this way, a hybrid zone acts as a semi-permeable barrier to
1159 the flow of alleles. Analysis of genetic and phenotypic variation across a
1160 hybrid zone of *Antirrhinum*, populations near the French-Spanish border
1161 is one such example demonstrating this (Whibley et al. 2006): The hybrid
1162 zone has a very steep cline in flower colour and morphology across the
1163 hybrid zone. After crossing plant morphs to determine the contribution of
1164 the EL, ROS, and SULF alleles to magenta and yellow flower colouration,
1165 they used image analysis to score the levels of pigment in the plant and a
1166 principal component analysis on pixel scores together allowed the creation
1167 of a 3D genotypic space or landscape controlling flower colour (Whibley
1168 et al. 2006). Sequencing of natural samples across the hybrid zone allowed
1169 them to identify three main haplogroups. One haplogroup was specific to
1170 the yellow morph, and the other two were found only in magenta morphs,
1171 the flower colour cline coincided with a cline in the frequency of these
1172 haplotypes. The researchers then sequenced loci not involved in flower
1173 colour determination, the PAL and DICH loci, which are linked to the ROS
1174 colour determination locus. They sequenced PAL and DICH loci from 18
1175 individuals either side of the hybrid zone. They found PAL alleles fell into
1176 two distinct haplogroups, whilst DICH had no haplogroup structure (Whibley
1177 et al. 2006). Sequencing PAL and DICH alleles from individuals across the
1178 hybrid zone revealed no cline in the frequencies of these alleles, showing

1179 they are subject to different evolutionary forces. These alleles also had
1180 no correlation with flower colour. They concluded the distribution of the
1181 two alleles reflects historical gene flow, thus the hybrid zone is a barrier
1182 to alleles determining flower colour, as F2 hybrids are less fit according
1183 to their 3D fitness landscape, but other alleles are able to pass through
1184 (Whibley et al. 2006).

1185 Hybridization and introgression, is thought to occur in roughly 10% of
1186 animal species and 25% of plant species (Mallet 2005). Hybridization may
1187 lead to hybrid speciation, which is where new hybrid lineages become
1188 reproductively isolated from parental populations, and so are considered
1189 separate species. Genomic studies have allowed determination of the
1190 sizes of parental chromosomal blocks in introgressed populations and
1191 hybrid species (Buerkle and Rieseberg 2008; Morrell et al. 2005), as they
1192 allow observation of associations among alleles of one species in the
1193 genetic background of another, indicating recent introgression. Genome-
1194 wide studies of introgression and hybridisation have also supported the
1195 conclusions supported by the work of Whibley et al. 2006, that there is
1196 variation in the amount of introgression across genomes, and so some
1197 regions of the genome are more permeable to foreign alleles than others
1198 (Martinsen et al. 2001; Macholán et al. 2007; Scotti-Saintagne et al. 2004;
1199 Turner, Hahn, and Nuzhdin 2005; Yatabe et al. 2007).

1200 Substantial changes can occur to a genome immediately after hybridisa-
1201 tion, such as gene loss or silencing, changes in expression of some genes
1202 (Adams and Wendel 2005). Analysis of three synthetic sunflower hybrids
1203 and three natural sunflower hybrid species has shown large karyotypic
1204 changes can occur over a handful of hybrid generations (Karrenberg, Lexer,
1205 and Rieseberg 2007; Lai et al. 2005). The natural hybrid species also ex-
1206 hibit increased genome sizes of up to nearly 50% compared to the parental
1207 species (Baack, Whitney, and Rieseberg 2005). All species showed similar

1208 increases in genome size because of the proliferation of retrotransposons
1209 (Ungerer, Strakosh, and Zhen 2006).

1210 The evolutionary consequences of hybridisation are complex. F1 hybrids
1211 are often larger and more fit than their parents due to the effects of heterosis
1212 (Lippman and Zamir 2007), due to either overdominance or the reciprocal
1213 complementation of deleterious alleles (Clark Cockerham and Zeng 1996),
1214 this explains the establishment of hybrids but does not determine the longer
1215 term evolutionary success or failure of hybrids, which is more complex
1216 and is discussed in more detail in chapter 3. In chapter 3, processes of
1217 hybridisation and introgression, and the evolutionary outcomes of such
1218 processes are discussed in more detail, and in the context of the work
1219 presented in that chapter, which focuses on the role of such processes in
1220 the adaptive evolution of a plant pathogen species as it adapted to many
1221 hosts.

1222 **1.2 The role of Bioinformatics in population ge-** 1223 **netics**

1224 Deoxyribonucleic acid was demonstrated as the genetic material by Oswald
1225 Theodore Avery in 1944 (Russell 1988). Watson and Crick demonstrated its
1226 double helix structure composed of four nucleotide bases in 1953 (Watson
1227 and Crick 1953). This led to the central dogma of molecular biology. In most
1228 cases, genomic DNA defined the species and individuals, which makes
1229 the DNA sequence fundamental to the research on the structures and
1230 functions of cells. Sequencing of genomes then is now an essential task to
1231 complete, yielding essential data biologists need to understand biology and
1232 evolution of organisms. The automated Sanger method was considered a
1233 first-generation sequencing technology (Sanger and Coulson 1975; Sanger,
1234 Nicklen, and Coulson 1977), and since then newer methods have been

1235 developed making sequencing cheaper and increasingly high throughput,
1236 these are referred to as next-generation sequencing (NGS) technologies
1237 (Goodwin, McPherson, and McCombie 2016).

1238 With the development of NGS technology, algorithms and tools for
1239 bioinformatics and evolutionary study have developed rapidly. Here, I
1240 present a brief overview of the principles of several key bioinformatics tasks
1241 that population genetic studies with NGS data require. The processes
1242 below assume quality control of NGS reads is completed.

1243 **1.2.1 Sequence Alignment**

1244 An alignment of two sequences aims to discover or highlight how similar
1245 the two sequences are. The concept of alignments is a natural one in
1246 settings where one sequence, changes over time into a second sequence,
1247 through a series of simple operations (called edit operations) like insertions
1248 of characters, deletions of characters, and a substitution of one character for
1249 another (Mäkinen et al. 2015). It is unsurprising therefore, that alignments
1250 are a common first step in many evolutionary analyses. An alignment of the
1251 characters in two sequences, which have stayed the same over time, could
1252 be defined as the list of pairs of positions (i, j) such that the i th position in
1253 the first sequence is considered a match to the j th positions in the second
1254 sequence (Mäkinen et al. 2015).

1255 In a practical setting, the two sequences (A & B) are typically short ho-
1256 mologous regions of the genomes of two different individuals, or species/taxa,
1257 and are considered to have evolved through a series of changes (edit op-
1258 erations), from some unobserved common ancestor (Lemey, Salemi, and
1259 Vamdamme 2009). DNA sequence alignment algorithms typically require a
1260 scoring matrix which which they score potential alignments. These matrices
1261 typically define scores for aligning any two characters in two sequences,

1262 and have some basis in biologically reality. For example, the BLOSUM scor-
1263 ing matrix was derived from data of conserved regions of protein families
1264 (Lemey, Salemi, and Vandamme 2009). The score of any given pairwise
1265 alignment is the sum of the scores that were assigned by the scoring matrix
1266 for each position of the alignment.

1267 A local alignment algorithm attempts to find the best alignments for sub-
1268 sequences of a query sequence with a reference sequence (Lemey, Salemi,
1269 and Vandamme 2009). Whereas global alignment algorithms attempt to
1270 find the best end to end alignment between a query sequence and a
1271 reference sequence. Traditionally, pairwise sequence alignments were
1272 computed using dynamic programming algorithms such as the Needleman-
1273 Wunsch (global sequence alignment) (Needleman and Wunsch 1970),
1274 and the Smith-Waterman (local sequence alignment) algorithms (Smith
1275 and Waterman 1981), but efficient and accurate techniques for sequence
1276 alignment is an active area of research, and so many advances, and
1277 different techniques and software packages have been developed. Multiple
1278 alignment is the generalisation of pairwise sequence alignment to more
1279 than two sequences, this is a hard problem which becomes computationally
1280 unfeasible for many sequences without use of heuristics, such as the
1281 progressive alignment method, which first constructs a guide tree (Löytynoja
1282 and Goldman 2005).

1283 Sequence alignments can be used to align multiple gene or protein
1284 sequences together, align reads from high throughput sequencing platforms
1285 to a reference genome assembly (Li and Durbin 2009), or to align different
1286 genome assemblies together (Paten, Earl, and Nguyen 2011). In all cases
1287 these alignments may be used to run variant calling algorithms to infer the
1288 presence of mutations and structural alterations that are present in the
1289 genomes of different taxa, individuals, or populations, and can be used to
1290 genotype individuals, and compute population genetics and evolutionary

1291 analyses.

1292 **1.2.2 Variant Calling**

1293 Variant calling yields genotype data which may then be used in population
1294 genetics study. Identification of SNPs (sometimes called Single Nucleotide
1295 Polymorphisms or simply mutations) can be done with a read pileup output
1296 after aligning reads to a reference genome (Li et al. 2009; Li 2011). If a
1297 position j in the reference genome is covered by n reads, and of those
1298 reads, p per cent of them indicate that position j is an A, and the rest
1299 indicate that position j is a G, then it is possible to reason whether this is
1300 because the sample that was sequenced is polymorphic, or because of
1301 an alignment error or sequencing error (Mäkinen et al. 2015). Such errors
1302 are easy to identify, as they are independent events, and as such exist in
1303 a very low frequency, because the probability of observing many errors
1304 in the same location decreases exponentially. Therefore, so long as the
1305 sequencing is done to a sufficient depth of coverage, one can identify the
1306 polymorphic positions in a genome and rule out the errors with reasonable
1307 accuracy (Mäkinen et al. 2015).

1308 Larger variants can also be detected from the read pileup. If there is a
1309 deletion in the genome of the sample from which the reads were sequenced,
1310 then if it is larger than the error threshold in the alignment, then there should
1311 be regions of the read pileup where the reference is uncovered by reads
1312 (Mäkinen et al. 2015). The region should have the same length as the
1313 deletion. If there is an insertion in the genome of the sample from which
1314 the reads were sequenced, then if it is longer than the error threshold of
1315 the alignment, then in the pileup there would be a series of consecutive
1316 positions $(j, j + 1)$ for which no read covers both j and $j + 1$ (Mäkinen et al.
1317 2015). This is a simplistic approach to indel detection because in reality
1318 software implementations and algorithms also take into account errors,

1319 noise, base call qualities, and can have additional complexities such as
1320 utilizing data from many samples, and from linked sites. (Li 2011; Nielsen
1321 et al. 2011; Mielczarek and Szyda 2016).

1322 Another approach to indel detection is to take advantage of sequencing
1323 technology platforms, which produce paired-end, or mate pair reads. Se-
1324 quencers can produce pairs of reads for each DNA molecule, one begins
1325 from one end of the molecule, and the other begins from the other end, and
1326 both extend towards the middle of the molecule. When paired-end read
1327 pairs are aligned to a reference genome, they have an expected distance k
1328 between them, this expected distance is known in advance according to
1329 the protocol used to prepare the DNA library for sequencing (Mäkinen et al.
1330 2015). Its possible to compute the actual distance for each paired-end read
1331 pair, and then compute the mean and variance of those distances. Once
1332 the mean distance k' and variance is known, each paired-end read pair
1333 can be tested to see if its distance is significantly different to the average
1334 distance. If the distance is significantly different an indel is inferred between
1335 those reads with length of $k - k'$ (Mäkinen et al. 2015).

1336 1.2.3 Haplotype phasing

1337 Genotypes are the unordered combination of alleles at each site of an
1338 organisms genome. The haplotype are the sequences of alleles that have
1339 been inherited together from one parent. For example, diploids possess
1340 two copies of each chromosome, therefore, in addition to being interested
1341 in which variants they possess (the genotype), one is also interested to
1342 know to which of a diploids two haplotypes each variant belongs is the
1343 variant in the organisms maternal copy of a DNA molecule, or is it in the
1344 paternal copy? The process of identifying all the variants which are situated
1345 along the same haplotype of an organism is called haplotype phasing. In
1346 an individual, variants which are clearly homozygous may be assigned to

1347 both haplotypes very simply as both haplotypes must possess them.

1348 Given that when there are N heterozygous sites in a sequenced DNA
1349 molecule, there are a total of $2^N - 1$ possible haplotypes, that could result
1350 in those haplotypes (Mäkinen et al. 2015). Haplotype phasing was known
1351 to be a hard problem even before the development of high throughput
1352 sequencing technology. However, advances have been made and several
1353 software packages now exist to perform this task. The most accurate and
1354 widely used methods employ Hidden Markov Models to infer haplotypes
1355 Mäkinen et al. 2015. For some time, a software implementation called
1356 PHASE was considered the superior method. PHASE took ideas from
1357 coalescent theory about the joint distribution of haplotypes (Marchini et al.
1358 2007; Marchini and Howie 2010; Howie, Marchini, and Stephens 2011).
1359 PHASE was limited by its speed however and since the development
1360 of PHASE other methods implemented in packages like IMPUTE2 and
1361 SHAPEIT1 & 2 have made improvements to the efficiency and accuracy of
1362 haplotype inference algorithms (Stephens and Donnelly 2003; Delaneau,
1363 Marchini, and Zagury 2012; Delaneau et al. 2013; Delaneau, Zagury, and
1364 Marchini 2013; O'Connell et al. 2014).

1365 The flow of aligning high throughput sequencing reads to a reference,
1366 running variant calling and possibly haplotype inference, followed by down-
1367 stream population genetic analysis on the genotype or haplotype data, is
1368 now a standard work-flow. The choice of which software packages and
1369 algorithms should be used for each task can be a subjective decision which
1370 should aim to follow best-practice for each case in question. For example,
1371 the best algorithm to use on human data, may not be the best one to use
1372 on an organism like wheat which has a radically different genome.

1373 CHAPTER 2

1374 HybridCheck

1375 This chapter is based on the published scientific paper:

1376 *Ward, B. J., & van Oosterhout, C. (2016). Hybridcheck: Software for the*
1377 *rapid detection, visualization and dating of recombinant regions in genome*
1378 *sequence data. Molecular Ecology Resources, 16(2), 534-539.*

1379 The project and items of work were initially set out by my supervisor,
1380 but the work I present in this chapter is entirely my own work. I drafted the
1381 pseudo-code for the project, improved the dating algorithm presented in
1382 the chapter from it's original inefficient design, wrote all the software code,
1383 documented the package, and conducted all simulations used to test the
1384 software package, and created a website, github repository, and a web-app
1385 which provides an interface for the package.

2.1 Introduction

Recombination is one of the five evolutionary forces and is important for the formation of novel genotypes, haplotypes and alleles, thereby playing a key role in adaptive evolution (Grauer and Li 2000). Recombination is also crucial for separating deleterious mutations from their genomic background, and in combination with purifying selection it helps to curtail the mutational load (Lynch and Gabriel 1990). Recombination plays a fundamental role in the repair of damaged DNA, when homologous recombination replaces a damaged DNA strand with its intact counterpart. In all likelihood, it was this function of recombination that was important in early prokaryotic life and evolution (Cavalier-Smith 2002). With respect to adaptive evolution, however, the principal consequence of recombination is that it generates novel combinations of nucleotides, which in turns allows for selection to act a much finer scale, i.e. at the level of nucleotides rather than the entire genome. Given its fundamental importance in the biology, various mechanisms have evolved that facilitate recombination; with some depending on sexual reproduction whereas others also occur in asexually reproducing taxa. As evolutionary biologists/molecular ecologists studying gene and genome sequences, it is important to understand how the various mechanisms can result in recombination.

Homologous recombination is a process that occurs in both eukaryotes and prokaryotes, and it is an essential process through which single strand and double strand breaks, as well as base mismatches in DNA molecules are repaired. With homologous recombination, there is an equal exchange of homologous DNA sequences between the two chromatids (Lemey, Salemi, and Vandamme 2009). In eukaryotes, this can occur through Double Strand Break Repair (DSBR) and Synthesis Dependent Strand Annealing (SDSA) (McMahill, Sham, and Bishop 2007; Sung and

1414 Klein 2006). In prokaryotes, the RecBCD pathway, and the RecF pathways
1415 are the primary mechanisms (Madigan et al. 2012; Smith 2012). Although
1416 these pathways differ mechanistically, they all result in the invasion of donor
1417 DNA into a recipient DNA molecule through the formation of Holliday junc-
1418 tions, branch migration, ligation, and the repair of the DNA strands (Alberts,
1419 Johnson, and Lewis 2002).

1420 The precise outcome of recombination and its effect on the donor and
1421 recipient DNA molecule depends on how the Holliday junctions are cut
1422 and resolved (Mimitou and Symington 2009). Crossing-over or reciprocal
1423 homologous recombination occurs when there is an equal exchange of
1424 sequence variation between the two homologous chromosomes (Grauer
1425 and Li 2000). Gene conversion is a type of non-reciprocal homologous
1426 recombination in which there is an unequal exchange of one sequence (the
1427 donor) to another (the recipient), such that the donor sequence replaces
1428 the recipient DNA (Grauer and Li 2000). Whereas crossing-over does not
1429 affect nucleotide variation, gene conversion tends to reduce nucleotide vari-
1430 ation by making the donor and recipient sequence identical to one another.
1431 However, even though gene conversion tends to homogenise nucleotide
1432 variation, this process too can increase haplotype and genotype variation in
1433 the population, just like crossing-over (Spurgin et al. 2011). Both reciprocal
1434 and non-reciprocal recombination can occur between non-homologous
1435 sequences (Lemey, Salemi, and Vamdamme 2009). In addition, recombi-
1436 nation can occur when distinct species or biotypes hybridise, in which case
1437 it is referred to as genetic introgression (McMullan et al. 2015). Genetic
1438 exchange between even more distantly related taxa can result in horizontal
1439 gene transfer (Eisen 2000; Ochman, Lawrence, and Groisman 2000). This
1440 too is considered a form of recombination, which occurs after gene flow
1441 between distinct taxa, and bacterial geneticists most commonly use the
1442 term 'horizontal gene transfer'.

1443 Recombination can complicate evolutionary genetic, phylogenetic and
1444 phylogenomic analyses because neighbouring nucleotides within a single
1445 genome can differ markedly in their ancestry and coalescence. In the ab-
1446 sence of recombination, the ancestry of a multiple alignment of homologous
1447 sequences can be represented by a single gene phylogeny. However, after
1448 a single recombination event, the sequences could have a different phylo-
1449 genetic history and thus different phylogenies either side of the breakpoint
1450 (Lemey et al. 2009). Each recombinant region between two breakpoints
1451 could have a distinct ancestry and be represented by a different phylogeny.
1452 With high recombination rates, the history of a set of sequences becomes
1453 increasingly complex as different portions of the genome are shuffled, re-
1454 sulting in overlapping regions with distinct coalescence (Jouet, McMullan,
1455 and Oosterhout 2015). If recombination occurs in a single panmictic pop-
1456 ulation, however, there will be relatively little variation in the ancestry of
1457 recombinant regions because all sequences coalesce relatively recently.
1458 On the other hand, recombination in structured populations (e.g. between
1459 distinct biotypes, strains or races) may result in the genetic introgression
1460 of diverged donor sequences, and this can lead to a mosaic-like genome
1461 structure (McMullan et al. 2015). In such cases, it is inappropriate to
1462 force a single phylogenetic tree onto a mosaic-like sequence, and it has
1463 been shown that this can significantly bias estimates of coalescent times
1464 (Jouet, McMullan, and Oosterhout 2015). Not only phylogenetic analyses
1465 are hindered by recombination, but also population genetic statistics can
1466 become biased if recombination is not accounted for, for example resulting
1467 in an upwards biased estimate of theta (and hence the effective population
1468 size) (McVean, Awadalla, and Fearnhead 2002; Watterson 1975), and the
1469 erroneous identification of positive selection (Shriner et al. 2003).

1470 Given that recombination can potentially affect population genetic, evo-
1471 lutionary genetic and phylogenetic analyses, it is important to examine

1472 whether recombination has left a signature in the sequence data. There are
1473 probably three questions one might address when analysing recombination
1474 in genome sequence data:

- 1475 1. Is there evidence of recombination?
- 1476 2. Where are the breakpoints / regions of recombination located in the
1477 sequence?
- 1478 3. What is the rate of recombination scaled relative to the mutation rate
1479 or theta?

1480 To detect the evidence for recombination, graphical exploratory tools can
1481 be used such as Splitstree (Huson and Bryant 2006), which visualises the
1482 impact of recombination on the phylogenetic relationship between alleles
1483 or sequences. However, to formally test the evidence of recombination,
1484 statistical tests need to be used, and many algorithms have been devel-
1485 oped for this purpose (Lemey, Salemi, and Vandamme 2009; Lemey et
1486 al. 2009). The general rationale of these tests is that recombination can
1487 insert novel nucleotides into a sequence alignment, making it appear that
1488 these polymorphisms have arisen there by mutation. A single nucleotide
1489 polymorphic (SNP) that is shared between two sequences, but which is
1490 not shared with their common ancestor is called a homoplasy. Such ho-
1491 moplasyes are explained either by recombination or convergent evolution
1492 (Maynard Smith and Smith 1998). Statistical methods for detecting recomb-
1493 ination are based on detecting phylogenetic incompatibilities that result from
1494 homoplasyes (Bruen, Philippe, and Bryant 2006; Posada, Crandall, and
1495 Holmes 2002), or by finding clusters of identical substitutions in sequences
1496 (Posada, Crandall, and Holmes 2002). Measures that are computed by
1497 such methods, such as for example the homoplasy test (Maynard Smith
1498 and Smith 1998), the informative sites test (Worobey 2001), the refined

1499 incompatibility score (Bruen, Philippe, and Bryant 2006), and the ABBA
1500 BABA test (Martin, Davey, and Jiggins 2014; Green et al. 2010) can be
1501 used to evaluate whether recombination has taken place. For example,
1502 ABBA BABA tests classify homoplasious SNPs as having one of two possi-
1503 ble parsimonious ancestries, and they calculate the Pattersons D statistic
1504 that is based on the ratio of both types of ancestries. In case there is a
1505 significant excess of one type of ancestry over the other, this is considered
1506 evidence of recombination.

1507 Once it has been established that recombination is affecting a nucleotide
1508 sequence, one can employ methods to identify where in the genome recom-
1509 bination has taken place. Those methods generally implement a scanning
1510 sliding window, and they calculate for each window the distribution of
1511 nucleotide substitutions or the genetic distance, or they assess the phyloge-
1512 netic relationships between sequences at the window (Lemey et al. 2009;
1513 Posada, Crandall, and Holmes 2002). The former two methods typically
1514 attempt to find inversions or sudden changes in substitution pattern or
1515 distance values across the windows, and they do not rely on a phylogeny.
1516 Phylogenetic methods, on the other hand, infer recombination by detecting
1517 changes in the topologies, i.e. the shape of the tree. If adjacent sections
1518 of DNA sequence are phylogenetically incongruent, this is evidence for a
1519 recombination event or breakpoint (Lemey, Salemi, and Vamdamme 2009).
1520 Methods that rely on sliding windows tend to be hampered by an increased
1521 false positive rate (Type I error rate) due to multiple testing (Lemey, Salemi,
1522 and Vamdamme 2009). Bayesian approaches (Paraskevis et al. 2005) have
1523 been developed to avoid such sequential testing problems, and in addition,
1524 they can identify breakpoint positions and the parental (donor) sequences
1525 (Suchard et al. 2002).

1526 One may also want to quantify the rate of recombination, either as a
1527 relative rate compared to the mutation rate, or as a measure of the number

1528 of bases or recombinant regions in a DNA sequence. Measures that as-
1529 sess the evidence of recombination like the homoplasmy test or the refined
1530 incompatibility score (Bruen, Philippe, and Bryant 2006; Maynard Smith
1531 and Smith 1998) can also be used to estimate the number of recombination
1532 events. For example, the refined incompatibility score for two sites in a
1533 sample can be interpreted as either the minimum number of convergent
1534 mutations, or the minimum number of recombination events that have oc-
1535 curred between a given pair of sequences (Bruen, Philippe, and Bryant
1536 2006). The homoplasmy test written by Maynard Smith and Smith 1998
1537 calculates whether there is a statistically significant excess of homoplasies
1538 derived from the dataset, compared to the number of homoplasies that
1539 would be expected by mutation, without the occurrence of any recombina-
1540 tion. Essentially then, simple measures and calculations of recombination
1541 rate estimation are based on trying to count the number of recombination
1542 events that have occurred during the evolutionary history of the collected
1543 sample (Stumpf and McVean 2003).

1544 However, given that these measures do not take into account the time
1545 to the most recent common ancestor of the sample, they simply count the
1546 number of recombination events rather than estimating the recombination
1547 rate (Posada, Crandall, and Holmes 2002). In addition, recombination
1548 events do not necessarily leave a detectable trace in the DNA sequences
1549 (Lemey, Salemi, and Vandamme 2009). To overcome this limitation, recom-
1550 bination can be modeled explicitly using coalescent approaches (Stumpf
1551 and McVean 2003). Using the coalescent as a framework, it is possible to
1552 estimate the population recombination rate ($\rho = 4Ner$) in software such as
1553 LAMARK (Hudson and Kaplan 1988; Hudson 2001; Kuhner 2006). This
1554 value is comparable to the population mutation parameter theta ($\Theta = 4Ne\mu$).
1555 Calculating ρ and Θ allows one to calculate the effect of recombination on
1556 nucleotide polymorphisms relative that of mutation (ρ/Θ).

1557 Having identified a recombination region or block between a recomb-
1558 nant sequence and its parental (donor) sequence, it is possible to estimate
1559 when recombination did occur. This is can be done by calculating a diver-
1560 gence time estimate of the block in the recombinant and parental (donor)
1561 sequence. The simplest estimates of divergence time assume a molecular
1562 clock (Li 2008; Metzgar, Scripps, and Jolla 2007), i.e. a mutation rate that
1563 is constant through time and across lineages. The nucleotide divergence
1564 between the two sequences is equivalent to $2\mu t$, in which μ is the base
1565 mutation rate and t the number of generations that have elapsed since
1566 divergence. Sequence evolution may deviate from a molecular clock, and
1567 hence, methods have been developed that can take into account variation
1568 in mutation rates between taxa, genes and evolutionary time (Brown and
1569 Yang 2011; Drummond et al. 2012; Drummond and Suchard 2010; Thorne,
1570 Kishino, and Painter 1998). The popular software BEAST allows dating
1571 estimates to be made using their Bayesian estimation framework using
1572 both strict and relaxed molecular clock models (Bouckaert et al. 2014).

1573 The HybridCheck project was created with the aim to help researchers
1574 understand the effects of recombination on genome sequence data. The
1575 software was written as a package for the R language, and it allows users
1576 to do the following.

- 1577 1. Evaluate the evidence of recombination in sequences.
- 1578 2. Identify recombination breakpoints and blocks.
- 1579 3. Estimate the age of recombinant blocks.
- 1580 4. Generate graphs to visualise the effects of recombination on the
1581 pattern of nucleotide similarity between sets of three sequences.

1582 The development of the package involved the following three stages:

- 1583 1. The R package was written to implement the functionality:

- 1584 (a) Conduct ABBA-BABA tests of introgression and calculate Patter-
1585 sons D, and Fd for four taxa or populations.
- 1586 (b) Scan alignments of 3 sequences for putative regions of recom-
1587 bination and generate plots of recombination signal from these
1588 Triplet Scans.
- 1589 (c) Automatically return putative regions of recombination from Triplet
1590 Scan data.
- 1591 (d) Calculate the probability that the high level of sequence similarity
1592 between two putative recombination regions is consistent with
1593 the mutation rate and sequence dissimilarity observed elsewhere
1594 in the sequence.
- 1595 (e) Estimate the 95% confidence interval for the coalescence time
1596 of a recombination region between two sequences (the donor
1597 and recipient). The algorithm assumes a molecular clock, and
1598 uses the binomial cumulative frequency distribution function.
- 1599 (f) Draw figures to visualise the (mosaic-like) genome structure
1600 and level of nucleotide (dis)similarity between sets of three se-
1601 quences.
- 1602 2. A user-friendly interface was developed by creating a web-app front-
1603 end for the R package. This used a framework called Shiny. This
1604 enables users that are unfamiliar with R to use the package as a
1605 web-app with a graphical interface, as well as an R code package.
- 1606 3. The performance of HybridCheck was evaluated using simulated data,
1607 and the package was assessed for the following criteria.
- 1608 (a) **False positive rate:** The detection of recombination regions in
1609 simulations without recombination.

1610 (b) **False negative rate:** A failure to detect recombination regions
 1611 or portions of recombination regions in simulated sequence data
 1612 with known recombination regions.

1613 (c) **Accuracy of block age estimates:** The accuracy of the esti-
 1614 mated coalescence time of detected recombinant blocks.

1615 2.2 Implementation

1616 2.2.1 Four Taxon Tests

1617 A Four Taxon Test (FTT) is implemented in HybridCheck to allow the user to
 1618 answer the question: *Is there evidence of recombination in my sequences?*
 1619 FTTs use two SNP patterns called ABBA and BABA to identify introgression
 1620 and require four sequences or populations, denoted as P1, P2, P3, and P4.
 1621 In addition, FTTs assume a phylogeny where P1 and P2 coalesce first to
 1622 form a taxonomic unit, which then coalesces with P3, and finally P4/A is
 1623 the out-group with the longest branch. The ABBA SNP pattern is expected
 1624 to be in abundance when introgression has occurred between P2 and P3
 1625 and the two populations share the derived allele i.e. the allele that is not
 1626 ancestral (the A in ABBA and BABA). Conversely, the BABA SNP pattern
 1627 is expected to be in abundance when introgression has occurred between
 1628 P1 and P3. Statistics computed for a FTT quantify the abundance of these
 1629 two SNP patterns. The FTT implemented in HybridCheck calculates two
 1630 statistics; Pattersons D, and F (Durand et al. 2011).

1631 Pattersons D in equation 2.1 tests for an excess of ABBA or BABA SNPs
 1632 between four populations:

$$D(P1, P2, P3, A) = \frac{\sum_{i=1}^n C_{ABBA}(i) - C_{BABA}(i)}{\sum_{i=1}^n C_{ABBA}(i) + C_{BABA}(i)} \quad (2.1)$$

1633 $C_{ABBA}(i)$ and $C_{BABA}(i)$ are defined as a binary count of whether the

1634 ABBA or BABA pattern is observed or not at site i if four sequences are
 1635 used. Alternatively if population samples are used $C_{ABBA}(i)$ and $C_{BABA}(i)$
 1636 are more generally defined using equations 2.2 and 2.3.

$$C_{ABBA}(i) = (1 - \hat{p}_{i1})\hat{p}_{i2}\hat{p}_{i3}(1 - \hat{p}_{i4}) \quad (2.2)$$

$$C_{BABA}(i) = \hat{p}_{i1}(1 - \hat{p}_{i2})\hat{p}_{i3}(1 - \hat{p}_{i4}) \quad (2.3)$$

1637 Where \hat{p}_{ij} is the frequency of the derived allele at site i in population
 1638 j . Pattersons D is expected to be 0 where no introgression has occurred
 1639 between the populations (Durand et al. 2011).

1640 The \hat{F}_d statistic is defined as the fraction of the genome shared through
 1641 introgression (Martin, Davey, and Jiggins 2014). The equation uses the
 1642 same numerator as that of the formula for Pattersons D , which is given the
 1643 name S and denotes the difference between the number of ABBA sites and
 1644 BABA sites, as per equation 2.4.

$$S = \sum_{i=1}^n C_{ABBA}(i) - C_{BABA}(i) \quad (2.4)$$

1645 The formula for the \hat{F}_d statistic compares this observed value of S , de-
 1646 noted as $S(P_1, P_2, P_3, P_4)$, with a value of S estimated under a scenario of
 1647 introgression (Martin, Davey, and Jiggins 2014). Specifically HybridCheck
 1648 considers two scenarios and computes \hat{F}_d for both: Complete introgression
 1649 between populations 2 and 3 and complete introgression between popula-
 1650 tions 1 and 3. These two scenarios are denoted as $S(P_1, P_D, P_D, P_4)$ and
 1651 $S(P_D, P_2, P_D, P_4)$ respectively. In both scenarios, P_D is the donor popula-
 1652 tion and is chosen by finding which of the introgressed populations has a
 1653 higher frequency of the derived allele (Martin, Davey, and Jiggins 2014).
 1654 Therefore, the two formulas for \hat{F}_d that are used by HybridCheck are given

1655 as equations 2.5 and 2.6.

$$\hat{f}_d = \frac{S(P_1, P_2, P_3, P_4)}{S(P_1, P_D, P_D, P_4)} \quad (2.5)$$

$$\hat{f}_d = \frac{S(P_1, P_2, P_3, P_4)}{S(P_D, P_2, P_D, P_4)} \quad (2.6)$$

1656 When calculating the FTTs, HybridCheck will break up the sequence
1657 alignment into a user definable number of blocks of a given length, and will
1658 compute for each block:

- 1659 1. Pattersons D.
- 1660 2. The two \hat{f}_d statistics (one for each of the two scenarios of introgres-
1661 sion).
- 1662 3. A Pvalue based on the binomial distribution.
- 1663 4. The number of sites that have a higher ABBA score.
- 1664 5. The number of sites that have a higher BABA score.

1665 These blocks are then used perform a jackknife to compute jackknife
1666 estimates, standard deviation, and Z scores for the four populations of the
1667 whole alignment. The binomial P-values computed for each block used with
1668 Fishers combined probability formula to calculate an overall binomial based
1669 P-value for the entire alignment.

1670 HybridCheck can be directed by the user to use certain populations in
1671 the place of P1, P2, P3, and P4. Alternatively it can automatically generate
1672 combinations of four populations and then decide which of the populations
1673 should be assigned which of the four positions, using the distances between
1674 the sequences. The statistics calculated in the four-taxon tests have been
1675 described and their performance evaluated in previous work by Martin,
1676 Davey, and Jiggins 2014. HybridCheck can be directed by the user to use

1677 certain populations in the place of P1, P2, P3, and P4. Alternatively it can
1678 automatically generate combinations of four populations and then decide
1679 which of the populations should be assigned which of the four positions,
1680 using the distances between the sequences.

1681 **2.2.2 Sequence triplet scans for recombination signal**

1682 A sliding window scan of pairwise sequence similarity for three sequences
1683 (hereafter referred to as a triplet) was implemented in HybridCheck to allow
1684 the user to answer the question: *Where are the breakpoints / regions of*
1685 *recombination located in the sequences?* HybridCheck was designed to
1686 generate and scan every possible triplet for a multiple sequence alignment.
1687 In addition, HybridCheck can be set to ignore triplets that include two or
1688 more sequences that are highly similar, reducing the number of scans to
1689 be performed. HybridCheck can also analyse a user-defined subgroup of
1690 sequences, or use the results of the four-taxon tests to generate the sets of
1691 triplets that need to be analysed. All non-polymorphic sites are removed
1692 from each triplet prior to the sequence scans.

1693 Potential recombinant regions are identified from the sliding window
1694 similarity scan data based on significantly elevated levels of sequence
1695 similarity. The cut-off point to identify elevated similarities is found by
1696 calculating the kernel density distribution of all raw sequence similarity data
1697 and identify peaks that fall outside this distribution. The start and end points
1698 of peaks are recorded (in base pairs) as well as the number of mutations
1699 within the block.

1700 The exact probability that the nucleotide similarity within a block is
1701 significantly higher than the overall sequence average can be calculated by
1702 modeling the accumulation of mutations as a Bernoulli trial. The probability
1703 of observing k or fewer mutations in a nucleotide sequence alignment of
1704 two sequences of length n is given by equation 2.7.

$$Pr(X \leq k) = \sum_{i=0}^{\lfloor k \rfloor} \binom{n}{i} p^i (1-p)^{n-i} \quad (2.7)$$

1705 In this equation (2.7), p is the proportion of observed single nucleotide
 1706 polymorphisms (SNPs) between the two aligned sequences (including
 1707 the non-informative sites). If the probability falls below the Bonferroni
 1708 corrected critical value $\alpha = 0.05$, the amount of polymorphism in the block
 1709 is inconsistent with the level of polymorphism that is expected from the
 1710 accumulation of mutations. In this case, recombination is taken to be a
 1711 valid explanation for the number of observed substitutions.

1712 2.2.3 Estimating the age of recombinant regions

1713 HybridCheck can estimate the coalescence times of the introgressed blocks.
 1714 This time is estimated assuming a strict molecular clock and using the
 1715 observed number of SNPs in the introgressed block. In order to correct for
 1716 mutation saturation, homoplasy, back mutations and transition / transversion
 1717 ratios, HybridCheck converts the number of SNPs into the number of
 1718 mutations using a JC (Jukes and Cantor 1969), K80 (Kimura 1980), F81
 1719 (Felsenstein 1981), HKY (Hasegawa, Kishino, and Yano 1985), or GTR
 1720 (Tavare 1986) correction.

1721 Considering the mutation accumulation process as a Bernoulli trial, and
 1722 the coalescence time can be found by finding the root of the equation 2.8.

$$f(n, k, 2t, Pr(X \leq k)) = \left(\sum_{i=0}^{\lfloor k \rfloor} \binom{n}{i} 2\mu t^i (1 - 2\mu t)^{n-i} \right) - Pr(X \leq k) \quad (2.8)$$

1723 In equation 2.8, μ is the mutation rate, t the time in generations, k the
 1724 observed number of SNPs, and n the total number of base pairs in the
 1725 block. The R function uniroot computes the value for $2\mu t$ by finding the root

1726 (i.e., the zero value) of function 2.8. (Brent 1973). In order to calculate the
1727 median and 5-95%CI, the function is solved for $2\mu t$ when Pr is set to 0.5,
1728 0.05 and 0.95.

1729 2.2.4 Performance Testing

1730 HybridCheck was tested on sequence triplets of 50kb in length which
1731 contained no introgression events to quantify its false positive rate α (i.e.
1732 erroneously identifying recombination). The simuPOP Python module
1733 (Peng and Kimmel 2005) was used to simulate three populations with 500
1734 individuals that derived from a single panmictic ancestral population, and
1735 which continued to evolve in genetic isolation. The populations diverged for
1736 between $0.01 \leq \mu t \leq 0.1$ generations (this is equivalent e.g. to $t = 1$ to 10
1737 million generations with $\mu = 10^{-8}$ base mutation rate). Sequence triplets
1738 were generated by randomly sampling one sequence from each of the three
1739 populations. A total of 100 independent sequence triplet replicates were
1740 generated for each simulated level of divergence (μt).

1741 HybridCheck was also tested on 50kb sequence triplets which contained
1742 set known introgression events of various ages to assess the sensitivity of
1743 the software to detect hybridization and the false-negative (β) rate. These
1744 triplets were also generated by simuPOP simulations in which two parental
1745 sequences diverged for between $0.02 \leq \mu t \leq 0.08$ generations, exactly as
1746 in the false positive error simulations described above. However, unlike
1747 the false positive error simulations, two subsequent steps were simulated:
1748 The parental sequences recombined at a user-defined breakpoint at 35
1749 kb, generating a third recombinant sequence. Then, in order to age the
1750 introgression blocks, the three sequences diverged for another $\mu t = 0:0.1$
1751 generations under a JC69 model, and during this time, the signal of intro-
1752 gression becomes eroded by mutation.

1753 Finally, the accuracy of the dating algorithm was tested using a regres-
1754 sion analysis. This used the same simulated data as was generated for
1755 evaluating the type II error rate. The estimated age calculated by Hybrid-
1756 Check was the response variable in the regression, and regressed the
1757 known coalescence time of the recombinant blocks in the simulations, was
1758 the explanatory variable.

1759 2.3 Results

1760 The false positive rate is presented in Figure 2.1, plotted on the y-axis
1761 against the amount of divergence (expressed as μt) between sequences on
1762 the x-axis. Depending on the divergence time of the populations, the false
1763 positive rate decreased with increasing sequence divergence but remained
1764 consistently less than $\alpha=0.05$. This means that if a triplet of sequences
1765 is analysed for recombination with HybridCheck, the more diverged they
1766 are from each other, the less likely it is that blocks will be falsely identified
1767 as putatively recombinant, when in fact no recombination has taken place.
1768 From this, one may conclude that recombination detection analyses can be
1769 confounded when populations or sequences analysed are not very diverged
1770 from each other, and that apparent recombination blocks or signals may be
1771 explained by other factors. Such facts include ancient population admixture
1772 or incomplete lineage sorting, and this will be addressed in more detail in
1773 the discussion.

1774 The false negative rate is presented in Figure 2.2. The false negative
1775 rate is plotted on the y-axis, against the amount of time since recombination
1776 occurred (expressed as μt). The data is partitioned into series, according to
1777 the amount of divergence (expressed as μt) between parental sequences
1778 prior to hybridisation. Figure n+1 shows that HybridCheck was able to detect
1779 >95% of recent introgression events even if the two parental populations

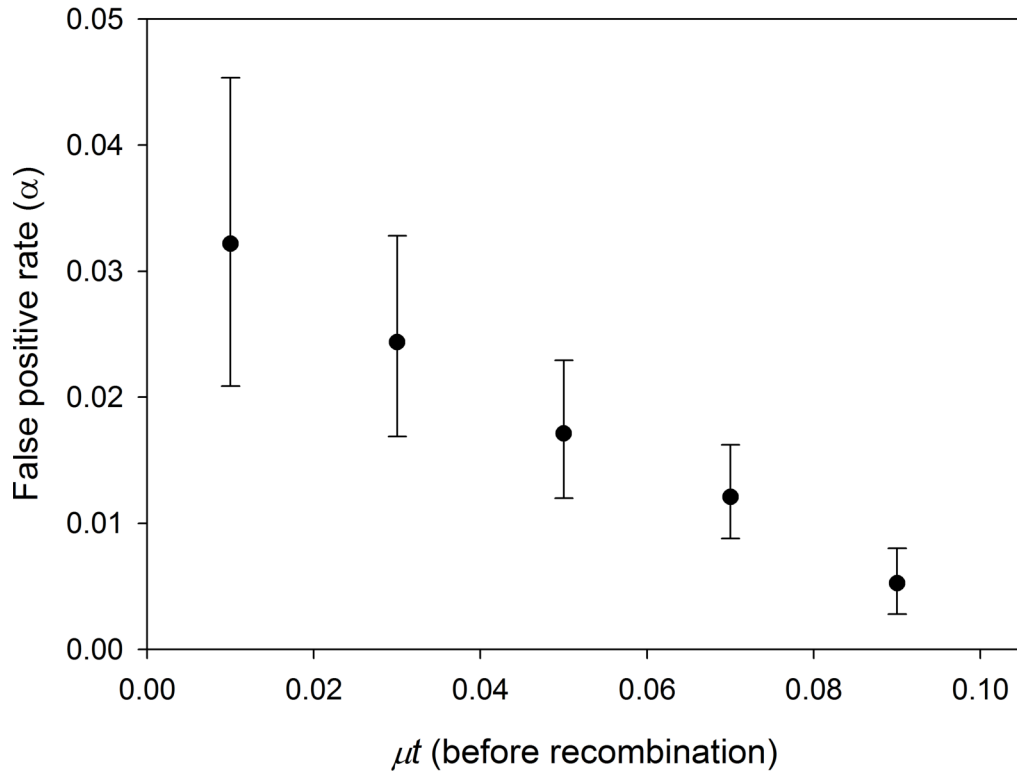


Figure 2.1: The mean(± 5 - 95%CI) false positive rate (α) of HybridCheck as a function of the ancestral divergence time μt (i.e. the amount of time of the sequences diverged before recombination). As sequences become more diverged, the false positive rate decreases.

1780 had diverged only moderately. However, more ancient introgression events
 1781 were detected only if both parental populations had significantly diverged.

1782 The accuracy of the dating estimates HybridCheck calculates for our
 1783 simulated scenario is presented in Figure 2.3. This analysis shows that
 1784 when the ancestral sequences had diverged significantly ($\mu t \geq 0.2$), the
 1785 age estimates calculated by HybridCheck are a good approximation of the
 1786 actual time passed since recombination (Linear Regression: Estimated age
 1787 = $0.000795 + 0.968 t$, $R^2=99.3\%$). However, when the exchanges occurred
 1788 between sequences that were only moderately diverged ($\mu t < 0.2$), the
 1789 age of the recombination events are underestimated when recombination
 1790 happened in the distant past ($\mu t > 0.05$) (see Figure 2.3). In such cases,
 1791 mutations accumulated after the recombination event fragmented the blocks,
 1792 resulting in an underestimate of the number of SNPs in the blocks that were

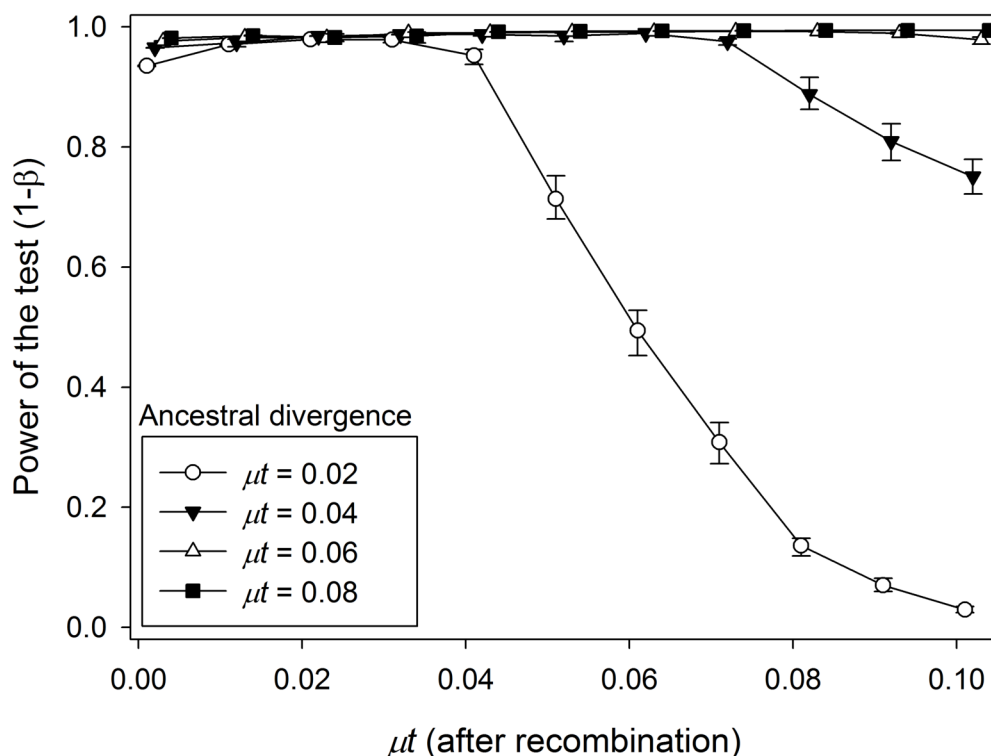


Figure 2.2: The mean(± 5 95%CI) statistical power ($1 - \beta$) of HybridCheck as a function of the divergence time of the sequences after recombination (expressed in μt) for sequences with ancestral divergence times $\mu t = 0.2, 0.4, 0.6$ and 0.8 generations. Recombination between moderately diverged sequences can be detected in $>95\%$ of the cases, as long as the recombination event was relatively recent.

1793 detected.

1794 2.4 Discussion

1795 In this project, the objectives were to create and test a software package for
 1796 the exploratory analysis of large sequences for evidence of introgression
 1797 and hybridization. The package is designed to take the researcher through
 1798 the following questions:

- 1799 1. Is there evidence of recombination / introgression?
- 1800 2. Where are the recombination regions in the sequences?
- 1801 3. What is the divergence time of recombinant blocks that are detected

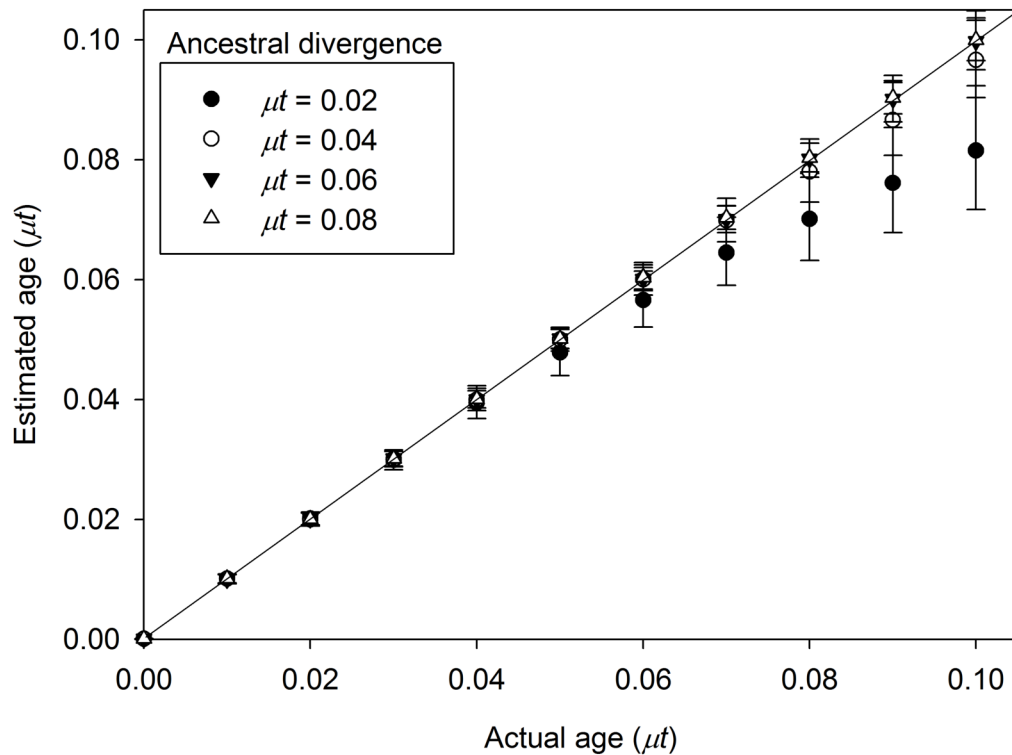


Figure 2.3: The mean (\pm SEM) estimated age (expressed in μt) of recombinant blocks calculated using the dating algorithm with a JC correction in HybridCheck, versus their actual age. In most of the scenarios, HybridCheck returns an unbiased estimate of the divergence time. However, the age is underestimated in cases of ancient recombination between populations that have ancestral divergence of 0.2.

1802 by the package?

1803 2.4.1 Performance of detecting recombinant blocks

1804 The data demonstrate that for the simulated scenarios, HybridCheck per-
 1805 forms best when sequences are diverged sufficiently prior to hybridization,
 1806 and the hybridization or recombination event was relatively recent. How-
 1807 ever, when the parental sequences of the hybrid sequence were sufficiently
 1808 diverged recombinant blocks were clearly detected long after the recombina-
 1809 tion event ($\mu t > 0.06$). In addition, when divergence between parental
 1810 sequences of a hybrid sequence was high then dating estimates of the
 1811 recombinant blocks remained more accurate for older recombination blocks.

1812 If two parental sequences are significantly diverged prior to hybridisation,

1813 the introgressed regions will be more apparent in the sequence similarity
1814 scans of HybridCheck because their high nucleotide similarity stands in
1815 sharp contrast with the genomic background. With a lower level of ancestral
1816 divergence, the increase in local sequence similarity caused by recombina-
1817 tion is more difficult to distinguish from stochastic variation in nucleotide
1818 divergence, around a higher average level of sequence similarity. As a re-
1819 sult, the algorithm HybridCheck employs to decide on a suitable sequence
1820 similarity threshold can be confounded as it tries to identify regions with
1821 sequence similarity that fall outside of the mean noise levels of sequence
1822 similarity. Therefore, HybridCheck would struggle to analyse a study system
1823 in which populations or taxa analysed are too closely related and have not
1824 diverged for long enough to accumulate unique polymorphisms which will
1825 be shared between parental and hybrid sequences.

1826 Previous studies have shown that many window based recombination
1827 detection methods perform better when the divergence is above 0.05
1828 (expressed as a proportion of the sequence length) (Posada and Crandall
1829 2001). Furthermore, simple implementations such as MAXCHI, and site
1830 incompatibility based methods usually perform better than phylogenetic
1831 based methods because the latter only tend to detect recombination if it
1832 changes the tree topology (Posada and Crandall 2001).

1833 HybridCheck window scans attempt to find elevated similarity between
1834 genome sequences / contigs of two taxa which are unrelated. Such eleva-
1835 tions in similarity are indicative of, and often coincide with incongruence
1836 between differing gene tree topologies. However, such signatures can
1837 have causes other than recombination, and elevated levels of sequence
1838 similarity could also be due to stabilizing selection conserving sequences
1839 between populations. Alternatively, diverging populations of organisms
1840 could show increased levels of divergence in regions of the genome that
1841 are under adaptive selection, and if there is gene flow between populations

1842 the background will be homogenized compared to the regions that are sub-
1843 ject to divergent selection (Nadeau et al. 2012). Such genomic islands of
1844 divergence appear to be less evident between populations that are further
1845 along the speciation process in butterfly species (Nadeau et al. 2012). A
1846 selective sweep is a phenomenon whereby positive selection in an allele
1847 reduces variation in neighboring regions due to linkage. This is also called
1848 hitchhiking (Hedrick 1980). If a selective sweep is strong and only one
1849 haplotype exists in high numbers in the population as a result, then a large
1850 reduction in variation is possible. Selective sweeps could create regions of
1851 sequence similarity similar to those created by hybridisation events. Note
1852 however, this scenario reduces variation around a positively selected allele
1853 within in a population.

1854 HybridCheck attempts to overcome these effects of selection in part by
1855 removing non-polymorphic sites prior to measuring the sequence similarity
1856 across sequences, but it is still possible that selection could be responsible,
1857 and the removal of informative sites by selection therefore reduces the
1858 power of HybridCheck to reliably identify introgression in those regions.
1859 Therefore HybridCheck is not recommended or useful if a researcher is
1860 interested in smaller regions subject to very strong selection, due to the
1861 resulting lack of information. If there are protein-coding regions in a detected
1862 recombinant region and selection is thought to be responsible, then the
1863 sequences should be analysed for evidence of purifying selection and/or
1864 selective sweeps within the detected region.

1865 Elevated sequence similarity and incongruent tree topologies can also
1866 be caused by incomplete lineage sorting or deep coalescence (Rogers
1867 and Gibbs 2014). This occurs when an ancestral species is polymorphic
1868 for a given gene before the species tree splits into two daughter species.
1869 After the first species split, if the polymorphism does not become resolved
1870 into two separate monophyletic lineages before the next speciation event,

1871 then the species tree will not match the gene trees of individual alleles
1872 (Rogers and Gibbs 2014). This problem is likely if a population size is very
1873 large, or if the time between branching events is low (Rogers and Gibbs
1874 2014). Much of the genome of *Homo sapiens* shows evidence of incomplete
1875 lineage sorting. As a consequence, parts of the genome supported the
1876 phylogeny (chimpanzee, (human, gorilla)), whereas other regions of the
1877 genome supported the phylogeny (human, (chimpanzee, gorilla)) (Galtier
1878 and Daubin 2008). Both these phylogenies disagree with the species
1879 phylogeny of homonids (gorilla, (human, chimpanzee)) (Galtier and Daubin
1880 2008; Rogers and Gibbs 2014). This discordance is because selection can
1881 cause similar sequences, or islands of divergence as previously described,
1882 and then incomplete lineage sorting results in gene trees that are discordant
1883 with the species tree and other gene trees, as a result of the incomplete
1884 and stochastic resolutions of polymorphisms, before subsequent speciation
1885 events (Sally et al. 2012).

1886 However, HybridCheck can help discern recombination from incomplete
1887 lineage sorting by comparing the coalescence time of recombinant regions
1888 with the split of the species. If the age of a recombinant region is significantly
1889 younger than the split of the ancestral species, the pattern is inconsistent
1890 with incomplete lineage sorting. In this case, genetic introgression after
1891 hybridisation is a more plausible explanation for the observed increase in
1892 local sequence similarity. HybridCheck makes this practically possible for
1893 the researcher to do, for many recombinant blocks.

1894 To summarise the performance of the HybridCheck when identifying
1895 recombinant regions, the HybridCheck use case is intended predominantly
1896 as an exploratory method to scan for signal between sequences from
1897 diverged populations or taxa, rather than within populations. Outside of this
1898 use case, HybridCheck may be unsuitable for some systems as a result
1899 of limited divergence between sequences, and selection, both of which

1900 result in reduced information for the HybridCheck analysis method. Recent
1901 speciation and large population sizes may result in incomplete lineage
1902 sorting, which can affect patterns of divergence and ancestry in similar ways
1903 to recombination, however coalescent times computed by HybridCheck
1904 may help distinguish incomplete lineage sorting from recombination. When
1905 using HybridCheck for a study system outside of its designed use case,
1906 whilst it is useful for highlighting the regions of the genome affected by the
1907 above factors, regions should not be uncritically considered the result of
1908 hybridisation or recombination, and the alternative causes e.g. selection
1909 should be followed up and ruled out before any such conclusion.

1910 **2.4.2 Performance of estimating the age of recombina-** 1911 **tion events**

1912 From the results it is evident that the dating algorithm used in HybridCheck
1913 tends to underestimate the divergence time of recombinant blocks in old
1914 recombination events. This is because recombination blocks can become
1915 fragmented by accumulation of subsequent mutations following the recombina-
1916 tion event. Consequently, older recombination blocks tend to be smaller,
1917 when they are actually larger. Thus, not all mutations are accounted for,
1918 resulting in an underestimate of the divergence time particularly for old
1919 recombination events.

1920 Furthermore, the dating algorithm used in HybridCheck makes several
1921 assumptions in order to be simple and fast. As a result however, if these as-
1922 sumptions are broken then this will affect how representative the estimates
1923 returned by HybridCheck are of the true age of a recombination event. The
1924 algorithm assumes that the mutation rate has been constant over time and
1925 identical in all taxa. This assumption is not always true, and more sophis-
1926 ticated approaches, such as the Fossilized-Birth-Death process allow for

1927 the calibration of divergence time estimates during Bayesian phylogeny
1928 estimation (Heath, Huelsenbeck, and Stadler 2014). It uses all available
1929 fossils, and considers extant species and fossils of species part of the same
1930 macro-evolutionary process (Heath, Huelsenbeck, and Stadler 2014).

1931 In addition, the algorithm uses a nucleotide substitution rate to infer
1932 the mutation rate. In order to correct for mutation saturation, homoplasy,
1933 back mutations and transition / transversion ratios, HybridCheck converts
1934 the number of SNPs into the number of mutations using a JC (Jukes and
1935 Cantor 1969), K80 (Kimura 1980), F81 (Felsenstein 1981), HKY (Hasegawa,
1936 Kishino, and Yano 1985), or GTR (Tavare 1986) correction. However,
1937 substitution rates do not solely depend on mutation rates, and they appear
1938 to be auto-correlated across sequences due to the effect of selection.
1939 Selection can vary between sites, genes and taxa, and selection and
1940 substitution rates can change through time as conditions change (Barrick
1941 and Lenski 2013; Bromham and Penny 2003).

1942 Furthermore, the size of populations must be taken into account (Bromham
1943 and Penny 2003). Bayesian coalescent approaches incorporated in soft-
1944 ware such as BEAST (Bouckaert et al. 2014) should be used when using a
1945 relaxed clock or more advanced method of dating. However, these methods
1946 are computationally more demanding and might become unfeasible when
1947 estimating the divergence time of a large number of recombination events.
1948 In such cases, the age estimate returned by HybridCheck offers a good
1949 approximation when recombination occurred relatively recently ($\mu t < 0.05$),
1950 and also when the ancestral sequences have diverged significantly before
1951 hybridizing.

1952 In conclusion, the HybridCheck project is intended as a simple all-
1953 inclusive tool to analyse recombination in genome sequence data. The
1954 implemented algorithms are not as sophisticated as methods that employ
1955 Bayesian estimation of parameters and coalescent simulations. However,

1956 this means that the package is computationally fast, which makes it a useful
1957 first port-of-call for identifying recombination and assessing whether other
1958 explanations such as incomplete lineage sorting may apply.

1959 CHAPTER 3

1960 The role of introgression in the adaptive 1961 evolution of the generalist plant pathogen, 1962 *Albugo candida*

1963 This chapter is based on the published scientific paper:

1964 *McMullan, M., Gardiner, A., Bailey, K., Kemen, E., Ward, B. J., Cevik, V.,*
1965 *... Jones, J. D. (2015). Evidence for suppression of immunity as a driver*
1966 *for genomic introgressions and host range expansion in races of Albugo*
1967 *candida, a generalist parasite. eLife, 4, 1-24.*

1968 This thesis chapter presents a research project that was a collaboration
1969 between many researchers. In this chapter in order to provide clear de-
1970 scription of the work involved, some details regarding some work that has
1971 not been performed by myself are presented. Specifically, work described
1972 in sections 3.2.1 and 3.2.2 were completed by collaborators and not myself.
1973 My contributions to the work are described in sections 3.2.3 and 3.2.4, and
1974 it is results of this work that is presented in this chapter.

1975 **3.1 Introduction**

1976 Host specificity is a defining feature of pathogens, and can be defined as
1977 the inverse of the number of hosts that a given pathogen can infect (Poulin
1978 and Keeney 2008). Host specificity is negatively correlated to the probability
1979 of parasite extinction, and positively correlated to the ability of a parasite to
1980 colonise and adapt to a new host (Poulin and Keeney 2008). Host specificity
1981 is constrained by the physiology of the pathogen. Therefore the host
1982 specificity of a pathogen is constrained by factors including (but not limited
1983 to) the pathogen's method of transmission, method of obtaining nutrients
1984 and energy from the host, and the ecology of the pathogen and host (Poulin
1985 2011). Such factors are proximal constraints on host specificity, but host
1986 specificity is ultimately constrained by the evolutionary and biogeographical
1987 history of the pathogen and its potential hosts (Poulin and Keeney 2008;
1988 Poulin 2011).

1989 A highly specialist parasite occurs in only a single host species. They
1990 often require host-host contact for transmission, and their longevity and
1991 future is strongly linked to that of their host species (Poulin and Keeney
1992 2008). Conversely, a parasite that is more generalist may survive the
1993 extinction of one host species, since there is another host species they can
1994 exploit to survive. Generalist parasite species may rely less on contact-
1995 transmission or close proximity between hosts. For example, they may be
1996 transmitted through food, or some other species vector (Pedersen et al.
1997 2005). However it should be noted that even if a pathogen has a very high
1998 mobility, and dispersal, its host-specificity can be high (Poulin and Keeney
1999 2008).

2000 The availability of ecologically and evolutionarily related or similar hosts
2001 cohabiting the same habitat, may cause differences in the host-specificity of
2002 two otherwise similar pathogen species (Jex, Schneider, and Cribb 2006).

2003 Furthermore, parasites put selective pressure on host populations to adapt
2004 and develop immunity, increasing the frequency of genetic and epigenetic
2005 variants that improve immunity, and pathogen detection in the host. As a
2006 result, these adapting host populations impose selection pressure on the
2007 pathogen populations, increasing the frequency of variants that maintain
2008 the pathogens efficiency of immune suppression. Over time, both host and
2009 parasite co-evolve and become intimately associated, as they both adapt to
2010 each other's latest antagonistic evolutionary innovations. This is called an
2011 evolutionary arms-race (Boutemy et al. 2011; Buckling and Rainey 2002;
2012 Cooper et al. 2008; Kemen and Jones 2012; Lamour et al. 2010).

2013 Overall, the general pattern observed in nature, is that most parasite
2014 species are largely specialised and co-evolve with only a few, if not one,
2015 host species (McMullan et al. 2015). It should be noted however, that this
2016 generalisation is based on a measure of host specificity that is based on
2017 a simple measure of host specificity, namely the number of host species
2018 that are colonised by a parasite in natural populations. However this metric
2019 makes an oversimplification that does not reflect biological reality. For
2020 example, two pathogen species may have the same number of host species,
2021 but if one of the pathogens infects species of one genus, and the other
2022 infects species of multiple genera, then it is not realistic to conclude both of
2023 the parasite species are equally specialised. It is because of this problem,
2024 that Poulin and Mouillot 2003 defined a host-specificity measure that takes
2025 into account the taxonomic or phylogenetic distances between the hosts
2026 colonised by a parasite. Later the authors published an improved metric
2027 that incorporated the phylogenetic or taxonomic distinctness of a pathogens
2028 host species, but also weighted for the prevalence of the parasite on its
2029 different host species (Poulin and Mouillot 2005). The rationale for such a
2030 weighting is that a pathogen that is largely concentrated on only one of its
2031 multiple hosts should be classified as more specialised than a pathogen

2032 that utilises and colonizes all of its host species evenly.

2033 The organism of interest in this work is the obligate biotrophic plant
2034 pathogen, *Albugo candida*. Plant pathogens have a parasitic relationship
2035 with their host, and are classified according to the nature of this relationship
2036 with the host. Pathogens which obtain nutrients from decaying plant matter
2037 are classified as necrotrophs, whereas pathogens which require living
2038 host tissue in order to obtain nutrients are classified as obligate biotrophs
2039 (Kemen and Jones 2012). These biotrophs don't typically secrete abundant
2040 lytic enzymes, and cause little physical or structural damage to the host
2041 plant (Kemen and Jones 2012). Pathogens with a combination of these two
2042 lifestyles are classified as hemibiotrophic (Kemen and Jones 2012; Lamour
2043 et al. 2012). *Albugo candida* is an obligate biotroph, and whilst *Albugo*
2044 *candida* is a generalist, infecting species of the Brassica family, obligate
2045 biotrophs are typically specialists (McMullan et al. 2015).

2046 After an obligate biotroph makes a host-jump, it is expected that selec-
2047 tion will increase any adaptive genetic or epigenetic variant in the population
2048 that results in more efficient immune suppression of the new host (Dong
2049 et al. 2014; Kemen and Jones 2012; Poulin and Keeney 2008; Raffaele
2050 et al. 2010; Thines 2014). Furthermore, host-parasite co-evolution over
2051 time will result in both the host and parasite constantly adapting to each
2052 others latest antagonistic adaptations, and they will become more intimately
2053 associated historically (Morgan and Kamoun 2007; Raffaele and Kamoun
2054 2012; Thines 2014). As both of these processes occur, new effectors
2055 and pathogenicity factors may be created, and existing ones may receive
2056 beneficial mutations, and they may also have their levels of expression
2057 changed epigenetic modification and inheritance (Dong et al. 2014; Gijzen,
2058 Ishmael, and Shrestha 2014; Raffaele and Kamoun 2012; Raffaele et al.
2059 2010; Win et al. 2012). These will be fixed due to selection pressure if
2060 they are beneficial. These modifications enable more efficient immune

2061 suppression and exploitation of one host species, but increase the risk of
2062 detection in other host species by triggering their immune system (Martin
2063 and Kamoun 2012). Thus, as obligate biotrophic pathogen populations
2064 become more adept at suppressing the immunity of one host, they will
2065 become less adept at infecting previous host(s) or other hosts it can infect.

2066 Therefore, obligate biotrophs are typically known for being intimately
2067 associated with their hosts i.e. they have a high host specificity (Thines
2068 2014). Yet there are generalist biotrophic parasites that appear to have
2069 overcome this evolutionary dilemma and show virulence on diverse hosts.
2070 *Albugo candida*, the organism that is the subject of this work, is one such
2071 generalist, but there are other generalist oomycetes, like *Phytophthora*
2072 *capsici* (Lamour et al. 2012).

2073 Some generalist parasite species have solved the dilemma by evolving
2074 multiple specialised races, and each specialised race can infect a different
2075 host. For example, the eukaryotic order *Albuginales*, of which *Albugo can-*
2076 *dida* is a member, is completely comprised of obligate biotrophic pathogens
2077 that cause disease on a broad range of plant hosts (Biga 1955; Choi and
2078 Priest 1955; Walker and Priest 2007).

2079 *Albugo* is the largest genus of the order *Albuginales*, and it was reported
2080 to consist of 33 specialist pathogens by Biga 1955. More recently, the
2081 estimate is that the genus comprises approximately 50 pathogens, and
2082 these are typically specialists. In addition new distinct *Albugo* species
2083 have been discovered that were previously thought to be members of
2084 *Albugo candida* (Pers) Roussel. (Choi et al. 2011; Choi, Shin, and Thines
2085 2009; Ploch et al. 2010; Thines et al. 2009). This is because in the past
2086 decades, classification was based largely on morphology, and this led to
2087 the application of a broad species concept, that resulted in *Albugo candida*
2088 (Pers.) Roussel being regarded as the causal organism of all incidents of
2089 white blister rust on all *Brassicaceae* hosts (Choi et al. 2011). As late as

2090 2011, it has been estimated that a dozen distinct species thought to be
2091 *Albugo candida* await discovery (Lamour and Kamoun 2009).

2092 *Albugo candida* (Pers.) Roussel can infect 241 species of plants in 63
2093 genera from the families of *Brassicaceae*, *Cleomaceae* and *Capparaceae*
2094 (Choi, Shin, and Thines 2009). *Albugo candida* infections are the causal
2095 agent of white blister rust disease, resulting in significant losses on *Brassica*
2096 crops of economical importance. For example, *Albugo candida* causes
2097 up to 56 of yield losses in Indian Mustard (Meena et al. 2002). *Albugo*
2098 *candida* consists of different physiological races, each usually featuring
2099 high host-specificity and approximately 24 races of *Albugo candida* have
2100 been defined, based on their host range (Saharan et al. 2014; Saharan and
2101 Verma 1992).

2102 *Albugo candida* reproduces both asexually and sexually (Holub et al.
2103 1995). During asexual reproduction, diploid zoospores are formed in
2104 zoosporangium beneath the leaf epidermis. The zoosporangium are visi-
2105 ble when dehydrated and in large numbers, as white blisters (Holub et al.
2106 1995). These sporangia then rupture the epidermis of the host leaf, to
2107 release zoospores for dispersal. During sexual reproduction, fertilization
2108 between two isolates creates non-motile, diploid, and thick-walled oospores
2109 (Holub et al. 1995). The oospores can resist extreme temperatures and
2110 desiccation. The relative importance of both reproductive modes is not
2111 well established, but the clonal (asexual) mode of reproduction allows rapid
2112 population expansion, especially given modern crop mono-culture growing
2113 practices. Although *Albugo candida* comprises distinct, specialised physi-
2114 ological races that colonize different host plants, and that distinct species
2115 have been identified that were initially thought to be *Albugo candida* (Choi
2116 et al. 2011), it is still considered a single species.

2117 According to evolutionary and population genetic theory, the trade-offs
2118 associated with adaptation and host-specialisation, coupled with strong

2119 population structuring, can result in adaptive radiation and speciation (Ab-
2120 bott et al. 2013; Stukenbrock 2013). *Albugo candida* then may be thought of
2121 as a currently ongoing adaptive radiation; The broad host range of *Albugo*
2122 *candida* is enabled by an ongoing specialisation of independent physio-
2123 logical races, and these races are likely heading for speciation (Dres and
2124 Mallet 2002). If strains or races of a parasite develop adaptations to specific
2125 hosts, and make trade-offs in doing so, specialising to the given host, does
2126 parasite specialization inevitably lead to speciation? Certainly, specialising
2127 on one or a few hosts, at the cost of being able to infect other hosts, will
2128 mean separation of specialised races, ecologically, and even geographically,
2129 over time such separation is expected to result in reproductive isolation.

2130 Compared to other microbial plant pathogens, *Albugo* species are no-
2131 table as infections strongly suppress host innate immunity. As a result,
2132 infections of *Albugo* species increase the susceptibility of the host to a sec-
2133 ondary infection by pathogens that would otherwise be avirulent, including
2134 downy mildews (Cooper et al. 2008). It has been suggested that this im-
2135 mune suppression caused by *Albugo* infections might allow an accelerated
2136 adaptation of other pathogen species to host that is susceptible to *Albugo*
2137 species (Thines 2014).

2138 However, whilst it has been suggested the immune suppression will
2139 accelerate the adaptation of other pathogens to the suppressed host, be-
2140 fore this project, no evolutionary rationale was proposed explaining why
2141 rendering a host susceptible to other pathogens could be adaptive for the
2142 various *Albugo* species. Hypothetically, a pathogen which colonizes and
2143 adapts to the hosts of *Albugo* species due to the immune suppression of
2144 *Albugo* species infections, will become competition against *Albugo* species
2145 for the same resources (Cooper et al. 2008).

2146 Suppression of host innate immunity would facilitate cohabitation of
2147 distinct physiological races that otherwise would not come into contact

2148 due to their specialisation and adaptive trade-offs, as previously discussed.
2149 When the distinct physiological races come into contact, genetic exchange
2150 including introgression and hybridisation may occur between them. Here,
2151 introgression is defined as the introduction of nucleotide variation from a
2152 parental donor race into the genome of a recipient race, through the mech-
2153 anism of recombination (Hedrick 2013). This flow of genetic variation from
2154 one donor physiological race, to a recipient physiological race, could slow
2155 down the genetic divergence of the races, and slow or prevent speciation.
2156 However, introgression between races that are specialised and adapted to
2157 exploit different hosts could be maladaptive, and therefore could be strongly
2158 selected against. This is because hybrids would inherit effector alleles
2159 derived from both parental races. Therefore, whilst the hybrid genomes
2160 would contain effectors that enable the immune suppression of multiple
2161 hosts, they could also contain effectors that trigger immunity on multiple
2162 hosts. Immune recognition of even a single effector is sufficient to trigger
2163 the immune response and stop an infection. Therefore any hybrid that pos-
2164 sess an expanded repertoire of effector alleles are likely to have a strong
2165 fitness disadvantage on most potential host plants, as with larger effector
2166 repertoire's comes an increased likelihood of one of them triggering host
2167 immunity.

2168 This chapter presents work conducted and contributed to a larger
2169 genome project analysis of *Albugo candida*, conducted by a team of scien-
2170 tists at the University of East Anglia, and The Sainsbury Laboratory. This
2171 project aimed to answer the following questions, in order to try and resolve
2172 this question of whether immune suppression and secondary infection is
2173 adaptive or maladaptive, and whether it is due to hybridisation:

- 2174 1. Are the distinct physiological *Albugo candida* races genetically iso-
2175 lated and on the road to speciation?

- 2176 2. Does suppression of host innate immunity enable cohabitation and
2177 growth of races with non-overlapping host ranges?
- 2178 3. Are the genomes of *Albugo candida* affected by recombination and
2179 hybridisation?

2180 The work presented in this chapter was primarily conducted with a goal
2181 of answering the third question of the project. During the collaborative
2182 project, genome sequence assemblies were created for five isolates that
2183 were collected from four host species (*Brassica oleracea*, *Brassica juncea*,
2184 *Capsella bursa-pastoris*, and *Arabidopsis thaliana*). This chapter presents
2185 analyses performed on the assembled sequence scaffolds for the detection
2186 of recombination, hybridisation, and mosaic genome structure.

2187 **3.2 Methods**

2188 In order to perform the analysis of genome structure that is the focus of
2189 this chapter, prior work was conducted to isolate the *Albugo candida* races
2190 used in this study, test for virulence, extract and sequence DNA, and RNA,
2191 and perform genome and transcriptome assemblies. These procedures are
2192 subsequently described in detail in (McMullan et al. 2015), and given that
2193 these procedures are not the focus of this chapter, the reader is referred
2194 to this paper for details on the wet lab and molecular methods. A brief
2195 summary of these methods is described below.

2196 **3.2.1 Isolation and cultivation of races used in the study**

2197 In order to address the research questions presented in the previous sec-
2198 tion, genome sequence assemblies were required of five isolates of *Albugo*
2199 *candida*, the white rust fungus. These isolates were collected from four dif-
2200 ferent host species: *Brassica oleracea*, *B. juncea*, *Capsella bursa-pastoris*,

2201 and *Arabidopsis thaliana*. The isolates were collected by Erik Kemen, prior
2202 to the evolutionary analyses that are the focus of the present chapter.

2203 The isolate designated AcNc2, was isolated from infected leaves of
2204 *Arabidopsis thaliana* Eri-1 field-grown plants in Norwich, England. The
2205 isolate was collected in 2007. The isolate AcEm2 was collected from wild
2206 *Capsella bursa-pastoris* in Kent, England in 1993. AcBoT was collected
2207 from infected cultivars of *Brassica oleracea* called 'Bordeaux F1', from
2208 Lincolnshire, England, in May 2009. AcBoL was harvested from infected
2209 *Brassica oleracea* leaves from Lincolnshire, but in the January of 2009. An
2210 isolate which is virulent on *Brassica juncea* called Ac2V was provided by M
2211 Borhan of Agriculture and Agri-Food, Canada. All of these isolates were
2212 single spore purified (Kemen et al. 2011).

2213 3.2.2 Genome assemblies of isolates

2214 The assembly of isolate AcNc2 was used as a reference. The assembly
2215 is 34Mb in size, and has 5212 contigs of approximately 160-fold coverage.
2216 The assembly was approximately 73% of an estimated genome size of
2217 45Mb. The unassembled part of the genome (approximately 11Mb) is likely
2218 to contain repeats, approximately 8% of which represent collapsed regions,
2219 since they have coverage that is several times higher than the average.
2220 For each isolate, several assemblies were constructed with different k-mer
2221 lengths. Each assembly was assessed according to number of contigs, N50
2222 (Bp and number), mean contig length, assembly size, GC content, average
2223 genome coverage, repeat content, and the number of predicted genes.
2224 High sequence similarity of the five *Albugo candida* isolates resulted in the
2225 conclusion that three races had been sequenced: AcNc2, and AcEm2 were
2226 isolates of the same race, and AcBoT and AcBoL were also two isolates
2227 which belonged to the same race. Therefore, detection of recombination
2228 and hybridisation in this chapter were first conducted on the three races

2229 AcNc2, Ac2V, and AcBoT, each of which had a 33-34Mb assembly.

2230 **3.2.3 Detection of recombination events**

2231 Recombination events were statistically identified on contigs $\geq 10,000$ Bp
2232 using the software RDP3 using five independent detection algorithms: RDP
2233 (Martin and Rybicki 2000), GENECONV (Padidam, Sawyer, and Fauquet
2234 1999), Maxchi (Smith 1992), Chimaera (Posada and Crandall 2001), and
2235 3Seq (Boni, Posada, and Feldman 2007). All of these tests are available
2236 in the Software Package RDP for Microsoft Windows (Martin et al. 2015).
2237 Tests were conducted using a critical value $= 0.05$ and p-values were
2238 Bonferroni corrected for multiple comparisons of sequences. Sequences
2239 were made linear using unphased base calling, i.e. where a sequence has
2240 a base position that is heterozygous, one of the nucleotides was assigned
2241 at random at that site.

2242 Recombination events were only considered genuine if they were sup-
2243 ported by at least three of the recombination detection methods in RDP,
2244 and recombination events detected using the methods in RDP were only
2245 counted if the parental sequences could be identified, and the start and
2246 end positions of recombination events were unambiguous.

2247 In order to visualise the effects of recombination and hybridisation on
2248 the genome structure of the *Albugo candida* races, the software package
2249 HybridCheck was developed for the R programming language. The devel-
2250 opment and testing of this software package is described in detail in chapter
2251 2, so only a brief description will follow. HybridCheck can analyse three
2252 sequences with a sliding window scan, and produce plots with use the RGB
2253 tricolour system to indicate where regions of hybridisation or recombination
2254 have occurred between sequences. Each sequence is designated one of
2255 the three primary colours, red, blue, or green. In regions of a given se-
2256 quence that are unique, then those regions are coloured in with the unique

2257 colour of that given sequence. However, in regions of the sequences in
2258 which all the SNPs are shared with another sequence, then the region is
2259 coloured with the hybrid colour of the two sequences (e.g. yellow if the two
2260 sequences have the unique colours red and green). All monomorphic sites
2261 are excluded in this computation. In cases where recombination is recent,
2262 the hybrid colouration is strong as most of the SNPs are shared between
2263 two sequences. However older events may have accumulated mutations
2264 since the recombination / hybridisation event. In such a case, there are
2265 less shared SNPs between two sequences, and the colour intensity is less
2266 strong.

2267 **3.2.4 Dating identified recombination events**

2268 Immediately after a recombination or hybridisation event has occurred, a
2269 hybrid or recombinant offspring's DNA sequence will have regions which
2270 are near identical to one parent, and regions which are near identical to
2271 the other parent. In those regions the molecular clock is effectively zeroed.
2272 Therefore, for a given recombinant region, the only substitutions which
2273 could be observed between the recombinant and the donor must have
2274 occurred since the recombination event took place.

2275 This divergence between a donor sequence region, and the same
2276 region in the recombinant offspring was used to estimate the time since the
2277 recombination event. Two methods were used to calculate the number of
2278 generations since individual identified recombination events occurred. A
2279 binomial mass function was used, which was developed for the HybridCheck
2280 R package. The equations are described more fully in chapter 2. Briefly, the
2281 method computes a window of time, within which the recombination event
2282 is most likely to have occurred. It does this by taking into consideration
2283 the cumulative probability of observing the number of mutations that have
2284 occurred in the recombinant region, between donor and parent, given

2285 the average mutation rate. The function assumes that the recombination
2286 event has evolved neutrally since the recombination event occurred, and
2287 that mutation rates between the two sequences were constant through
2288 time, and equal in both sequences. The mutation rate in oomycetes is
2289 unknown, and therefore the binomial mass function was used with two
2290 different mutation rates: $\mu = 10^6$ and 10^7 per base per generation. This
2291 binomial mass function was used to analyse all detected recombination
2292 events.

2293 In addition to the binomial mass function, an analysis was conducted
2294 in BEAST (Drummond and Rambaut 2007). Phylogenetic trees were esti-
2295 mated with a HKY + G model, a Yule tree prior, and a strict molecular clock
2296 assumption, where the mutation rate was assumed to be $\mu = 10^6$. Ten
2297 independent analyses were run, with an MCMC of 10 million steps, with a
2298 burn-in of 10%. Because of the computational complexity and time required
2299 for BEAST analyses, 20 recombinant regions were analysed in this manner.
2300 The results were compared to the date that was estimated for the recom-
2301 binant region by the binomial mass function, and this confirmed that the
2302 binomial mass function provides a good approximation of the divergence
2303 time.

2304 **3.3 Results**

2305 **3.3.1 Distribution of polymorphisms across races**

2306 Polymorphisms were found to be unequally distributed across the genomes
2307 of the *Albugo candida* isolates analysed. In some regions of the genome,
2308 there are stretches of identical sequence which are as long as 10kb in
2309 length. In other regions of the genome, stretches of lower sequence
2310 similarity may be found. For example, between the isolates AcBoT and

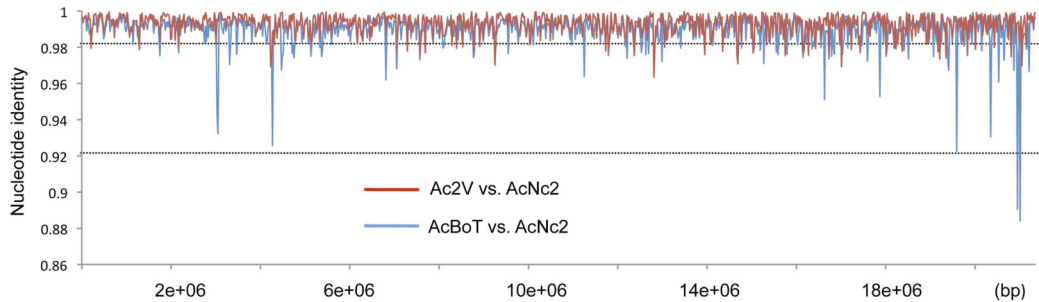


Figure 3.1: Nucleotide identity amongst the homologous genomic regions of Ac2V, AcBoT and AcNc2. The mean identity was calculated for the sliding window of 20 Kb.

2311 AcNc2, a region of approximately 5kb was observed with 89% sequence
2312 similarity. This is demonstrated in Figure 3.1.

2313 The distribution of the polymorphisms is highly suggestive of a mosaic-
2314 like genome as the polymorphisms are not only distributed unevenly, but
2315 they were distributed in a block-like manner. Stretches of nucleotide similar-
2316 ity are arranged in a block like structure; there are regions where AcNc2 is
2317 highly similar to AcBoT (and therefore diverged from isolate Ac2V), followed
2318 by regions where isolate AcNc2 is highly similar to Ac2V (and therefore di-
2319 verged from isolate AcBoT). The HybridCheck software package visualises
2320 such genome structure in Figure 3.2. The figure visualises the effect on
2321 the genome by colouring regions yellow where races AcNc2 and AcBoT
2322 show near sequence identity, cyan where races AcBoT and Ac2V show
2323 near sequence similarity, and purple where races AcNc2 and Ac2V show
2324 near sequence similarity. Note that in the figures, there are also regions
2325 of unique colouration (red, green, and blue), and such regions represent
2326 diverged parts of the genome where the three races have large proportion
2327 of unique (races-specific) polymorphisms (Figure 3.2).

2328 This observation of alternating blocks of high sequence identity between
2329 otherwise diverged (as represented by areas of red, green, and blue)
2330 genomes, provides supporting evidence for genetic introgression between
2331 diverged races that show a considerable (yet still incomplete) level of

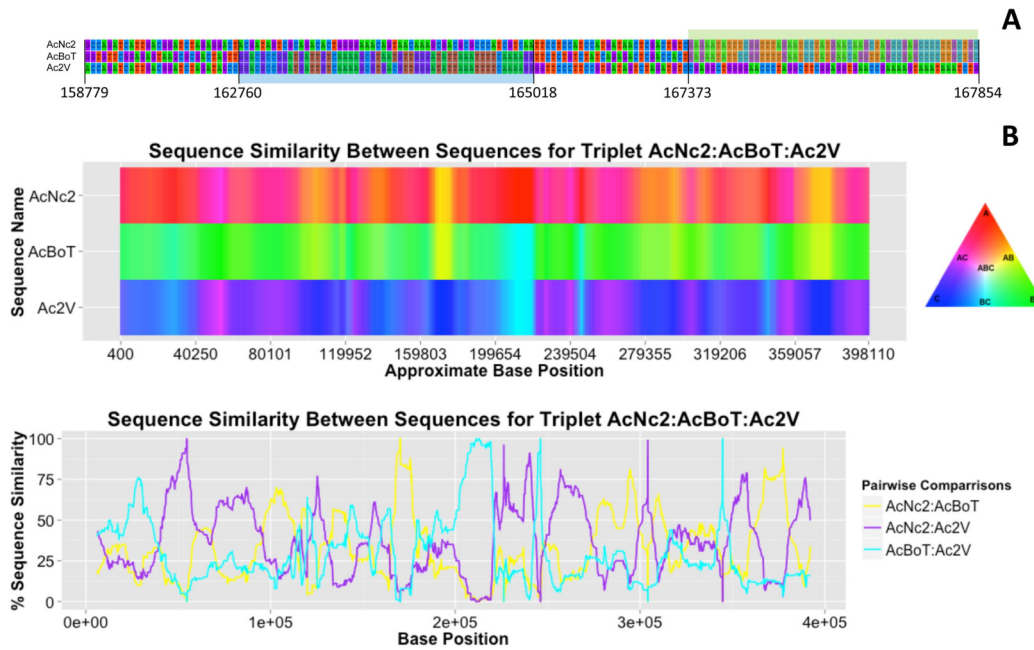


Figure 3.2: Extensive variation in sequence similarity between *Albugo candida* races. **A)** An sequence alignment between base positions 158,779 and 167,382 within contig 1 of *A. candida* races AcNc2, AcBoT and Ac2V. Two recombination blocks coloured blue and green are visible, displaying high sequence similarity between races. **B)** The sequence similarity across the length of contig 1, amongst three *A. candida* races. Similarity is visualised using the colours of a RGB colour triangle in the software HybridCheck. Areas where two contigs have the same colour (yellow, purple or turquoise) are indicative of two races sharing the same polymorphisms. The linear plot of the proportion of SNPs shared between the three pairwise comparisons between the races. Shown on the X-axis is the actual base position.

2332 reproductive isolation. The recombination detection methods described in
 2333 the previous section test for recombination blocks visualised here, formally.

2334 3.3.2 Recombination blocks identified using RDP

2335 All 133 contigs were analysed for presence of recombination blocks using
 2336 algorithms in the software package RDP. Recombination analysis with
 2337 these algorithms identified 675 recombination blocks on 127 sequence
 2338 contigs which were significant, even following correction of the alpha with a
 2339 Bonferroni correction. These identified blocks were reported as significant
 2340 for at least three different recombination detection tests. If the length of
 2341 all the significant blocks is summed in a linear fashion, then approximately

2342 25% of the total length of all contigs analysed is identified as recombinant,
2343 this is equal to 3Mb. These blocks represent regions of the genome which
2344 are derived from either another race, or the ancestor of another race.
2345 Algorithms in RDP were able to report such donor sequences in some
2346 cases. The full data-set from the RDP output is publicly available from
2347 <http://dx.doi.org/10.7554/eLife.04550.015>.

2348 **3.3.3 Estimated ages of recombination events**

2349 Dating analysis of the significant recombination blocks using the Hybrid-
2350 Check binomial algorithm indicated that the recombination events detected
2351 occurred at a range of different dates. If one assumes a $\mu = 10^{-8}$ sub-
2352 stitution rate which is constant across cell cycles, and that there are 100
2353 cell cycles per year, then the most recent introgression event occurred
2354 approximately 220 years ago, and the oldest detected event occurred al-
2355 most 200,000 years ago. The mean age for all the detected recombination
2356 events is approximately 6237 years ago, with a standard error of 12,594
2357 years. Furthermore, there is no significant difference between the average
2358 estimated dates across different contigs.

2359 The wide range in age estimates of the introgressed regions provides
2360 evidence for the hypothesis that recombination and hybridisation between
2361 diverging *Albugo candida* races has been a consistent and ongoing evolu-
2362 tionary process, affecting the entirety of the genome. This finding rules-out
2363 the hypothesis that one or a few recombination/hybridisation events in the
2364 distant past are responsible for creating the mosaic structure observed.
2365 This also helps explain the cause of the mosaic genome structure that has
2366 been observed: occasional introgression events across a range of evolu-
2367 tionary times is expected to result in genome containing introgressed blocks
2368 of sequence from a donor race, interspersed inside the distinct genomic
2369 background of the recipient race (i.e. the very pattern observed in *Albugo*

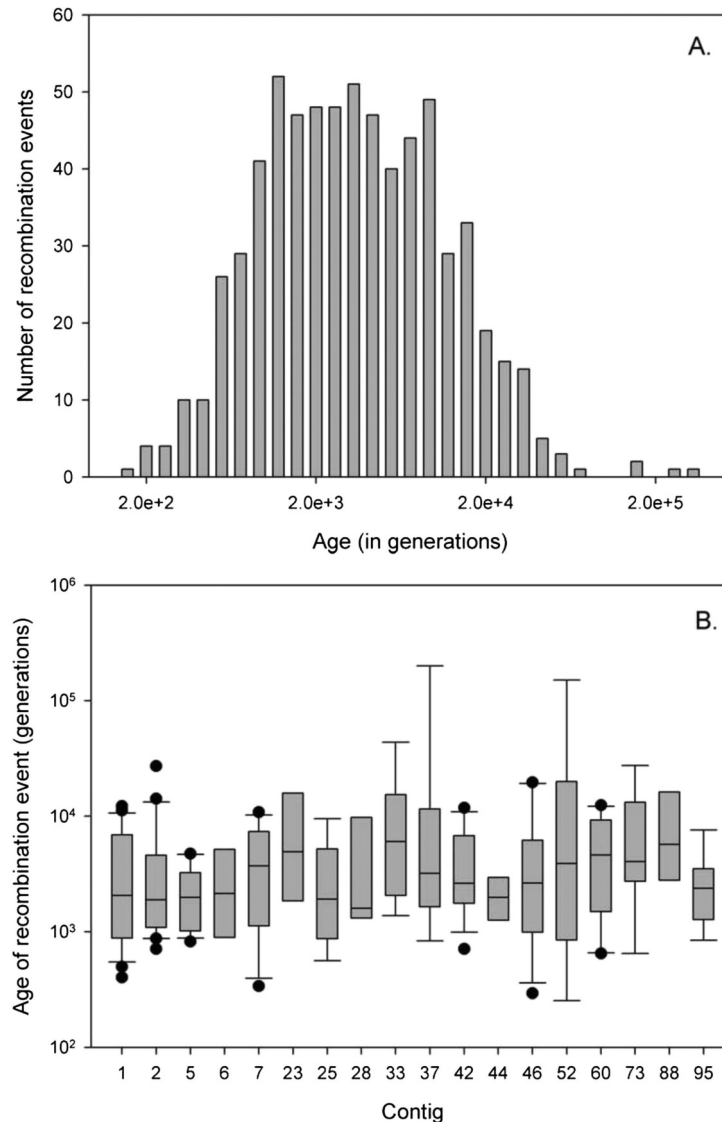


Figure 3.3: **A)** Age of the 675 recombination blocks, identified across the whole genome, estimated using the HybridCheck binomial mass function, assuming a substitution rate of $\mu = 10^6$; **B)** A box plot of the median (plus first and third quartile) log-age of recombination events in contigs. Only contigs with eight or more events are shown. There is no significant difference in age of events between contigs (GLM: $F_{22, 233} = 1.06$, $p = 0.387$).

2370 *candida*).

2371 3.4 Discussion

2372 The genome of *Albugo candida* appears to have a mosaic-like genome
 2373 structure: 675 regions were identified in 127 analysed contigs, which were
 2374 consistently identified by multiple and independent recombination detection

2375 methods. The mosaic-like structure reflects discordant phylogenetic signals
2376 of genomic regions with distinct coalescence, and this suggests that intro-
2377 gression has occurred at a range of time points throughout the evolutionary
2378 history of the *Albugo candida* races.

2379 **3.4.1 Hybridisation and clonal reproduction of *A. can-*** 2380 ***dida***

2381 *Albugo candida* is an obligate biotroph, growing and reproducing on living
2382 plant tissue, and virulence experiments confirm that the *Albugo candida*
2383 races isolated in this study are indeed host specific (McMullan et al. 2015).
2384 To explain the observed mosaic genomes, two distinct and host specialised
2385 *Albugo candida* races would have to make contact by colonizing the same
2386 host plant in order to hybridize, although ex-situ hybridisation cannot be
2387 ruled out. Yet, any *Albugo candida* race landing on a non-host plant is
2388 likely to trigger host immunity before it can mate with another distinct race.
2389 So, given that the genome structure expected from recent introgression
2390 between distinct races is observed, how have they made contact? One
2391 potential explanation was that infected host plants could form secondary
2392 contact zones for *Albugo candida*: if a host plant was infected by a com-
2393 patible (infectious) *Albugo candida* race its immunity would be suppressed.
2394 With a suppressed immune system, non-specialised *Albugo candida* races
2395 might be able to colonise the already infected host, enabling both races to
2396 make contact and hybridise through sexual reproduction. This hypothesis
2397 was tested with experimental infections of host plants with multiple races.
2398 These experiments confirmed that a virulent race of *Albugo candida* could
2399 suppress the immunity of its host plant, such that other non-virulent races
2400 of *Albugo candida* could co-colonise it (Cooper et al. 2008; McMullan et al.
2401 2015).

2402 Following formation of a viable hybrid, clonal reproduction would allow
2403 fast dispersal of the pathogen and population expansion. This aspect of the
2404 model was supported by analysing genomic identity between isolates which
2405 infect the same host species (i.e. within different races) and quantifying the
2406 shared proportion of heterozygous sites. Genotypic similarity at heterozy-
2407 gous sites of pairs of independent isolates that infect the same host plant
2408 was exceptionally high; AcBoT and AcBoL shared 97% of their heterozy-
2409 gous sites in common, and AcEm2 and AcNc2 shared 99.95%. Sharing
2410 of this proportion of heterozygous sites rules out Mendelian segregation
2411 and sexual reproduction, and confirms that these isolated were reproduced
2412 clonally. Given that AcEm2 and AcNc2 were sampled 100 miles apart geo-
2413 graphically, and ten years apart in time, clonal reproduction appears to be
2414 the principal mode of reproduction of this race of agronomically important
2415 pathogens.

2416 The largest contig of the reference assembly, (contig 1; 400kb) was
2417 used to analyse polymorphism distribution and detect recombination blocks.
2418 The proportion of heterozygous sites in contig 1 was calculated for each
2419 isolate. Very few sites of contig 1 were heterozygous within AcNc2 (0.03%),
2420 AcEm2, and Ac2V (0.01%). Within isolates AcBoT and AcBoL, the pro-
2421 portion of heterozygous sites was higher (both 0.65%). The high levels of
2422 genotypic identity observed between isolates which infect the same the
2423 host species would not be expected if sexual reproduction and Mendelian
2424 segregation was the primary mode of reproduction, especially given that
2425 isolates AcEm2 and AcNc2 are separated by approximately 100 miles and
2426 10 years. Furthermore, the high proportion of heterozygous sites (for contig
2427 1) in isolates AcBoT and AcBoL is more consistent with asexual population
2428 expansion: A diploid organism reproducing asexually/clonally most of the
2429 time will accumulate mutations between each pair of homologous chromo-
2430 somes. This will generate more heterozygous sites over time, resulting in

2431 allelic divergence and increased observed heterozygosity. However, the
2432 observation of a low level of observed heterozygosity in AcEm2 and AcNc2
2433 is not expected in organisms where asexual and clonal reproduction is
2434 the primary method of reproduction. Given there is no evidence of self-
2435 fertilisation (or any other form of asexual reproduction), it is likely that gene
2436 conversion has been operating to reduce within genome diversity in the
2437 races over time. The phenomenon is called Loss of Heterozygosity or
2438 LOH, and it has been observed in other plant pathogen species such as
2439 *Phytophthora capsici* (Lamour et al. 2012), as well as at a whole genome
2440 scale in yeast (Diogo et al. 2009). In both studies it was hypothesized the
2441 Loss of Heterozygosity observed has facilitated rapid adaptive evolution
2442 and genome plasticity.

2443 To summarise, it appears that the generalist pathogen *Albugo candida*
2444 is comprised of distinct physiological races, which are diverging as they spe-
2445 cialise on different host species. Secondary contact between distinct races
2446 on an immunosuppressed hosts results in inter-specific sexual reproduction
2447 between races, producing new hybrid offspring. These hybrids may be able
2448 to spread rapidly by clonal reproduction on their own, or introgression may
2449 occur.

2450 3.4.2 Biology of genetic introgression and hybridisation

2451 Introgression is defined as the transfer of genetic information (DNA or
2452 RNA) from one species (or OTU, race, or biotype) to another as a result
2453 of hybridization between them followed by repeated backcrossing (Ridley
2454 2004; Abbott et al. 2013).

2455 Hybridisation and introgression can lead to a mosaic-like genome struc-
2456 ture, with regions of different parental lineages interspersed throughout
2457 the genome (Baack and Rieseberg 2007; Stukenbrock et al. 2012). Those

2458 regions will have different ancestry or coalescence, and hence, be rep-
2459 resented by different phylogenetic trees. Introgression has the potential
2460 to augment the adaptive evolutionary potential of populations and intro-
2461 duce a source of genetic variation into genomes. As a source of genetic
2462 variation, mutations have longer waiting times, and lower initial frequen-
2463 cies. In contrast, introgression can occur multiple times, thereby increasing
2464 the probability of fixation of the variant. Furthermore, whereas mutations
2465 tend to be neutral (Kimura 1968), or have (slightly) deleterious fitness ef-
2466 fects (Ohta 1973), introgression inserts pre-selected variation of one of
2467 the parental (donor) lineages into the hybrid line (Hedrick 2013). Adaptive
2468 introgressed variants can be new, have less pleiotropy, less strong linked
2469 effects, and less recessivity (Hedrick 2013). In contrast to mutation, multiple
2470 simultaneous changes across multiple loci are possible with hybridisation
2471 and introgression, but whether these multiple changes are deleterious or
2472 not depends on the details of the molecular interactions within the hybrid.

2473 The view of Wright is that selection favours favorably interacting gene
2474 combinations, resulting in a highly integrated genome which contains coad-
2475 apted gene complexes (Wright 1931; Wright 1932; Dobzhansky 1970).
2476 However, Fisher argued that selection acts on individual genes, and would
2477 favour genes which increase fitness on average across all possible genetic
2478 backgrounds of a given lineage, such genes were called "good mixers"
2479 (Fisher 1930). Both of these views are compatible with the concept of
2480 negative epistasis (Hedrick 2013; Burke and Arnold 2001) in a hybrid ge-
2481 netic background (also called hybrid incompatibility): In any two separated
2482 lineages, fixation of alleles in one lineage occurs independently and there is
2483 no selection for compatibility with any other lineage. Hybridisation produces
2484 novel genotypes which have not previously been subject to selection, and if
2485 they are less well adapted than the parental genotypes, selection would act
2486 against such less fit hybrids. This reduction in fitness of segregating hybrids

2487 has been taken as evidence for unfavorable interactions between genomes
2488 of parental individuals, negative epistasis, and hybrid incompatibility. The
2489 most widely accepted model of such incompatibility was developed by Bate-
2490 son, Dobzansky and Muller (Dobzhansky 1936; Muller 1942). Negative
2491 epistasis has been confirmed empirically in several animal and plant organ-
2492 isms in the past, including (but not limited to) *Drosophila spp.* (True, Weir,
2493 and Laurie 1996; Palopoli and Wu 1994; Hollocher and Wu 1996; Cabot
2494 et al. 1994), *Helianthus spp.* (Rieseberg et al. 1996; Rieseberg, Whitton,
2495 and Gardner 1999), *Tigriopus californicus* (Burton 1990b; Burton 1990a;
2496 Burton, Rawson, and Edmands 1999), and *Iris spp.* (Cruzan and Arnold
2497 1994; Burke, Voss, and Arnold 1998), and is a primary cause of hybrid
2498 inferiority.

2499 However, hybrids can be superior to their parental lineages. Hybrid fit-
2500 ness can occur by several means. F1 hybrids are commonly larger in body
2501 size and have higher growth rates and yields (Baack and Rieseberg 2007;
2502 Hedrick 2013; Burke and Arnold 2001). Such vigour is called heterosis,
2503 and is explained by the dominance and the over-dominance hypotheses
2504 (Baack and Rieseberg 2007; Lippman and Zamir 2007). Other explanations
2505 posit that synergistic interactions between different alleles at different loci
2506 (i.e. positive epistasis and inheritance of complete co-adapted linkage
2507 blocks), and changes in gene expression can also contribute to heterosis
2508 (Baack and Rieseberg 2007; Swanson-Wagner et al. 2006). Heterosis may
2509 contribute towards the establishment of an asexual or allopolyploid hybrid.
2510 Fitness resulting from Heterosis may be short lived, for introgressed hybrid
2511 lineages. This is because sexual reproduction over several generations
2512 would cause loss of heterozygosity in the subsequent (backcrossed) gener-
2513 ations. Instead, long term success depends largely on the fixation of novel
2514 favorable gene combinations from the two parents (Baack and Rieseberg
2515 2007; Burke and Arnold 2001). The genes in such combinations must either

2516 interact favorably with other genes in the combination to increase fitness,
2517 or increase fitness in an additive way, with little or no interaction. Thus,
2518 selection and niche differentiation play a central role in the establishment of
2519 these relatively fit hybrids, because otherwise competition and gene flow
2520 with parental populations may overwhelm them (Buerkle et al. 2000; Riese-
2521 berg, Archer, and Wayne 1999). Just as evidence of negative epistasis has
2522 been found empirically in several species, empirical evidence of epistasis
2523 producing relatively fit hybrids has also been found for several species. For
2524 example, in addition to confirming cases of hybrid inferiority in *Helianthus*
2525 *spp.*, Rieseberg and colleagues also found beneficial epistatic interactions
2526 in hybrid of *Helianthus annuus* and *Helianthus petiolaris* (Gardner et al.
2527 2000; Rieseberg et al. 1996). Evidence of favorable cytonuclear interactions
2528 was found in hybrids of *Iris fulva* and *Iris brevicaulis*, indicating that as well
2529 as interactions between genes, interactions between the nucleus and the
2530 cytoplasm can also determine the success of a hybrid (Burke, Voss, and
2531 Arnold 1998). Hybrid lineages may also exhibit transgressive segregation
2532 i.e. they may have more extreme trait values than either of the parents,
2533 when the parents possessed alleles of opposing effects. This may be bene-
2534 ficial or deleterious, depending on the nature of the trait and may be caused
2535 by epistasis, or, as QTL analyses have demonstrated, through additive
2536 effects (Baack and Rieseberg 2007; Burke and Arnold 2001). Hybridisation
2537 could also help purge mutational load by the masking deleterious alleles
2538 in heterotic F1 individuals, followed by introgression of favorable alleles
2539 (Ingvarsson and Whitlock 2000).

2540 **3.4.3 Introgression and evolution of *Albugo candida* in** 2541 **the wider context**

2542 Given the potential advantages of introgression, it has been hypothesised
2543 that introgression it is instrumental in generating novel combinations of pre-
2544 selected virulence effectors from different diverged races in *Albugo candida*
2545 (McMullan et al. 2015). Not all such combinations may be successful
2546 or viable, but successful genotypes would be important in facilitating the
2547 colonisation of new hosts i.e. a host jump. As a hypothetical example, the
2548 *Albugo candida* race Ac2V is proposed to possess an effector allele, which
2549 interacts with an *Arabidopsis* R gene called WRR4. This prevents Ac2V
2550 from colonising *Arabidopsis*. It is unknown which effector interacts with
2551 WRR4, but if the effector allele segregated away in hybrid offspring, or was
2552 removed through loss of heterozygosity, the hybrid offspring may be able to
2553 overcome *Arabidopsis* resistance.

2554 The impact of introgression and hybridisation has been demonstrated
2555 in other species. For example, in sunflower species *Helianthus anomalus*
2556 (Ungerer et al. 1998). *Helianthus anomalus*, like *Albugo candida*, has a
2557 genome which appears to be composed of distinct parental blocks. How-
2558 ever, unlike *Albugo candida*, the introgression was dated as occurring over
2559 a short timespan of 10 - 60 generations, which provides support for the
2560 idea that hybrid speciation is a punctuated process (Ungerer et al. 1998).
2561 The dating analysis of blocks present in *Albugo candida* suggests that
2562 introgression has occurred between different races at different times, and
2563 repeatedly throughout the evolution of the species. Furthermore, unlike *Al-*
2564 *bugo candida*, introgression in the sunflower species occurred between two
2565 different species, and resulted in a new hybrid species. For *Albugo candida*,
2566 whilst the races are isolated from each other most of the time, repeated in-
2567 trogression between them during secondary contact on immunosuppressed

2568 host plants likely acts to prevent them becoming completely isolated, new
2569 species. A classic example of an adaptive radiation is Darwin's Finches
2570 (*Geospiza*, *Certhidea*, *Pinaroloxias*, and *Camarhynchus/Platyspiza* spp.),
2571 and even here hybridisation has been demonstrated (Lamichhane et al.
2572 2015): Recent whole-genome resequencing, and phylogenetic analysis
2573 based on autosomal, mtDNA, and sex-linked loci of 120 birds representing
2574 all of the Darwin finch species and two other related species revealed dis-
2575 cordant phylogenies (Lamichhane et al. 2015). Calculations of Patterson's
2576 D, supported the hypothesis of gene flow and hybridisation throughout the
2577 radiation (Lamichhane et al. 2015). Rare introgression is thought to have
2578 facilitated the exchange of mimicry genes between *Heliconius* butterfly
2579 species, post isolation (Martin et al. 2013).

2580 Studies from hybridisation with yeast provide findings which corroborate
2581 the findings of this study. For example, genetic exchange between 3 strains
2582 of *Saccharomyces cerevisiae* has been quantified, and indicates that for
2583 these strains out-crossing has only occurred 314 times during approxi-
2584 mately 16 million cell cycles (Ruderfer et al. 2006). This is approximately
2585 one out-crossing event per 50,000 cell cycles. Thus while the strains of
2586 yeast do mate and recombine in the wild, this is not a frequent occurrence
2587 (Ruderfer et al. 2006). This is also what has been inferred for *Albugo*
2588 *candida* as the result of this study. In addition, the genomes of wine strains
2589 of *Saccharomyces cerevisiae* contain introgressed blocks from the species
2590 *Saccharomyces paradoxus*, *Saccharomyces kudriavzevii kudriavzevii*, *Sac-*
2591 *charomyces uvarum uvarum*, and *Zygosaccharomyces bailii* (Dujon 2010).
2592 The blocks in the genome of *Saccharomyces cerevisiae* are almost iden-
2593 tical to the corresponding regions in the genomes of the donor species,
2594 indicating that the introgression events have been recent (Dujon 2010).
2595 This is similar to what this study has demonstrated for *Albugo candida*. It
2596 appears that introgression is a general phenomenon in yeast genomes,

2597 but one review concluded that its importance in its evolution has yet to be
2598 determined.

2599 The importance of introgression in the evolution of *Albugo candida* is
2600 hypothesized to be as follows: Isolation, divergence and specialisation of
2601 races will generate repertoires of tried and tested effectors for a specific
2602 race. Those adapted race-specific repertoires are then brought together
2603 when two races hybridize to generate novel repertoires of novel combi-
2604 nations of these effectors. Specific avirulence effectors that trigger host
2605 immunity may be lost through segregation and the Loss of Heterozygos-
2606 ity (LOH) effect hypothesized to be taken place in oomycetes by Lamour
2607 et al. 2012, and documented here and in McMullan et al. 2015. These
2608 hybrids, having new combinations of effectors, and having lost effectors
2609 which impeded their colonisation of other hosts previously, may expand
2610 their geographical range and population size clonally. Such new hybrids
2611 may be able to colonise new hosts, explaining the phenomenal host range
2612 of species such as *Albugo candida* (and possibly other generalists). Hy-
2613 bridisation between races has been shown to expand host range in other
2614 plant pathogen species such as *Phytophthora* spp. (Ersek, English, and
2615 Schoelz 1995), and the transfer of virulence genes leading to host range
2616 expansion has also been demonstrated in bacterial and fungal pathogens
2617 (Ford Doolittle 1999; Mehrabi et al. 2011). Sexual oospores of *Albugo*
2618 *candida* are tolerant of strong environmental pressures, which raises the
2619 prospect: might hybrid spores produced by reproduction between two races
2620 lie dormant, forming banks of hybrid genotypes, waiting for conditions better
2621 suited to their genotype and phenotype?

2622 The ability to expand host range and generate novel genotypes through
2623 hybridisation, and then reproduce rapidly clonally may be especially fa-
2624 vored in a monoculture based agro-ecological environment, characterized
2625 by different, large, homogeneous regions of (often clonal) host plants of

2626 one species (Stukenbrock and Bataillon 2012). Recently Stukenbrock et al.
2627 2012 demonstrated that the plant pathogen species *Zymoseptoria pseu-*
2628 *dotritici* was formed by the hybridisation of two distinct fungal individuals,
2629 and that the genome is characteristic of bottleneck and selection following
2630 the hybridisation event which occurred approximately 380 sexual genera-
2631 tions ago, resulting in the generalist grass pathogen. The obligate biotroph
2632 and powdery mildew, *Blumeria graminis f. sp. Hordei* also has a mo-
2633 saic genome of alternating monomorphic and polymorphic DNA sequence
2634 blocks (Hacquard et al. 2013). Pathogen adaptation to agro-ecological
2635 environments is characterized by high genome plasticity of pathogens (a
2636 successful pathogen needs to keep up in the co-evolutionary arms race
2637 with its host), but a reduction in diversity for recently emerged lineages
2638 (selection is strong and new and recently emerged lineages are often bottle-
2639 necked) (Stukenbrock and Bataillon 2012). Pathogens such as late blight of
2640 potato, *Phytophthora infestans*, wheat yellow rust *Puccinia striiformis*, and
2641 *Magnaporthe oryzae*, which are specialised, may represent an end-result
2642 of a much broader process of pathogen adaptation and evolution. The
2643 results gained from this work provide insight into how recombination and
2644 hybridisation plays a role in generating novel virulent races, and into their
2645 subsequent spread and geographical range expansion by clonal propaga-
2646 tion. These findings are of particular relevance to modern, monoculture
2647 based agriculture.

2648 CHAPTER 4

2649 Allelic divergence in the polar diatom 2650 *Fragilariopsis cylindrus*

2651 This chapter is based on a submitted scientific paper:

2652 *Mock, T., Otilar, R. P., Strauss, J., Allen, A. E., Dupont, C. L., Fricken-*
2653 *haus, S., ... Grigoriev, I. V. (Submitted). Extensive genetic diversity and*
2654 *differential bi-allelic expression in a Southern Ocean diatom. Nature.*

2655 This project was a very large collaboration spanning many years to
2656 sequence the genome of the *Fragilariopsis cylindrus* organism. In order
2657 to clearly describe my work and set it in context, some work that was not
2658 performed by myself is described. In particular, any work mentioned in
2659 the introduction is not my contribution to the work, but was completed by
2660 colleagues. My contributions to the work are described in sections 4.2.2.1,
2661 4.2.2.2, 4.2.2.3, 4.2.2.4, and the results section presents data that was the
2662 outcome of my work only. In the discussion some further preliminary work
2663 is described. A figure showing this work is provided as an appendix, and
2664 this work was done jointly and equally between myself and a colleague.

4.1 Introduction

4.1.1 Sexual reproduction and recombination

Sex as a mode of reproduction has a two-fold cost. Firstly, most sexually reproducing species only have one gender capable of bearing offspring (Visser and Elena 2007). Secondly, in sexually reproducing organisms, any individual will only contribute approximately half of its genetic information to each offspring; i.e. in diploid sexuals, gametes are haploid (Agrawal 2001). In contrast, an asexually reproducing, clonal organism contributes all of its genetic information to each offspring, and every individual is typically capable of bearing young (Schlupp, Taebel-Hellwig, and Tobler 2010). This generalization applies to most sexual organisms however, there are exceptions. For example, not all sexually producing organisms have the two-fold cost problem. Yeasts are sexual organisms with two mating types and both types are capable of producing offspring. In addition, a species of poecilids can reproduce through a process of gynogenesis; a process similar to asexual reproduction through parthenogenesis, but is distinct as the presence of sperm is required to stimulate egg development (Schlupp, Taebel-Hellwig, and Tobler 2010). Hybridisation has also given rise to a Hermaphroditic Cichlid individual which can self (Svensson et al. 2016). In addition, some species shuttle between asexual and sexual reproduction, and the frequency at which this happens directly affects the factors raised above.

All else being equal, an asexual species should outperform a sexual species over time because of its faster population growth rate. However, sexual and asexual species do co-exist together, sometimes with similar fecundity (Schlupp, Taebel-Hellwig, and Tobler 2010). However, despite this, sexual reproduction is very widespread, especially among the eukaryotes.

2692 These observations led researchers to think that the benefits of sexual
2693 reproductions must be evolutionary and lead to the production of offspring
2694 with benefits that outweigh to costs. To summarize most of the commonly
2695 cited reasons sexual reproduction is maintained, it may be described as a
2696 mechanism, through which:

- 2697 1. Beneficial mutations can spread through a population more quickly.
- 2698 2. Novel genetic combinations are generated.
- 2699 3. Deleterious mutations can be purged or masked.

2700 These benefits are possible because sexual reproduction brings to-
2701 gether into one individual, the chromatids (and alleles they contain) in
2702 the gametes of two parental individuals from separate genealogical lines
2703 (out-crossing). In addition, when parental individuals generate gametes,
2704 meiotic recombination will result in new combinations of genes (Felsenstein
2705 and Yokoyama 1976). This in turn contributes to the generation of novel
2706 genetic (or rather, genotypic) variation. As a result, two or more beneficial
2707 mutations from separate genealogical lines may occur together within the
2708 same individual, thus facilitating the spread of beneficial mutations through
2709 the population to fixation.

2710 This is formalized by the Hill-Robertson effect (Hill and Robertson 1966),
2711 and is demonstrated by considering two loci with the haplotype A_2B_2 with
2712 a fitness of 1. It is then assumed two mutants at both loci (A_1B_2 , A_2B_1)
2713 can occur after a time period with fitnesses of $1 + s$, and that fitnesses are
2714 multiplicative such that A_1B_1 has fitness $(1 + s)^2$. With no or low recombina-
2715 tion, the ancestral haplotype is lost by selection and both advantageous
2716 mutants will exist in the population for some time until one is lost by drift
2717 (Coop and Przeworski 2007). But with recombination, a haplotype A_1B_1 is
2718 possible, bringing both mutants together in one haplotype before one of the

2719 mutants is lost by drift, thus both mutants get fixed rather than one Coop
2720 and Przeworski 2007. With low recombination rates selection increasing the
2721 frequency of the mutant alleles is less effective, this is the Hill-Robertson
2722 effect (Hill and Robertson 1966).

2723 The effect is more likely to occur when selection is not too strong, re-
2724 combination rates are low, and when the favorable mutants have negative
2725 disequilibrium i.e. they initially occur on different haplotypes (Hedrick 2010).
2726 An asexual lineage, in contrast would have to acquire one beneficial muta-
2727 tion, followed by another, a limitation called clonal interference (Gerrish and
2728 Lenski 1998).

2729 Similarly, deleterious mutations accumulating throughout the population
2730 in different genealogical lines may occur together within one individual,
2731 which suffers stronger negative selection pressure and is eliminated from
2732 the population (Crow 1994). A third possibility is a deleterious allele is
2733 inherited from one parent, and the corresponding allele inherited from the
2734 second parent is not deleterious. In that case, the affects of the delete-
2735 rious allele may be alleviated or masked, as the offspring individual still
2736 possesses a non-deleterious copy. Chromosomal crossover during meiosis
2737 may also result in the removal of deleterious mutations (Crow 1994).

2738 The maintenance of sexual reproduction has also been attributed to its
2739 role in DNA mismatch repair (Bernstein, Bernstein, and Michod 2011). The
2740 repair and complementation hypothesis proposes that sexual reproduction
2741 is an adaptive response to incorrect DNA replication, through mutation
2742 and damage to the DNA molecule (Bernstein et al. 1984; Bernstein 1985;
2743 Bernstein, Hopf, and Michod 1987). Recombination repair is the only
2744 mechanism currently known which removes double stranded damages to
2745 the DNA molecule and such double strand damage is common and could
2746 be lethal if not repaired: in human cells such damage occurs approximately
2747 50 times per cell cycle (Vilenchik and Knudson 2003).

2748 Recombination and sexual reproduction also plays a role in eliminating
2749 detrimental variation from the population, which otherwise would accumu-
2750 late over time and decrease the fitness of the population (Muller's ratchet)
2751 (Muller 1932). Recombination produces individuals containing fewer delete-
2752 rious mutants, helping to reverse the decline in fitness.

2753 The Red Queen Hypothesis also offers an explanation as to why sex
2754 has repeatedly evolved in all life forms (Paterson et al. 2010). It states that
2755 in a rapidly changing environment, alleles that were previously neutral or
2756 deleterious and the rapid change makes sexual reproduction advantageous.
2757 Such rapid changes are proposed to be particularly evident during co-
2758 evolution between a parasite and its host (Decaestecker et al. 2007).

2759 However, despite the advantages of sex, evidence of ancient asexuality
2760 has been identified in the genomes of some organisms including root-
2761 knot nematodes and bdelloid rotifers (Lunt 2008; Welch and Meselson
2762 2000; Meselson and Welch 2007; Pouchkina-Stantcheva et al. 2007). The
2763 classic hallmarks of ancient asexuality are diverged alleles and a lack of
2764 phylogenetic incongruence caused by recombination (Schurko, Neiman,
2765 and Logsdon 2009).

2766 **4.1.2 *Fragilariopsis cylindrus* and Diatoms**

2767 *Fragilariopsis cylindrus* is a species of Diatom: microscopic eukaryotic
2768 phytoplanktons, which are found throughout all the worlds oceans wher-
2769 ever there is sufficient light and nutrients to support them (Armbrust 2009).
2770 Diatoms are so named because of their shape and method of reproduc-
2771 tion: Their cells are covered by a silica cell wall made of two halves, and
2772 they reproduce by asexual mitotic division, decreasing in size each time.
2773 Diatoms occasionally reproduce by forming an auxospore, which reverses
2774 the decline in size resulting from reproduction by mitotic division (Armbrust
2775 2009). Auxospores also play a role in sexual reproduction, forming after

2776 haploid gametes fuse to form a diploid zygote. Diatoms are an important
2777 group of organisms of study because of their role in the ecosystem and in
2778 marine biogeochemical cycles (Assmy et al. 2013; Thomas and Dieckmann
2779 2002; Pondaven et al. 2000).

2780 Diatoms provide an important ecosystem service by performing photo-
2781 synthesis. It has been estimated that of all photosynthesis that occurs on
2782 earth, one fifth is performed by Diatom species. Each year diatoms gener-
2783 ate as much organic carbon as that produced in total by all the terrestrial
2784 rainforests on Earth (Armbrust 2009). The organic carbon that is produced
2785 by diatoms by photosynthesis is input into food webs: in coastal regions
2786 diatoms support fisheries (such as anchovies in the Peruvian ocean) and in
2787 the open-ocean, much of the organic matter produced sinks and becomes
2788 food for deep-sea organisms (unless it reaches the ocean floor, where it
2789 may become sequestered in sediment and rock) (Armbrust 2009; Bowler,
2790 Vardi, and Allen 2010). As a result, a significant amount of petroleum
2791 deposits under the ocean floor are derived from diatoms sinking.

2792 As Diatoms are found throughout all the world's oceans, they popu-
2793 late interesting and dynamic environments in which environmental factors
2794 change and can become extreme. They are known to be adapted to limited
2795 iron, extremes in temperature (Arrigo et al. 2012; Bayer-Giraldi et al. 2011;
2796 Bowler, Vardi, and Allen 2010), salinity (Krell 2006), and temporal variation
2797 in the environment: seasons cause rises and falls in temperature, and
2798 freezing and melting sea ice also means the environment's structure can be
2799 heterogeneous through time. All these extremes occur in the environment
2800 of *Fragilariopsis cylindrus*, which is particularly successful in the Southern
2801 Ocean, and is often found to form large populations in the bottom layer of
2802 sea ice and the wider sea-ice zone including open waters (Kang and Fryxell
2803 1992). Such ice is characterized by temperatures below the freezing point
2804 of sea water, high salinity caused by the semi-enclosed pores within the ice,

2805 and low diffusion rates of dissolved gases and inorganic nutrients (Thomas
2806 and Dieckmann 2002). The environment is not limited in dissolved iron
2807 however, unlike the surface ocean (Wang et al. 2014). Furthermore, the
2808 environment is dynamic: every winter, phytoplankton in the Southern Ocean
2809 get locked into sea ice and are released again in the following summer,
2810 when most of the sea ice melts (Vancoppenolle et al. 2013). However, only
2811 a subset of these phytoplankton them have evolved adaptations to cope
2812 with this dramatic environmental change, including *F. cylindrus*, which is
2813 known to thrive in both habitats (Bayer-Giraldi et al. 2011; Vancoppenolle
2814 et al. 2013).

2815 How Diatoms have adapted to such conditions, and become so suc-
2816 cessful in the oceans, is of interest to evolutionary biologists and genome
2817 sequencing has provided insight. Complete genome sequences are avail-
2818 able for two Diatom species (*Thalassiosira pseudonana* and *Phaeodactylum*
2819 *tricornutum*), containing between 10 and 14 thousand genes. However, of
2820 those genes only approximately half can be assigned a putative function
2821 based on experimental knowledge (Bowler, Vardi, and Allen 2010). Further-
2822 more, approximately 35% of the genes found are specific to each Diatom,
2823 which suggests some of them encode adaptations to specific environmental
2824 conditions (Bowler, Vardi, and Allen 2010). As secondary genome se-
2825 quences became available, the origin of Diatoms seems to be a secondary
2826 endosymbiosis between red algae and a heterotrophic eukaryote, and sur-
2827 prisingly many bacterial genes were identified, highlighting the role of HGT
2828 in the evolution of Diatom species (Bowler, Vardi, and Allen 2010; Raymond
2829 and Kim 2012).

2830 Diatom specific genes were found to have high diversification rates, and
2831 since *Thalassiosira pseudonana* and *Phaeodactylum tricornutum* diverged
2832 approximately 90 million years ago, and the two have diverged as much as
2833 metazoans had diverged in approximately 550 million years (Bowler, Vardi,

2834 and Allen 2010). It is thought that diversification in Diatoms has been driven
2835 by transposable elements, which increased the rate of insertion, deletion,
2836 and recombination events (Bowler, Vardi, and Allen 2010). In contrast,
2837 diversification of genes in metazoan genomes during the aforementioned
2838 550 million years, is thought to have occurred largely through whole and
2839 segmental gene duplication events (Bowler, Vardi, and Allen 2010). Some
2840 of the diatom specific transposons are activated in response to stresses
2841 such as Nitrogen starvation, suggesting diversification of Diatom genes
2842 may be stimulated by environmental cues (Bowler, Vardi, and Allen 2010).
2843 The resulting mix-and-match genomes (Armbrust 2009) of Diatom species
2844 has brought together unique combinations of genes facilitating adaptation
2845 to a range of environments, including that encode unique pathways of
2846 nutrient assimilation. comparing the genome of a psychrophile such as
2847 *F. cylindrus* with that of diatoms evolved in temperate oceans provides
2848 an opportunity to obtain first insights into how this species has evolved to
2849 conditions of Southern Ocean waters, and managed to persist for millions
2850 of years, underpinning the ecology of an unique food web.

2851 Recently the first large-scale genomic sequencing of *Fragilariopsis*
2852 *cylindrus*, a eukaryotic psychrophilic organism of ecological importance,
2853 including whole-genome sequence, transcriptome and population genetic
2854 analyses, was completed. In this thesis chapter I present my contribution
2855 to the population genetic analyses of this large body of collaborative work.
2856 This goal of the work described in this chapter was conducted in order
2857 to evaluate hypotheses about the evolutionary history of *Fragilariopsis*
2858 *cylindrus*. These hypotheses were proposed during the genome project,
2859 to explain observations about the genome data, and the hypotheses that I
2860 tested in this project.

2861 **4.1.3 The *Fragilariopsis cylindrus* genome project**

2862 The draft of the *F. cylindrus* genome was approximately 60Mb in length,
2863 which is larger than the sequences for the nuclear, plastid and mitochondrial
2864 genomes of the cosmopolitan diatom *T. pseudonana* (34Mb), and the whole-
2865 genome sequence of *P. tricornutum* (27Mb) (Armbrust 2009; Mock et al.
2866 0). The draft genome of *F. cylindrus* is smaller in size compared to the
2867 toxigenic coastal species *Pseudo-nitzschia multiseriis* (300 Mb) (Armbrust
2868 2009).

2869 Assembler programs typically use single end or paired end reads to find
2870 overlaps in sequence fragments, joining them to form contigs. Since it is
2871 known that paired end reads are generated from the same DNA fragment,
2872 this can help link contigs onto scaffolds, which are ordered assemblies of
2873 contigs, with gaps in between them (Baker 2012). However, assemblers
2874 are not always accurate: one common problem is that if one suspects that
2875 the read depth for an assembled region is too high, then it may be that the
2876 assembler has merged multiple regions because of their high sequence
2877 similarity (typically these are repeat rich regions or duplications) (Baker
2878 2012). A second problem is if one suspects that regions of an assembly
2879 have a lower read depth than the rest of the assembly, then it may be
2880 that those regions represent single polymorphic loci, which have been
2881 assembled as two distinct loci (Baker 2012). 30.2Mb of the scaffolds of *F.*
2882 *cylindrus* could not be collapsed into a single haplotype, because they had
2883 greater than 1.5% nucleotide discrepancies. The genome contains just over
2884 20,000 protein-encoding genes, and of those, 28% of them represent alleles
2885 that could not be collapsed (Mock et al. 2017). The genome contains 46 Mb
2886 of collapsed haplotype and 15.1 Mb of diverged haplotype that represents
2887 the diverged alleles of the same genetic loci.

2888 The genome contains 21,066 predicted protein-encoding genes, 6,071

2889 genes were represented by diverged alleles, and each pair of diverged
2890 alleles had both coding and non-coding regions, and were up to 6% poly-
2891 morphic in the non-coding regions. Comparison of the diverged allele, and
2892 non-diverged allele gene ontologies (GO) revealed that genes in the cate-
2893 gories catalytic activity (GO:0003824), transporter activity (GO:0005215),
2894 metabolic process (GO:0008152), transport (GO:0006810) and integral to
2895 membrane (GO:0016021) were significantly enriched in the diverged alleles
2896 set (Mock et al. 2017). Furthermore, biological process GO categories
2897 metabolic process (summarising lipid-catabolic process (GO:0016042),
2898 glucose metabolic process (GO:0006006), oxidation-reduction process
2899 (GO:0055114) and translation (GO:0006412)) as well as GO category
2900 transport-related categories protein transport (GO:0015031) and proton
2901 transport (GO:0015992) enriched in metatranscriptome sequences from
2902 Southern Ocean sea ice, and these sequences had high similarity to se-
2903 quences contained in the diverged alleles of *F. cylindrus* according to
2904 BLASTX analyses (Mock et al. 2017).

2905 Differential expression experiments and RNA-Sequencing suggested
2906 that 40% of the non-collapsed, diverged allelic pairs showed a 4 fold unequal
2907 bi-allelic expression (Mock et al. 2017). This suggested an allele-based
2908 adaptation to different environmental conditions. The differential expression
2909 in alleles suggested they were controlled by separate regulatory systems.
2910 Alleles showing the strong unequal bi-allelic expression were found to have
2911 an elevated rate of non-synonymous mutations, which suggests significant
2912 positive / adaptive selection and evolution of these allelic pairs (Mock et al.
2913 2017). It was concluded therefore, that positive selection has been a driving
2914 force in the evolution of these alleles and hence the adaptation of this
2915 diatom to the environmental conditions it faces.

2916 An evolutionary explanation of the 28% of genes that could not be
2917 collapsed (i.e. diverged genes) is desired, as it would explain one of the

2918 mechanisms through which this diatom appears to have adapted to its
2919 polar environment. However, this signature of positive selection alone does
2920 not provide a sufficient evolutionary explanation: Meiotic recombination,
2921 which occurs during sexual reproduction, should act to homogenize any
2922 two alleles of one gene in the diatom genome.

2923 Allelic divergence is a classic signature in genomes of organisms called
2924 ancient asexuals (Little and Hebert 1996; Pouchkina-Stantcheva et al.
2925 2007; Schurko, Neiman, and Logsdon 2009). By its definition asexuality is
2926 a negative proposition, based on an apparent lack of sexual reproduction in
2927 an organism, and since absence of evidence is not equivalent to evidence
2928 of absence, ancient asexuality is a difficult proposition to demonstrate in an
2929 organism absolutely (Schurko, Neiman, and Logsdon 2009). Indeed the
2930 existence of ancient asexuals has been debated and doubted in the past
2931 (Judson and Normark 1996; Little and Hebert 1996), and this is perhaps
2932 unsurprising considering current theory explaining the benefits of, and
2933 maintenance of sexual reproduction.

2934 If the divergence of alleles is due to ancient asexual reproduction, then
2935 the recombination rate between these alleles should be reduced. It was
2936 also expected that phylogenetic networks would have a very clear structure,
2937 with deep branches. To test these predictions and evaluate empirical
2938 data I performed population genetic simulations. More detail is presented
2939 in the methods section, but briefly, sequence data was available to test
2940 for the evidence of recombination based on an environmental sample of
2941 *F. cylindrus*, that was amplified by PCR and sequenced using Sanger
2942 sequencing. It resulted in 200 high quality sequences from alleles of
2943 Ferrichrome ABC transporter and Large Ribosomal Protein L10, and the
2944 signature of recombination between these alleles was analyzed as well as
2945 several other population genetic parameters.

2946 This project had the aim of establishing whether ancient asexuality and

2947 a lack of recombination is evident, by establish whether recombination has
2948 occurred by analyzing the aforementioned DNA sequences.

2949 The specific aims were:

- 2950 ● Use LAMARC to establish a population recombination rate and popu-
2951 lation Theta parameter.
- 2952 ● Use the incompatible sites test to detect evidence of phylogenetic
2953 incompatibility (and therefore recombination) between closely related
2954 sequences.
- 2955 ● Visualize recombination signal of choice sequences with the Hybrid-
2956 Check package.
- 2957 ● Conduct a comparative phylogenetic network analysis.
 - 2958 – Construct un-rooted phylogenetic networks of alleles present in
2959 the natural sea-ice populations.
 - 2960 – Construct un-rooted phylogenetic networks from *silico* popula-
2961 tions simulated using simuPOP. Some of these *silico* populations
2962 were simulated under asexual (clonal) regimes of reproduction,
2963 and some were simulated under a sexual reproduction regime,
2964 with different mutation and recombination rates.
 - 2965 – Compare the empirical networks with those simulated, to try
2966 and suggest the mutation and recombination rates the Diatom
2967 population may have in nature.

2968 **4.2 Materials and Methods**

2969 **4.2.1 Materials**

2970 **4.2.1.1 Sequence Data: PCR Amplified Alleles**

2971 In this study, subsequently described analyses were performed using the
2972 same dataset. Two genes (ABC Iron Transporter (Protein ID 240308)
2973 and Large Ribosomal sub-unit (Protein ID 240308)) of an environmental
2974 sample of *F. cylindrus* were amplified by PCR and sequenced using Sanger
2975 sequencing to yield high quality sequences. A total of 93 and 103 alleles
2976 were found in both genes, respectively. The DNA extraction, and PCR
2977 amplification, was completed by Dr. Jan Strauss. Sanger sequencing was
2978 performed by (Mock et al. 2017). These two sequence datasets shall be
2979 referred to hereafter as FcABC (ABC Iron transporter), and FcLR (Large
2980 Ribosomal Subunit).

2981 **4.2.1.2 Sequence Data: Allelic pairs from the genome**

2982 Previously, a set of diverged alleles was defined for any downstream analy-
2983 ses: The genome assembly was aligned against itself using BLAST, with a
2984 95% nucleotide identity threshold, and greater or equal to 50% alignment
2985 coverage for smaller scaffolds. Syntenic scaffolds that were homologous
2986 across their whole length were analyzed with Mauve. Diverged alleles on
2987 large scaffolds were referred to as allele 1, the corresponding allele on the
2988 smaller scaffold was referred to as allele 2. For more details, the reader
2989 is referred to the paper (Mock et al. 2017). The allelic pair set was used
2990 to estimate coalescence times between alleles, as described in the next
2991 section.

4.2.2 Methods

4.2.2.1 Estimating Coalescence times of alleles

Because the FcABC and FcLR sequences were used for recombination detection, and the calculation of networks for the simulation and network analysis portion of this study, it was important to determine the two sequence datasets were representative of the allelic pairs identified in the genome data. Therefore, coalescence times were calculated A) Between the two sequences of each allelic pair identified from the genome data (see above), B) between pairs of FcABC sequences, C) between pairs of FcLR sequences. If the distributions of coalescence times for A) FcABC, and B) FcLR, overlap the distribution of coalescence times calculated for the genome data, then the FcABC and FcLR sequence datasets could be considered representative of the allelic pairs from the genome data.

Coalescence times were estimated using the algorithm available in the HybridCheck R package (<https://github.com/Ward9250/HybridCheck>). The algorithms and design of HybridCheck is described in chapter 2 of this thesis. Briefly, the algorithm used estimates coalescence time of two aligned sequences based on the number of mutations that are observed between two sequences. HybridCheck models a Bernoulli trial with a strict molecular clock, which assumes a constant mutation rate ($\mu = 10e-9$) and a Jukes and Cantor model for base substitutions.

Coalescence time estimates calculated by the HybridCheck algorithm are expressed in terms of generations, as described in chapter 2. An estimate in terms of real time (years) was desired to attempt to put the divergence of the allelic pairs into a historical context. Estimates were converted to years using an estimated division rate of 12.472 per year. This yearly division rate assumed a division rate of 0.1 per day, and a growing season of four months per year, where each month consisted of 30.4368

3020 days. 946 allelic pairs were successfully pulled, aligned, and dated from
3021 the genome sequence data.

3022 **4.2.2.2 Testing for recombination in the PCR amplified alleles with** 3023 **the PHI-test**

3024 We tested for recombination in both the FcABC and FcLR sequence
3025 datasets using the PHI-test for recombination (Bruen, Philippe, and Bryant
3026 2006). The test accepts a multiple sequence alignment and is based on
3027 the principle of refined compatibility: For a given pair of informative sites in
3028 a multiple sequence alignment, they are deemed compatible if there is a
3029 phylogenetic history that can be inferred parsimoniously, on the condition
3030 that there is no recurrent mutation, or convergent mutations (Le Quesne
3031 1969).

3032 If the condition is not satisfied then the sites are classified as incom-
3033 compatible. Incompatible sites are explained either by homoplasies, or by
3034 recombination. The PHI-test extends this notion by using the refined in-
3035 compatibility score, which allows for consideration of situations in which
3036 multiple homoplasies can be parsimoniously inferred a pair of sites (Bruen,
3037 Philippe, and Bryant 2006). The PHI-test then computes the mean refined
3038 compatibility scores of nearby sites and a p-value is calculated parametri-
3039 cally (Bruen, Philippe, and Bryant 2006). The analyses were repeated with
3040 window sizes of 100, 50, and 10 base pairs.

3041 **4.2.2.3 Population recombination rate and theta parameter estima-** 3042 **tion with LAMARC**

3043 A population recombination rate, and the population mutation rate Θ (Theta),
3044 was inferred for the FcABC and FcLR sequence datasets, using the LAMARC
3045 software for coalescent analysis (Kuhner 2006). Five independent runs
3046 were run for both datasets, in which 20 sequences were randomly sampled

3047 from each sequence dataset, and analysed with LAMARC, using uninfor-
3048 mative priors and default settings, as much about *F. cylindrus* populations in
3049 the wild is unknown. These results informed the choice of the Θ parameter
3050 used in simulations as described below.

3051 **4.2.2.4 Comparative Phylogenetic Network Analysis**

3052 *Population Genetic Simulations.*

3053 All simulation scenarios were written as simuPOP scripts (Peng and
3054 Kimmel 2005). Since we are interested in assessing whether *F. cylindrus*
3055 has an asexual past causing allelic divergence, when the word recombina-
3056 tion is used in the section is specifically refers to meiotic recombination
3057 unless otherwise stated.

3058 Two scenarios were simulated:

- 3059 1. A scenario in which individuals reproduced clonally (i.e asexually) and
3060 no recombination could take place.
- 3061 2. A scenario in which individuals reproduced sexually every generation,
3062 and in which the rate of meiotic recombination could be specified.

3063 In all three of these simulations, individuals in the simulated population
3064 were diploid and so contained one pair of chromosomes each (two homol-
3065 ogous copies). The chromosomes were 750bp in length and the pairs of
3066 chromosomes begin as identical. By initializing individuals in this manner
3067 and then evolving them, each individual containing a pair of 750bp acted
3068 as an evolving allelic pair.

3069 When running each simulation design, various combinations of effective
3070 populations size, and mutation rates were used in a balanced manner such
3071 that Θ for the simulated populations should result in a similar Θ estimated for
3072 the FcABC and FcLR sequences by the LAMARC analysis. This permitted
3073 the preservation of the Θ parameter of the population but allowing more

3074 reasonable compute time. Θ values of 0.66, 0.066, 0.0066, were chosen
3075 based on the LAMARC analysis, with the value 0.066 being closest to the
3076 estimates returned by LAMARC.

3077 It was assumed that the census size set in the simulations is a rea-
3078 sonable approximation for the effective population size, given that in our
3079 simulations the population was panmixtic, i.e.:

- 3080 • There are always an equal number of males to females.
- 3081 • No one individual is more likely to produce offspring than any other.
- 3082 • Mating is random when sexual reproduction occurs any male can
3083 potentially be paired with any female.
- 3084 • The number of breeding individuals is always the same for all genera-
3085 tions.

3086 For the simulations where recombination occurs, various recombination
3087 rates (relative to μ) were used, from no recombination ($r = 0$), to $r = 0.1\mu$,
3088 $r = 0.5\mu$, $r = \mu$, $r = 5\mu$, and $r = 10\mu$.

3089 All simulations ran on the computer for a number of generations equal to
3090 the intended effective population size multiplied by 20. The mating scheme
3091 kept the population size constant during mating, one male and one female
3092 virtual diatom is randomly picked from the population. The number of
3093 offspring they produce is drawn from a Poisson distribution with mean and
3094 variance equal to 2. This is repeated over and over until the new offspring
3095 population is of equal size to the parental population. Individuals could be
3096 randomly selected for mating more than once.

3097 In every simulation performed, 96 individuals were randomly sampled
3098 and exported at various time points throughout all the simulation runs, and
3099 converted to FASTA sequence files. These FASTA files could then be used
3100 for generation of networks with SplitsTree (Huson 1998).

3101 *Preparation of PCR amplified allele sequences.*

3102 The population genetic simulations described above were simulated
3103 with the absence of selection pressure. Therefore, before constructing
3104 phylogenetic networks of the FcABC and FcLR sequences to compare with
3105 networks constructed from the simulated sequences, it was necessary to
3106 reduce the influence of selection as much as possible. Therefore, when
3107 constructing phylogenetic networks for the FcABC and FcLR sequences,
3108 only the 3rd codon positions were utilized. To do this, a script translated
3109 every sequence in every possible reading frame and scored the number of
3110 stop codons or unknown proteins present in the translation. It is assumed
3111 the correct reading frame for the alleles is the one in which there are no
3112 stop codon in the middle of the sequence. Furthermore, this reading frame
3113 should be the same for almost all sequences. Sequences that resulted
3114 in uncertain translations in every reading frame were not used, and only
3115 sequences that had showed one reading frame with no stop codons were
3116 used to build networks.

3117 *Calculating Phylogenetic Networks.*

3118 Phylogenetic networks were computed for the FcABC, FcLR, and simu-
3119 lated sequence datasets generated by each of the population genetic simu-
3120 lation scenarios previously described. All networks have been generated
3121 with the SplitsTree software (Huson 1998), and the methods used in the
3122 package to compute and draw the networks were the *Uncorrected_P* char-
3123 acter transform, the *NeighbourNet* distances transform, and the *EqualAngle*
3124 splits transform.

3125 These networks constitute an expectation of what may be seen in the
3126 networks of the *F. cylindrus* alleles under various scenarios of sexuality or
3127 asexuality. If *F. cylindrus* has a past history of asexual reproduction, we
3128 would expect networks of sequences generated by an asexual simulation
3129 to show greater similarity to the networks of the *F. cylindrus* alleles. If *F.*

3130 *cylindrus* has a past history of low levels of sex then its network would show
3131 more similarity to the network derived from the model in which there is lower
3132 levels of recombination, and so on. By comparing the *F. cylindrus* networks
3133 to those modeled networks it is possible to assess whether strict asexuality
3134 or infrequent sex is a likely possibility. It is important to note any simulated
3135 scenario with sexual reproduction with a zero recombination rate is not the
3136 same as asexual reproduction as the clonal reproduction scenario as the
3137 latter does not follow Mendelian inheritance, whereas sexual reproduction,
3138 with a recombination rate of zero, does follow Mendelian inheritance.

3139 In comparing networks of simulated allelic pairs and networks of the
3140 sequenced *F. cylindrus* sequences, characteristics regarding the structure
3141 of the network, can be expressed quantitatively. To quantitatively assess
3142 the networks, we calculated the p-distance matrices for all the sets of
3143 simulated scenario sequences, and for the real *F. cylindrus* sequences.
3144 In particular we calculated the mean and the variance both of which
3145 were expected to be higher for networks of sequences evolved with lower
3146 recombination rates, showing signs of allelic divergence. The distances
3147 reflect the mean branch length in the network and are principally affected by
3148 the mutation-drift equilibrium, and hence Θ . In order to assess the effect of
3149 recombination relative to the mutation rate (R/μ), we quantified the number
3150 splits in the network, again comparing the simulated networks with those of
3151 the *F. cylindrus* alleles.

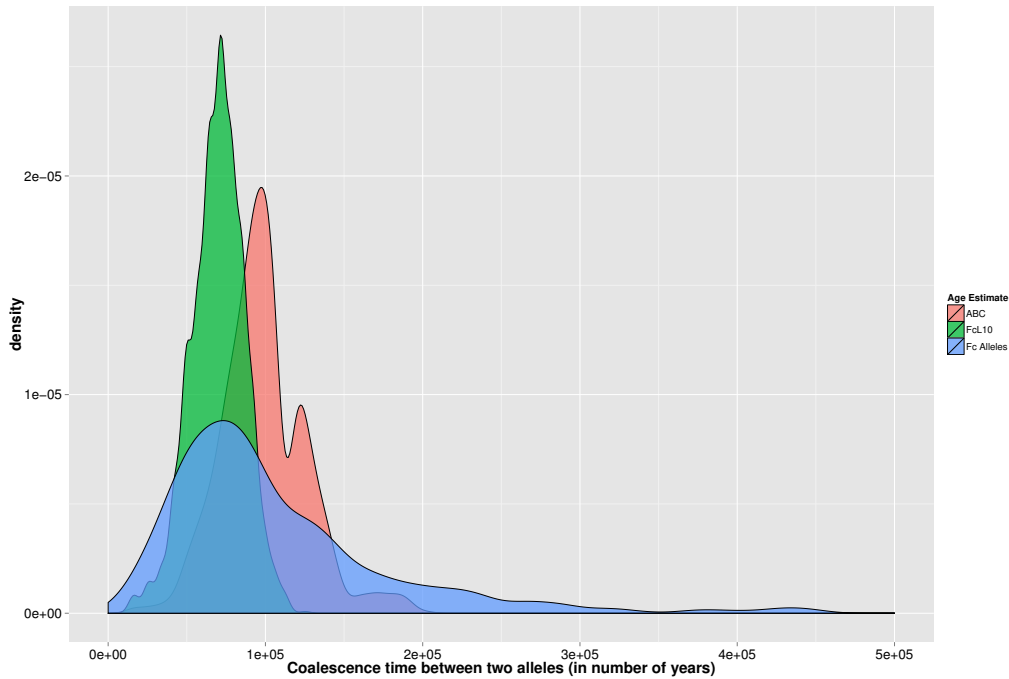


Figure 4.1: Smoothed density plot of the maximum coalescence times (in generations) calculated for allelic pairs of the ABC Iron Transporter (red), Large Ribosomal Subunit (green) and allelic pairs from the genome (blue).

4.3 Results

4.3.1 Estimating coalescence times of alleles

Figure 4.1 shows the distances calculated between the allelic pairs simulated from the ABC Iron Transporter and Large Ribosomal Subunit sequence pools, and between the allelic pairs identified from Fc Alleles RNAseq data. The three distributions show considerable overlap, which implies that the divergence between allelic pairs identified from the genome is representative of the divergence between alleles from two known genes (Figure 4.1).

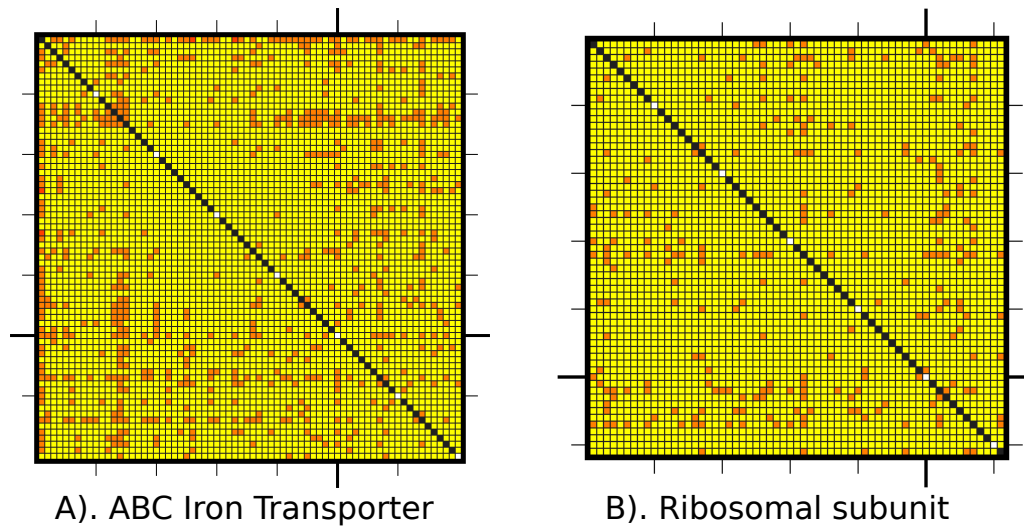


Figure 4.2: Incompatibility score matrices computed for A). The ABC iron Transporter and B). The Large Ribosomal Subunit. Yellow boxes indicate two informative sites are compatible, and darker boxes indicate the two sites are incompatible. The presence of incompatible sites in the alignments is suggestive of recombination.

3161 **4.3.2 Testing for recombination in the PCR amplified al-** 3162 **leles with the PHI-test**

3163 PHI Scores calculated for the sequences of the ABC Iron transporter and
3164 the Large Ribosomal Subunit (Table 4.1), and Figure 4.2 shows the refined
3165 incompatibility matrices between informative sites computed for the ABC
3166 Iron Transporter (A.), and the Large Ribosomal Subunit (B.). Yellow squares
3167 indicate pairs of informative sites that are compatible, darker squares
3168 indicate a pair of sites that are incompatible. The presence of incompatible
3169 sites in these sequences, and the PHI-Scores and NSS scores shown in
3170 Table 4.1 suggests recombination has indeed affected these sequences.

3171 **4.3.2.1 Comparative analysis of phylogenetic networks**

3172 Figure 4.3. Shows an example network generated from sequences pro-
3173 duced by the population genetics simulation scenario, in which individuals
3174 reproduced by asexual (clonal) reproduction. This network is clearly char-
3175 acterized by two distinct clades, separated by long branches.

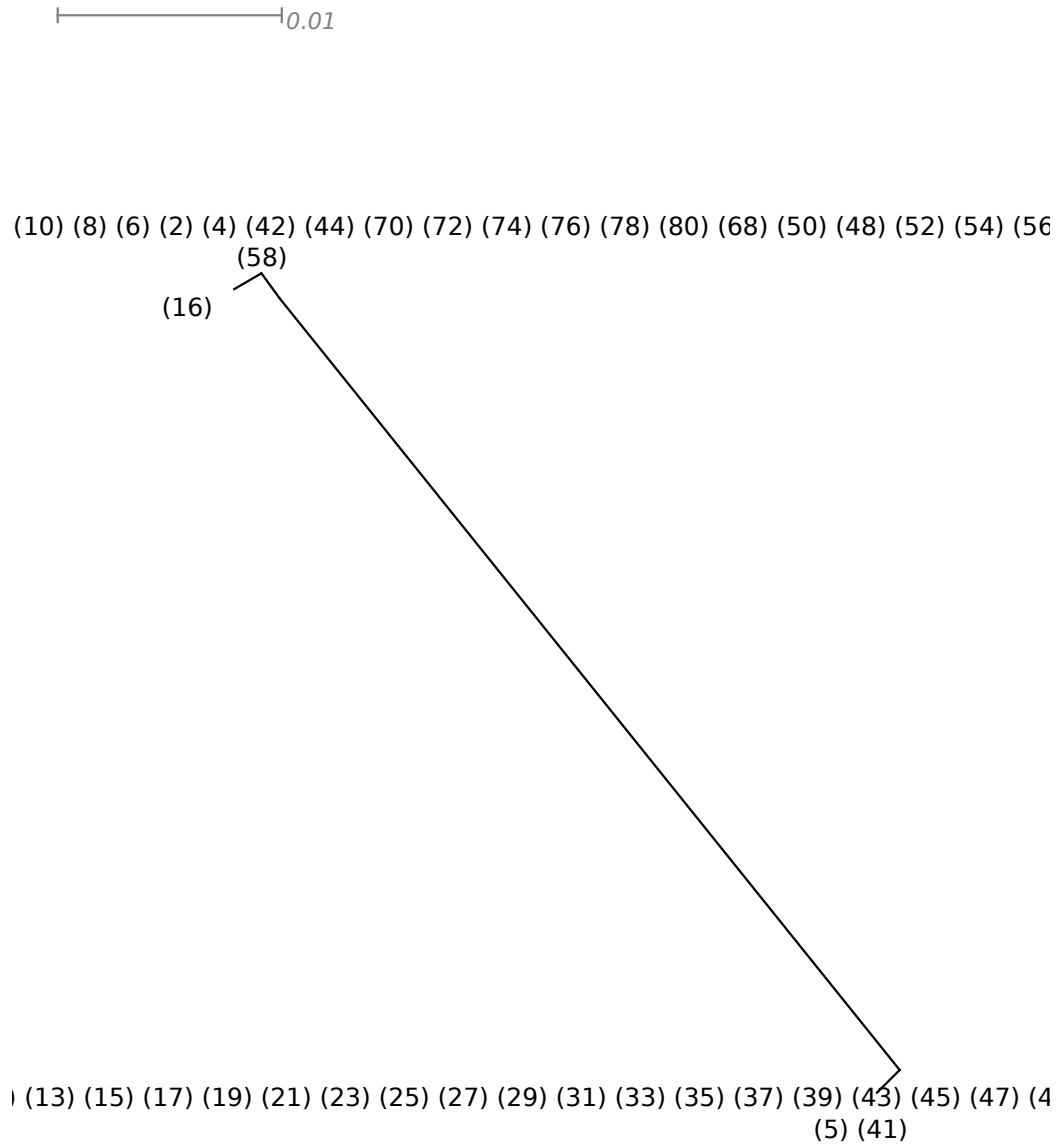


Figure 4.3: Network of simulated allelic pairs, evolved under an asexual reproduction scheme. The first copies of each allelic pair form a clade, and the second copies of each allelic pair form a clade. This is because there is no recombination during gamete formation, as with clonal reproduction, offspring are clones of their parent.

Table 4.1: PHI-Score and Neighbor Similarity Scores of the PCR amplified sequences for three different window sizes.

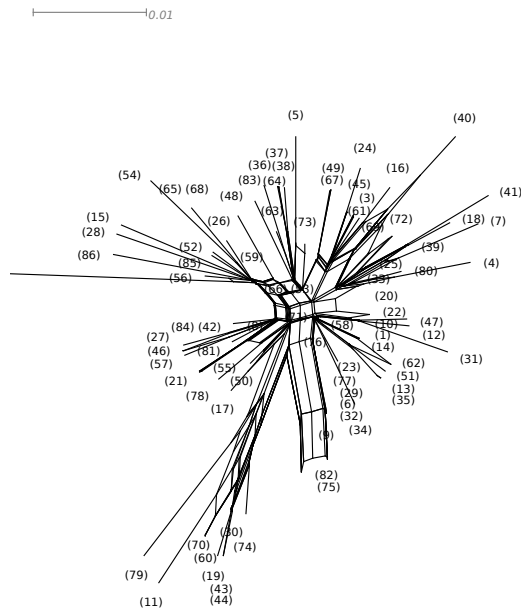
Sequences	Window Size	PHI Score	P-Value	NSS	NSS P-Value
FeABC	100	0.0930	0.0000405	0.81056	0.005
FeABC	50	0.0955	0.0041100	0.81056	0.004
FeABC	10	0.0870	0.0814000	0.81056	0.006
Fcl10	100	0.0930	4.0500000	0.81056	0.005
Fcl10	50	0.0385	0.0184000	0.88306	0.342
Fcl10	10	0.0500	0.2650000	0.88306	0.338

3176 If *F. cylindrus* has a history of asexual reproduction and ancient allelic
 3177 divergence, then it is expected that the networks calculated for the PCR
 3178 amplified sequences of the ABC Iron Transporter and the Large Ribosomal
 3179 Subunit will have a similar structure to that of the network in Figure 4.3.

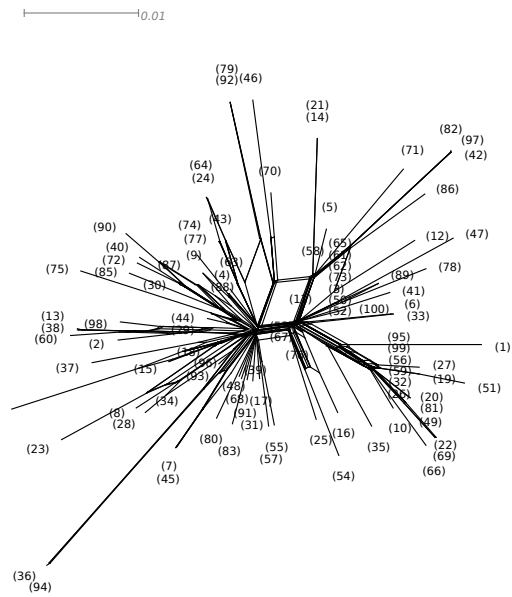
3180 Panels a and b in Figure 4.4 show the phylogenetic networks calculated
 3181 for the PCR amplified sequences of the ABC Iron Transporter (a), and the
 3182 Large Ribosomal Subunit (b). These two networks are clearly different
 3183 qualitatively to the kind of network in Figure 4.3 that would be expected if
 3184 *F. cylindrus* had a history of asexual reproduction without meiotic recomb-
 3185 nation. They do not show a clear partition between two clades or clusters,
 3186 instead they have average branch lengths of around 0.1, and contain around
 3187 255 splits.

3188 Panel a of Figure 4.5, demonstrates the effect of increasing or decreas-
 3189 ing θ in population genetic simulations, on the resulting sequences, and
 3190 thus the networks produced: The average branch lengths in networks, is
 3191 positively related to the θ parameter set in the simulation.

3192 Figure 4.6 presents this relationship qualitatively with the networks
 3193 produced by Splitstree. From figures 4.5 and 4.6 it can be seen that the
 3194 networks best matching the real sequence networks (figure 4.4) in terms of
 3195 branch lengths, are those produced by simulations where $\theta = 0.066$, which
 3196 is close to the value which LAMARC has estimated.

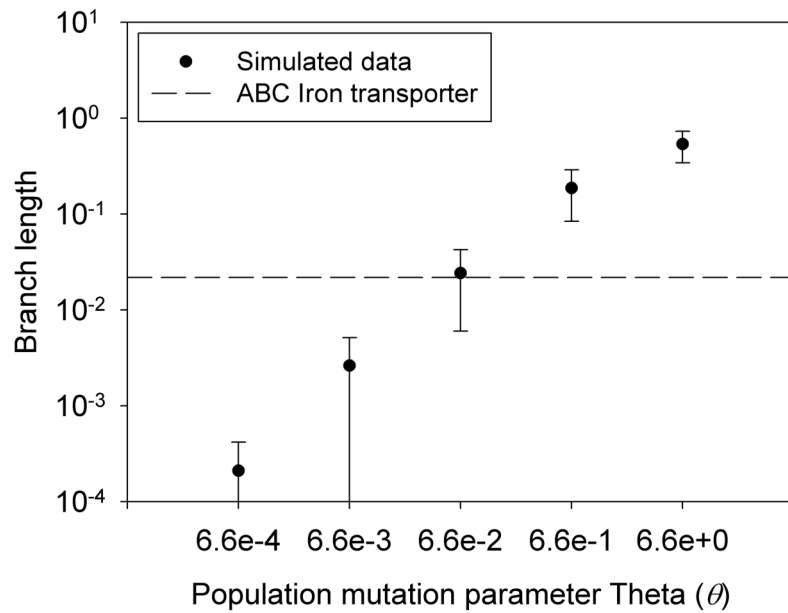
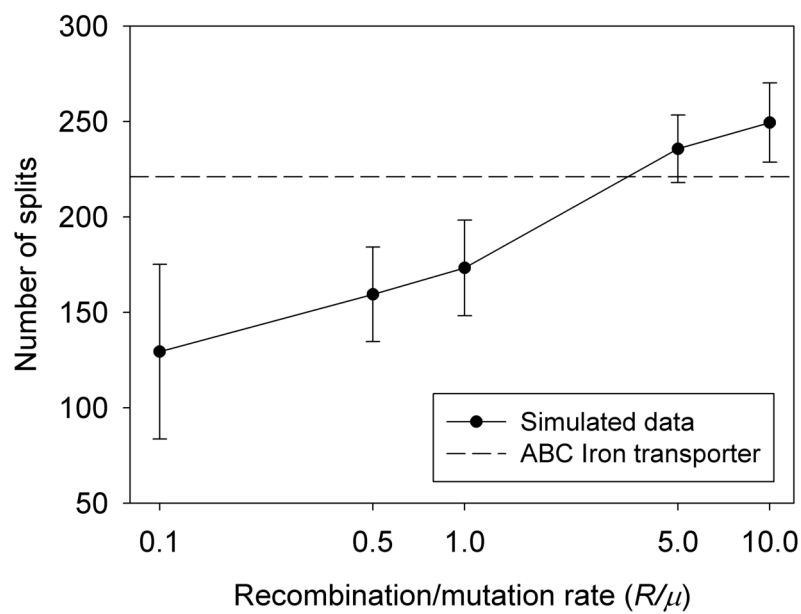


(a) ABC Iron Transporter



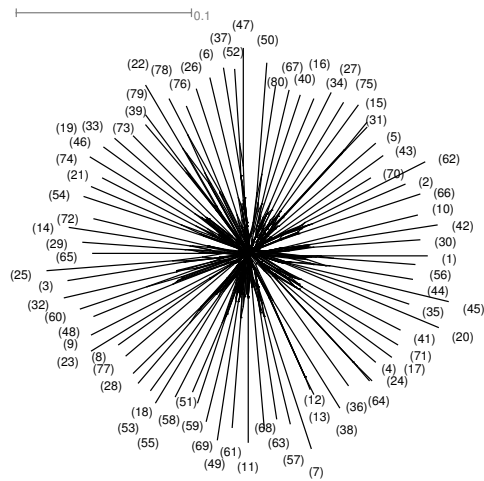
(b) Large Ribosomal Subunit

Figure 4.4: Split Networks of the ABC Iron Transporter and Ribosomal Subunit sequences have average branch lengths close to 10^{-2} and contain 225 splits.

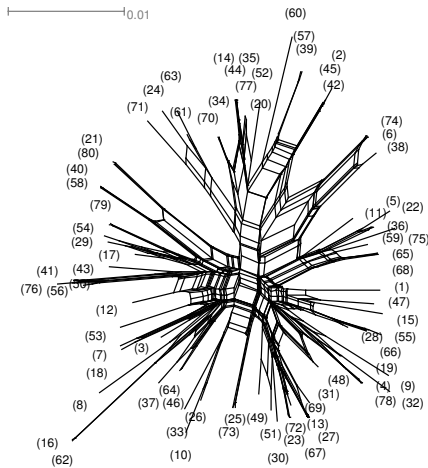
(a) The effect of θ on network branch lengths

(b) The effect of recombination rate on splits

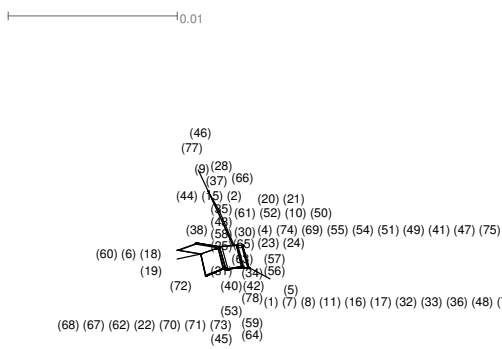
Figure 4.5: Quantifying the branch lengths and number of splits in networks produced from simulations with varying levels of recombination and values of θ . Larger values of θ cause longer branches (a), and higher recombination rates result in more splits (b).



(a) $\theta = 0.66$



(b) $\theta = 0.066$



(c) $\theta = 0.0066$

Figure 4.6: Networks computed from simulations with three different values of θ . Larger values of θ result in longer outer branches of networks.

3197 Panel b of figure 4.5 demonstrates the effect of varying the recombina-
3198 tion rate relative to the mutation rate in population genetic simulations, on
3199 the sequences and networks produced: The number of splits in networks
3200 is positively related to the recombination rate, relative to the mutation rate.
3201 This relationship is shown qualitatively in the networks drawn in figure 4.7.

3202 4.4 Discussion

3203 The phylogenetic networks resulting from population genetic simulations
3204 support several assumptions we had about how recombination, and popu-
3205 lation mutation rate (θ) may be inferred from phylogenetic networks. Specif-
3206 ically, (1) the levels of Theta affect the average branch lengths of the
3207 networks, and (2) the extent of recombination affects the number of splits
3208 in phylogenetic networks. These two assumptions are not controversial: a
3209 higher population mutation rate leads to more mutations in a population
3210 the same amount of time, and thus would lead to longer branches in any
3211 phylogeny or network computed for sequences sampled from the popu-
3212 lation (Frankham 1996; Hein, Schierup, and Wiuf 2004; Wakeley 2009).
3213 Phylogenetic Split Networks (Huson 1998) were conceived of as a way to
3214 detect and represent reticulate evolution. Wherever there is a non-tree like
3215 structure or loops, recombination may be inferred. The networks resulting
3216 from the simulations confirm these assumptions, and so give confidence
3217 in any inferences made about the population and evolution of *F. cylindrus*
3218 from the networks of the ABC Iron Transporter sequences, and the Large
3219 Ribosomal Subunit Sequences.

3220 Secondly, from comparisons between the networks of the ABC Iron
3221 Transporter sequences, Large Ribosomal Subunit Sequences, and simu-
3222 lated networks, it was concluded that LAMARC (Kuhner 2006) estimate
3223 of Θ was a reasonable estimate for the population of *F. cylindrus*. It was

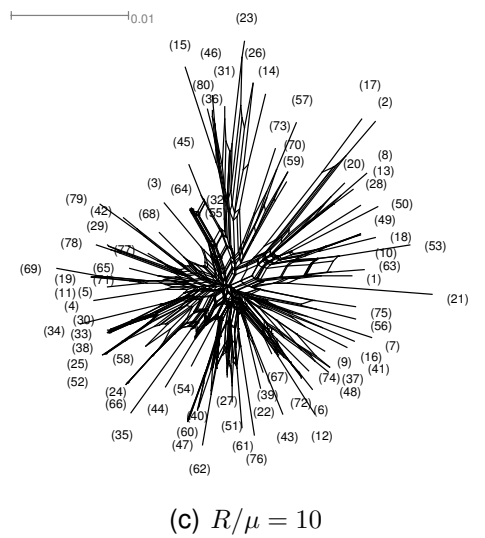
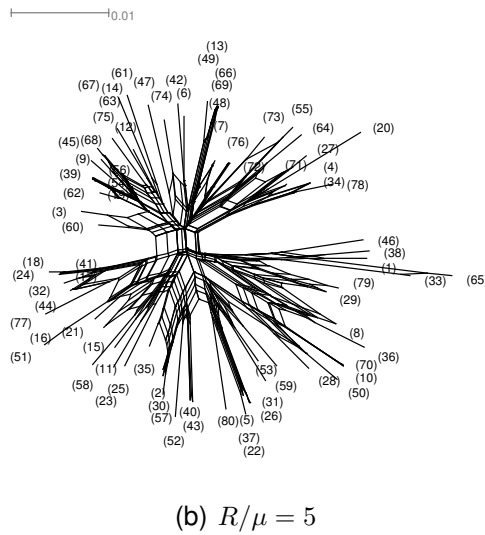
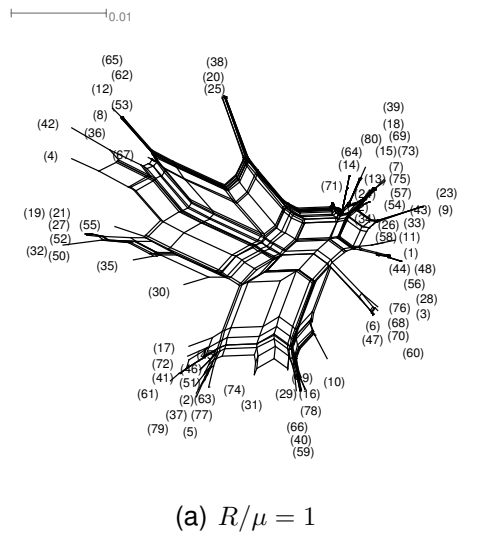


Figure 4.7: Networks computed from simulations with three different levels of recombination, relative to the mutation rate μ . Larger values of R result in more splits in networks.

3224 also concluded that these networks provide evidence of recombination for
3225 the sequences of the ABC Iron Transporter, and in the sequences of the
3226 Large Ribosomal Subunit. Evidence of recombination does not have to
3227 mean that an organism is reproducing sexually, meiotic recombination is
3228 associated with sexual reproduction, but mitotic recombination could also
3229 explain the recombination signal detected in these sequences. However,
3230 whilst mitotic or meiotic recombination may explain the recombination signal
3231 in the sequences, it was concluded that ancient asexuality is not a likely
3232 explanation, because of the lack of similarity of the ABC Iron Transporter,
3233 and the Large Ribosomal Subunit networks, to the networks generated by
3234 simulations of ancient asexual evolution.

3235 **4.4.1 Sex and the diatom reproductive cycle**

3236 Even though sexual reproduction has not been observed in the lab cultures
3237 of *F. cylindrus*, this diatom does not appear to be an ancient asexual. This
3238 might not be surprising given what is already known about Diatom biology
3239 and sexual reproduction. The typical cell cycle of Diatoms is diplontic i.e.
3240 the vegetative cells are diploid, and the haploid gametes are short lived
3241 (Chepurnov et al. 2004).

3242 The Diatom life cycle features two key phases which may be summarized
3243 by the following:

3244 The first phase is a long vegetative phase; this phase can last for months
3245 or years. During this phase, vegetative cells divide by mitosis, gradually
3246 becoming smaller. The cell size decrease during the vegetative phase of
3247 the diatom life cycle is due to the shape and structure of Diatom cell walls
3248 and the division pattern of the Diatoms. The cell wall is made of silicated
3249 components, which together are termed the frustule. The frustule is made
3250 of two overlapping halves or thecae (Chepurnov et al. 2004; Davidovich and
3251 Bates 1998; Poulickova 2008). These thecae are not the same size, the

3252 larger of the two thecae is called the epitheca, and the smaller of the two is
3253 called the hypotheca. When mitosis occurs, cytokinesis splits the diatom
3254 where the two thecae overlap. The two resulting daughter cells inherit one
3255 of the parent cells two thecae as its own epitheca, and they grow their own
3256 hypotheca (Chepurnov et al. 2004). Since one of the daughter cells inherits
3257 a hypotheca as its epitheca, it will be smaller in size to its parent cell. Thus
3258 the average cell size of a population of diatoms decreases as mitotic cell
3259 division occurs.

3260 The second phase is shorter, and includes sexual reproduction and the
3261 production of new vegetative cells, restoring the cell size (Chepurnov et al.
3262 2004). Production of gametes during the sexual reproduction phase has
3263 been demonstrated to occur by classical meiosis in many Diatom species.
3264 Diatoms restore their cell size through the production of auxospores, which
3265 result from sexual reproduction (Davidovich and Bates 1998). During aux-
3266 osporation, recombination and cell size restitution occurs: gametes fuse
3267 to form the auxospore, which expands and a new cell is produced within.
3268 The cell walls of the gamete producing cells are lost, and so the auxospore
3269 must then form the shape of the vegetative cells de novo (Chepurnov et al.
3270 2004). If a population of Diatom cells did not undergo sexual reproduction
3271 to produce the auxospores to restore their cell size, the population would
3272 gradually decrease in cell size until they become critically small. At this
3273 point the population would die, and this has been observed in experimental
3274 cultures. Diatom cells can only become sexualized when they are suffi-
3275 ciently small, but they may also not be able to become sexualized if they
3276 become too small or hit the critical cell size before they die (Chepurnov et al.
3277 2004; Davidovich and Bates 1998; Poulickova 2008). The maximum size
3278 of initial diatom cells, the maximum and minimum sizes of cells capable of
3279 sexual reproduction, and the minimum size before death are strict for each
3280 diatom species and are termed cardinal points (Chepurnov et al. 2004).

3281 However, despite the role of sex in the restoration of cell size in diatoms,
3282 it is not always necessary for cell size restoration. For some diatom species,
3283 asexual auxosporulation is a possibility, presumably it is some secondary
3284 modification of a developmental pathway that was sexual, and some species
3285 do not even undergo auxosporulation and exist as entirely as asexual
3286 populations, and their cell size is restored by vegetative cell enlargement
3287 (Chepurnov et al. 2004; Gallagher 1983; Nagai et al. 1995; Sabbe et
3288 al. 2004; Werner 1977). Species such as *Caloneis amphisbaena* and
3289 *Sellaphora pupula* "lanceolate" have been found to exist in populations of a
3290 very limited range of cell size, and this cell size has remained unchanged
3291 after many generations of observation (Mann 1989; Mann et al. 2004).

3292 Therefore, whilst sex is a common feature of the diatom life cycle, and
3293 is important for cell size restoration in many species, it is not unreasonable
3294 to suggest the hypothesis that a diatom like *F. cylindrus* could have evolved
3295 asexually for a long period of time. However, the network reconstructions
3296 and evidence of recombination demonstrated by this study cast doubt on
3297 that hypothesis as an explanation for the diverged alleles.

3298 **4.4.1.1 Allelic Divergence in diatoms may be explained by popula-** 3299 **tion size**

3300 If the ancient asexuality hypothesis is rejected as the explanation of the
3301 diverged alleles in *F. cylindrus*, then an alternative explanation of how this
3302 diatom evolved diverged and functionally differentiated alleles is desired.
3303 These alleles show signatures of positive selection, and they are differ-
3304 entially expressed. The question is; assuming sexual reproduction and
3305 recombination, why does recombination not homogenize the sequence
3306 variation between two alleles over time?

3307 An alternative hypothesis explaining the adaptive evolution of *F. cylin-*
3308 *drus* is a large population size, which would lead to bigger coalescence

3309 times between maternal and paternal loci. In combination with a low re-
 3310 combination rate, this would result in independent adaptive evolution and
 3311 divergence of the different haplotypes. This is intuitive if one considers a
 3312 coalescent process back through time of an idealized population, because
 3313 the coalescent relates genetic diversity to demographic history. In such a
 3314 process, the probability that any two lineages extant at time t , coalesce in
 3315 the previous generation $t-1$, is the probability that they share a parental DNA
 3316 sequence. For a diploid population there are $2N_e$ alleles in every generation,
 3317 assuming a constant population size (Hein, Schierup, and Wiuf 2004). As-
 3318 suming random mating and neutral evolution, the probability any two alleles
 3319 coalesce in the previous generation (i.e. they share the same parental
 3320 sequence) is $1/(2N_e)$. Therefore, the probability those two alleles do not
 3321 coalesce, is $1 - 1/(2N_e)$. These probabilities are dependent on the size of
 3322 the population in question (Wakeley 2009). Larger populations, result in
 3323 a smaller probability that two alleles coalesce in the previous generation,
 3324 and a greater probability that they do not. With each successive previ-
 3325 ous generation, the probability of coalescence is geometrically distributed
 3326 (Hein, Schierup, and Wiuf 2004; Wakeley 2009). This means that it is the
 3327 product of coalescence at the generation of interest and the probability of
 3328 non-coalescence at the preceding generations i.e.

$$P_c(t) = \left(1 - \frac{1}{2N_e}\right)^{t-1} \left(\frac{1}{2N_e}\right) \quad (4.1)$$

3329 From this equation, it can be seen that with larger populations, the prob-
 3330 ability that two alleles coalesce further back in time is greater i.e. the
 3331 expected coalescence time between two alleles is larger, therefore alleles
 3332 are expected to be more diverged.

3333 This explanation is consistent with the estimation of a Θ of 0.066 by the

3334 LAMARC (Kuhner 2006) analysis, which is also supported by the simula-
3335 tions. The population mutation rate Θ , is proportional to the product of the
3336 mutation rate and the effective population size and so the value predicted by
3337 LAMARC could be the result of a very large population. Furthermore, prior
3338 research has been performed to estimate the abundance of *F. cylindrus* in
3339 water columns around the Antarctic (Kang and Fryxell 1992). During the
3340 summer, numbers of 7.9×10^{10} cells m^{-2} were observed, and during the
3341 winter, numbers of 1.1×10^8 cells m^{-2} were observed. Marginal ice zones
3342 are known to be sites with much dynamic activity such as jets, eddies,
3343 currents, melting, freezing, and upwelling (Kang and Fryxell 1992). They
3344 are also known to be sites of increased phytoplankton biomass and primary
3345 productivity, due to their light levels, ice-distribution, and vertical stability.

3346 Therefore, the hypothesis that a large population size explains the
3347 levels of diversity is consistent with both population genetic (coalescent)
3348 theory, results of this study, as well as the findings of other research. It is
3349 also attractive, because of its simplicity. It is much more plausible that a
3350 phytoplankton species has very large populations; than it is that the species
3351 had abandoned sex as a reproductive strategy: Sex is a common aspect
3352 of the diatom life cycle and is often essential for cell size restoration and
3353 population survival. Furthermore, as was explored in the Introduction, there
3354 is a substantial body of theory explaining why sexual reproduction evolved
3355 two become a widespread reproductive strategy, and is advantageous,
3356 despite the apparent costs.

3357 4.4.1.2 Study limitations and subsequent FALCON assembly

3358 However, this study has some limitations which should be acknowledged
3359 when considering these results. First, whilst evidence of recombination
3360 in the form of the splits networks and the presence of incompatible sites
3361 is obtained from these sequences, it was not possible to examine any

3362 larger recombination events or blocks as was possible for *Albugo candida*
3363 in chapter 3. Indeed, the number of informative sites was too small for the
3364 HybridCheck software (which was implemented to analyse large contigs) to
3365 effectively run. Secondly, the analyses were only performed on two genes.
3366 Whilst it was concluded that the two genes were representative of the larger
3367 set of diverged alleles (figure 4.1), we do not know if the PCR primers
3368 amplified only maternal and paternal alleles of those genes or if some of
3369 the sequences amplified also represent paralogues.

3370 The question of whether the diverged alleles observed in the assembly
3371 were truly diverged alleles was resolved for the assembly described in
3372 the introduction experimentally. Single haplotyped fosmids were Sanger
3373 sequenced by collaborators, providing contiguity information and they were
3374 compared with the assembled genomic scaffolds, and an annotated protein
3375 set from the diverged regions in the genome. Data from these comparisons
3376 revealed a clear separation between allelic pairs and gene duplications
3377 based on 100% identity to the haplotyped Sanger sequenced fosmids.
3378 Additionally, the nucleotide similarity of the diverged alleles (mean = $97.01 \pm$
3379 0.03%) is significantly (p -value $< 10^{-09}$) higher than for gene duplicates
3380 (mean = $84.07 \pm 0.36\%$). Therefore, whilst it may be that some uncollapsed
3381 regions of the assembly could be duplicates, there is high confidence that
3382 the allelic pairs identified are indeed diverged alleles and not duplicates.

3383 Since this work has been completed, an assembly has been completed
3384 using PacBio long read sequencing technology, which has also supported
3385 that true diverged alleles have been identified and that they are not dupli-
3386 cated sequences (although duplicated sequences are indeed present in
3387 the genome). The sequencing work and library preparation was completed
3388 by the platforms and pipelines team. A 20kb fragment length library was
3389 constructed, and a 4kb insert size library was also created. Both libraries
3390 were sequenced using the PacBio RS2 instrument, using SMRT cells with

3391 the c4P6 chemistry. The 20kb fragment length library yielded 1.37Gb of
3392 data, and the 4kb insert size library yielded 3.85Gb of data. The final N50
3393 of read length varied between 8215 to 8898bp for the 20kb fragment length
3394 library, and the N50 ranged from from 2558 to 2680bp for the 4kb insert
3395 size library.

3396 Assembly was completed by collaborator Pirita Paajanen, who combined
3397 the data from the SMRT cells and filtered the shortest reads, yielding 3.8Gb
3398 of data which gave 63x coverage. Assembly was completed using the
3399 diploid aware PacBio assembler, Falcon 0.3.0. The output of the Falcon
3400 assembler was divided into two parts. The haploid assembly resulted
3401 in primary contigs from which a genome size of 59.7Mb was deduced.
3402 However, the assembler also produced alternate contigs which were the
3403 result of the assembler being unable to decide between two possible routes
3404 through the genome graph the genome. Such 'bubbles' in the genome
3405 graph represent diverged haplotypes, containing diverged alleles.

3406 The haplotype divergence differed between chromosomes: The longest
3407 chromosome 000000F had only one alternate contig with a length of 6047bp.
3408 In contrast, contig 000002F was 1246645bp long and had 14 associated
3409 alternative contigs, of a total sequence length of 633764bp. For each of the
3410 14 alternative contigs of chromosome 000002F, I extracted and aligned the
3411 two haplotype sequences using the pairwise alignment algorithm available
3412 in the Bio.jl software package (<https://biojulia.github.io/Bio.jl>), using an
3413 EDNA scoring matrix. Once aligned, a non-overlapping sliding window was
3414 moved across the sequences, and the p-distance between the sequences
3415 within each window was calculated. For each computation, the width of
3416 the sliding window was set as 1% of the width of the pairwise alignment in
3417 (bp). The results of this analysis are included as extra information in the
3418 appendix, figure A.1. The figure demonstrates different levels of divergence
3419 across the diverged haplotype pairs, including the appearance of indels

3420 between some haplotype pairs. Further work will show how the sequences
3421 of allelic pairs align to the FALCON assembly, revealing which pairs align
3422 to different haplotypes of a FALCON 'bubble' (true allelic pairs), and which
3423 pairs align to the same haplotype of a FALCON 'bubble' (potentially gene
3424 duplicates).

3425 At the time of writing, multiple population samples of *F. cylindrus* are not
3426 available, and so analyses presented here used sequences from cultures,
3427 and so further population genetic analyses should be conducted in the fu-
3428 ture as more data becomes available, for example to assess the population
3429 structure of *F. cylindrus* and investigate if gene flow is occurring between
3430 subpopulations of *F. cylindrus*.

3431 The fact that the genome assembly contains some duplicates, and
3432 that some of the allelic pairs analysed in this study may be found to be
3433 duplicates is not problematic for the hypothesis that this Diatom species has
3434 adapted through allelic divergence, as it may be argued allelic divergence
3435 could lead to gene duplication and the conditions for the divergence of
3436 alleles and the divergence of duplicates overlap: When diverged alleles are
3437 maintained in a population due to heterozygote advantage, duplications
3438 may rapidly spread through the population, causing an individual to act
3439 as a genetic heterozygote yet still breed true. Proulx and Phillips 2006
3440 argued that genetic redundancy is the mechanism usually cited as allowing
3441 duplicate genes to diverge, but redundancy is present in a diploid before
3442 duplication: Dominance creates the same kind of redundancy duplicates
3443 have, but for alleles of single copy genes. Therefore mode of inheritance is
3444 the thing then which most distinguishes duplicates from single copy genes:
3445 Segregation prevents the fixed inheritance of alternative allelic variants at
3446 a single locus (Proulx and Phillips 2006). In other words, heterozygotes
3447 at one locus are broken up by segregation during sexual reproduction,
3448 whereas duplicate loci in an individual can carry copies of alternate alleles

3449 at different loci. Their results show that fitness relationships that allow
3450 divergent alleles to evolve at one locus overlap significantly with those that
3451 allow the divergence of previously duplicated genes at two different loci
3452 (Proulx and Phillips 2006).

3453 4.4.1.3 Conclusion

3454 The genome of the polar diatom *Fragilariopsis cylindrus* contains diverged
3455 alleles that are differentially expressed in different environmental conditions.
3456 Evidence of recombination was found which contradicts the ancient asexual-
3457 ity hypothesis explaining how these diverged alleles may have evolved. An
3458 alternative, competing hypothesis is proposed, supported by the evidence
3459 presented, that a large population size has allowed diversifying selection to
3460 differentiate the alleles of genes despite the homogenizing effect of recomb-
3461 nation. Additional population samples, and analysis of larger contigs made
3462 possible by improved genome assembly for recombination, will help answer
3463 the question of how *F. cylindrus* has evolved this remarkable strategy to
3464 cope with varying environmental conditions.

3465 CHAPTER 5

3466 General Conclusion

3467 5.0.1 Summary and Conclusions

3468 In this thesis, work focused on how recombination facilitates the adaptive
3469 evolution of a plant pathogen and a polar marine diatom. Both of these
3470 organisms were of evolutionary interest due to aspects of their lifestyles
3471 and/or physiology: The plant pathogen *Albugo candida* was of interest
3472 because whilst it was an obligate biotroph, it has a very large host range,
3473 and the diatom *Fragilariopsis cylindrus* was of interest because the genome
3474 sequencing project and differential expression experiments revealed genes
3475 with diverged alleles that were differentially expressed in different environ-
3476 mental conditions.

3477 Recombination is important for the formation of novel genotypes, haplo-
3478 types and alleles, therefore it plays a key role in adaptive evolution (Grauer
3479 and Li 2000). Recombination separates deleterious mutations from their
3480 genomic background, in combination with purifying selection this reduces
3481 the mutational load (Lynch and Gabriel 1990). Recombination also brings
3482 beneficial mutations from separate lineages into one individual or lineage.
3483 However, recombination also plays a fundamental role in the repair of
3484 damaged DNA, when homologous recombination replaces a damaged
3485 DNA strand with its intact counterpart, and it was likely this function of

3486 recombination that was important in early prokaryotic life and evolution
3487 (Cavalier-Smith 2002). With respect to adaptive evolution, however, the
3488 principal consequence of recombination is that it generates novel combina-
3489 tions of nucleotides, which in turns allows for selection to act a much finer
3490 scale, i.e. at the level of nucleotides rather than the entire genome.

3491 The potential of recombination to generate novel allelic combinations
3492 is important for host and pathogens which are engaged in an evolutionary
3493 arms race to adapt and counter adapt to each others molecular mecha-
3494 nisms of pathogenicity or immunity. The red queen hypothesis explains the
3495 advantage of sexual reproduction in such terms. The variability generated
3496 by sexual reproduction (and meiotic recombination) results in genetically
3497 unique offspring, which permits a faster response to selection (Paterson
3498 et al. 2010). As a result sexually reproducing species are able to improve
3499 their genotype in changing conditions. Co-evolutionary interactions be-
3500 tween host and parasite select for sexual reproduction in hosts in order
3501 to reduce the risk of infection. Oscillations in genotype frequencies are
3502 observed between parasites and hosts in an antagonistic co-evolutionary
3503 way without necessitating changes to the phenotype, and in host-parasite
3504 co-evolution systems with multiple hosts, Red Queen dynamics may affect
3505 which host and parasite types become common (or rare) (Charlesworth
3506 and Charlesworth 2010).

3507 It was hypothesized that the *Albugo candida* species was composed of
3508 several host-specialised races, each locked in an evolutionary arms race
3509 with their specific host. Such a race with a specific host would lead to
3510 further divergence and possibly speciation of the races. However, *Albugo*
3511 is known to be able to suppress non-host resistance. Infections of *Albugo*
3512 *sp.* could suppress the runaway cell death phenotypes of plants, allowing
3513 formerly avirulent strains of downy mildew to infect (Cooper et al. 2008).
3514 Assuming that this ability extended to other non-host species, *Albugo*

3515 may be modeled as a 'microbial hub': taxa that are integral and highly
3516 connected to the network of a hosts microbial community. Such hubs may
3517 affect community compositions through microbe-microbe interactions or, as
3518 seems to be the case with *Albugo*, suppression of host defense responses
3519 (Agler et al. 2016). Therefore, non-host immune suppression would enable
3520 host-specific races of *Albugo candida* to overcome the ever increasing
3521 barrier to gene flow that specialisation imposes, and sexual reproduction
3522 between races, followed by introgression by back-crossing, would permit
3523 the generation of a range of novel genotypes. Consequently the species
3524 could evolve its wide host range.

3525 To assess this hypothesis it was necessary to scan the genome of
3526 *Albugo candida* isolates to identify recombinant regions. Furthermore, to
3527 distinguish such regions as recombinant and not the result of incomplete
3528 lineage sorting due to rapid divergence, the regions identified needed to be
3529 tested for significance and the coalescence times estimated.

3530 Scans of the genomes for recombination revealed a highly recombined
3531 mosaic genome, and therefore a rapid coalescence estimation method
3532 for all of the recombination blocks was desired, in addition to a method of
3533 plotting which effectively demonstrated the high degree of mosaic-ism in the
3534 *A. candida* genome. Therefore, rapid detection and dating of recombination
3535 blocks was implemented, and the software package HybridCheck was
3536 created and tested using simulated data as in chapter 2. HybridCheck
3537 was also tested for consistency with RDP analyses of *A. candida*, which
3538 identified recombination, and BEAST estimates of coalescence times for a
3539 subset of the identified recombination regions (chapter 3). The evidence
3540 presented in chapter 3 confirmed the model of *Albugo candida* evolution:
3541 Isolation, divergence and specialisation of races generates repertoires of
3542 effectors for a specific race. Those adapted repertoires are then brought
3543 together when two races hybridize. The result is the generation of novel

3544 repertoires of novel combinations of these effectors. Specific avirulence
3545 effectors that trigger host immunity may be lost through segregation and
3546 through loss of heterozygosity (Lamour et al. 2012; McMullan et al. 2015).
3547 Hybrids, with new combinations of effectors, and having lost effectors which
3548 impeded their colonisation of other hosts previously, may expand their
3549 geographical range and population size clonally. Some of these hybrids
3550 may be able to colonise new hosts, expanding the host range.

3551 The genome assembly project of *F. cylindrus* revealed that the genome
3552 contained 21,066 predicted protein-encoding genes, 6,071 genes were
3553 represented by diverged alleles, and each pair of diverged alleles had both
3554 coding and non-coding regions, and were up to 6% polymorphic in the
3555 non-coding regions. Furthermore, differential expression experiments and
3556 RNA-Sequencing suggested that 40% of the non-collapsed, diverged allelic
3557 pairs showed a 4 fold unequal bi-allelic expression (Mock et al. 2017).

3558 Alleles showing the strong unequal bi-allelic expression were found
3559 to have an elevated rate of non-synonymous mutations, which suggests
3560 significant positive / adaptive selection and evolution of these allelic pairs
3561 (Mock et al. 2017). It was concluded therefore, that positive selection has
3562 been a driving force in the evolution of these diverged alleles and hence
3563 the adaptation of this diatom to the environmental conditions it faces.

3564 An evolutionary explanation was hypothesized: The alleles of an allelic
3565 pair could diverge as a result of positive selection because there was a long
3566 history of asexual reproduction in the organism, and hence an absence of
3567 recombination acting as a homogenizing force between alleles.

3568 However, results from recombination detection analysis, and phyloge-
3569 netic network construction of PCR amplified sequences from DNA extracted
3570 from *F. cylindrus* cultures conflicted with results of the same analyses per-
3571 formed with DNA sequences obtained by population genetics individual
3572 based simulations of ancient asexuality. Indeed the results for *F. cylindrus*

3573 were more consistent with those of simulations of a scenarios of sexual
3574 reproduction, and a large Θ value. This result suggests an alternative
3575 competing hypothesis, that very large effective population sizes could have
3576 led to the divergence of the alleles in each allelic pair as a result of posi-
3577 tive selection, in the face of the homogenizing influence of recombination
3578 through sexual reproduction.

3579 **5.0.2 Impact and potential future directions**

3580 **5.0.2.1 *Albugo candida***

3581 A paper describing the extent of the introgression identified within the *A.*
3582 *candida* genome was published in eLife (McMullan et al. 2015). According
3583 to Google Scholar, the study has been cited 11 times at time of writing.
3584 Citations include reviews of the role of hybridisation and introgression in
3585 the adaptive evolution and emergence of new fungal and filamentous plant
3586 pathogen strains (Depotter et al. 2016; Dong, Raffaele, and Kamoun 2015;
3587 Stukenbrock 2016), research demonstrating the role of recombination in
3588 the evolution of the Rp1 resistance genes in grasses (Jouet, McMullan, and
3589 Oosterhout 2015), and a study presenting evidence that for *Coleosporium*
3590 *ipomoeae*, any genotypes can infect multiple hosts from non-local commu-
3591 nities, but only are highly host specific when tested on hosts from local
3592 communities, calling into question theoretical results of single-pathogen
3593 single-host studies which suggest that selection favours genotypes with a
3594 broad host range (Chappell and Rausher 2016). Following the 2015 eLife
3595 paper, Belhaj et al. 2015 published a more extreme example of the ability
3596 of *Albugo spp.* to suppress the host immune system. They found that
3597 *Phytophthora infestans*, which is typically a potato and tomato specialist
3598 pathogen, was capable of infecting the plant model organism *Arabidopsis*
3599 *thaliana* when *Albugo laibachii* has also colonized the plant. The nature of

3600 the *P. infestans* infection was similar to that of an *Albugo laibachii* infection:
3601 Transcription profiling of *P. infestans* infections revealed a significant overlap
3602 between the sets of secreted proteins of *P. infestans* during infection of
3603 *Arabidopsis thaliana* and during infections of potato. This suggests there
3604 is similar gene expression dynamics on the two species, and it raises the
3605 question. Is gene flow between two different Oomycete species possible?
3606 And could this contribute to adaptive evolution of these pathogens.

3607 It is well established that *Albugo* suppresses non-host immunity in hosts
3608 it infects, and as a result of work presented in this thesis it was concluded
3609 that this lowers barriers to gene flow and permits introgression, facilitating
3610 the generation of novel pathogen haplotypes and enabling *Albugo can-*
3611 *dida* to evolve a wide host range. However, this model of *Albugo candida*
3612 evolution raised a conceptual problem: This phenomenon appears to ex-
3613 tend to other pathogen species that were not *Albugo spp.* (Belhaj et al.
3614 2015), and therefore *Albugo spp.* may act as a microbial hub as previously
3615 noted. If this is the case, how is it that *Albugo spp.* (obligate biotrophs
3616 with a vital dependence on the host) can compete in this limited niche,
3617 whilst at the same time enable non-host colonization for other pathogen
3618 species who are then presumably competitors for the same resource. An
3619 answer to this problem was provided by a paper from Ruhe et al. 2016.
3620 Shotgun proteomics was completed of the apoplastic fluid of samples of
3621 lab-grown *Arabidopsis thaliana* that were infected with *Albugo spp.*, and
3622 samples which were uninfected. Work was repeated for wild-grown *Ara-*
3623 *bidopsis thaliana* and they found that whilst both lab-grown and wild-grown
3624 *Arabidopsis thaliana* supported extensive *Albugo* colonization (Ruhe et al.
3625 2016). However, no or low levels of defense-related proteins were detected
3626 in lab samples, but regardless of *Albugo spp.* infection status, wild plants
3627 showed a broad spectrum of defense-related proteins at high abundances
3628 and lab-grown plants did not. These results suggest that *Albugo spp.*

3629 do not strongly affect immune responses and leave distinct branches of
3630 the immune signaling network intact (Ruhe et al. 2016). This suggests
3631 that the pathogens of the *Albugo* genus, including *Albugo candida* in the
3632 wild are fine tuned to avoid triggering strong host defense reactions, but
3633 also to avoid a broad-spectrum host defense suppression, thus allowing
3634 them to avoid competition from other species growing in the same niche
3635 (Ruhe et al. 2016). Since races of *Albugo candida* are members of the
3636 same species, they may still colonize the same host plant at the same
3637 time, allowing introgression to occur (explaining the introgression signal
3638 observed), but other more distantly related competing pathogens may be
3639 excluded by this precise host immunity manipulation observed by Ruhe
3640 et al. 2016, and so may not get to compete with *Albugo spp.*. However
3641 this experiment only examined *Arabidopsis thaliana* as a host, and crops
3642 grown in monoculture are often uniform and subject to artificially maintained
3643 conditions and treatments, and this may be considered analogous to plants
3644 grown in laboratory conditions. So it is uncertain whether in monoculture
3645 environments *Albugo spp.* manipulate their host immune systems subtly
3646 and precisely, thus avoiding colonization of competition, or whether as with
3647 lab-grown *Arabidopsis thaliana* they do significantly affect the secretome of
3648 the host, allowing competitors to colonize.

3649 In the future, additional study of more strains and population samples of
3650 *Albugo candida* is desirable, since the study presented in this thesis only
3651 examined the genomes of three 'races', and more samples might increase
3652 the number of *Albugo candida* races we can analyse. Future potential work
3653 also includes disentangling the true branching order of *Albugo candida*
3654 races, and improving the detection and dating methods used to analyse
3655 *Albugo candida* genomes (see below).

3656 5.0.2.2 HybridCheck

3657 The HybridCheck software package was initially created out of a need
3658 specific to the *A. candida* project in chapter 3. Following the *A. candida*
3659 project, the HybridCheck software was published in a short software note
3660 in Molecular Ecology Resources (Ward and Oosterhout 2016), and other
3661 groups across the Norwich Research Park became interested in using it
3662 with their own study systems.

3663 In particular, researchers at Norwich Medical School working on *Cryp-*
3664 *tosporidium* used HybridCheck to perform chronological assessment of
3665 recombination events identified in the genomes of three trains of *C. parvum*
3666 (IIaA15G2R1, IIcA5G3j, IIcA5G3a), and a single *C. hominis* (IbA10G2)
3667 GP60 sub-type strain (Nader 0). They found 104 unique recombination
3668 events, and a skewed distribution of recombination events across chromo-
3669 somes. More recombination events were identified on chromosome 6, and
3670 a greater number of events was observed for *C. parvum anthroponosum*
3671 sub-type IIcA5G3a than for any other strain. More than 90% of all recomb-
3672 nation events occurred proximal to loci suspected to drive virulence or play
3673 a major role in host-parasite interactions in human cryptosporidiosis. There-
3674 fore it appears that in this pathogen too, recombination is an important force,
3675 generating novel gene combinations and driving the adaptive evolution of
3676 a pathogen to its host (Nader 0). The estimated divergence dates calcu-
3677 lated in their study provide the first chronological description for genetic
3678 introgression between human-infective *Cryptosporidium spp.* HybridCheck
3679 analyses revealed a chromosome-wide consensus that places a majority of
3680 introgression events between zoonotic (IIaA15G2R1 and IIcA5G3j) and an-
3681 throponotic (IIcA5G3a) *C. parvum* sub-type strains at approximately 10-15
3682 thousand generations ago, while genetic introgression (or recombination)
3683 between the two more closely related zoonotic strains appears to be more

3684 recent (between approximately 3 to 5 thousand generations ago) (Nader
3685 0).

3686 Based on infectivity studies in healthy adult volunteers, the average
3687 generation time within a host is 14.8 hours, and assuming a steady rate of
3688 transmission within host populations, they derived a minimum estimate of
3689 the recombination events of around 5.9 (zoonotic vs. zoonotic *C. parvum*),
3690 17.6 (zoonotic vs. anthroponotic *C. parvum*), and 176.7 (*C. hominis* vs.
3691 *C. parvum*) years ago (Nader 0). In other words, they estimate that the
3692 evolutionary split between the two primary human-infective species appears
3693 to have occurred at the turn of the second industrial revolution, around
3694 1840 (Nader 0).

3695 Whilst this result is putative and needs validation with other dating
3696 methods before publication submission, it is a clear demonstration of the
3697 utility of HybridCheck for researchers in estimating coalescence times
3698 rapidly, across many recombination affected genomic regions.

3699 Future directions for work involving HybridCheck include its continued
3700 use in other organisms. For example HybridCheck is already being used
3701 to generate preliminary results for population genomic data for mice (*Mus*
3702 *spp.*), being generated at the Earlham Institute, with the aim of confirming
3703 hypotheses of genetic isolation between species, and identifying potential
3704 introgressions between populations. Future work involving HybridCheck
3705 may also involve programmatic work. Bioinformatics methods and the
3706 detection of introgression is an active area of research, and more algorithms
3707 and methods will likely be created in the future. Therefore, HybridCheck
3708 would have greater utility as a provider of different methods for the detection
3709 and dating of recombinant and introgressed regions, that are able to work
3710 on multiple different data sources or formats. As a programming problem,
3711 such software code might be best implemented, using multiple dispatch, to
3712 make it more easily maintained, and more easily used. Multiple dispatch is

3713 a feature of some programming languages in which a function (sometimes
3714 called a method) can be dynamically dispatched based on the type of more
3715 than one of its arguments. This thesis author has already co-founded,
3716 develops, and maintains a new bioinformatics infrastructure and community
3717 called BioJulia, based around a modern new programming language for
3718 scientists and technical programmers, called Julia. The language is high-
3719 level, implements a flexible type and multiple dispatch system, and can
3720 achieve speeds matching those of compiled software written in the C
3721 language, with less lines of code. These features make it ideal for the kind
3722 of rapid and flexible development that Bioinformaticians often do, and should
3723 development of HybridCheck continue towards this goal, the framework
3724 already has many high performance code modules and features that a
3725 BioJulia port of HybridCheck could take advantage of.

3726 In the near future, approaches to recombination detection may also
3727 change. Currently, HybridCheck and other methods typically analyze DNA
3728 or protein sequences and identify regions that are phylogenetically incon-
3729 gruent i.e. where computed phylogenetic topologies change or there is
3730 a change-point in computed genetic distances. After the identification of
3731 these regions, it may be assumed they are recombination, or incomplete
3732 lineage sorting, and subsequent analyses, such as the dating method in
3733 HybridCheck, may be employed to try to distinguish whether the cause
3734 is recombination or incomplete lineage sorting. The cause may also be
3735 assumed based on rates of speciation or population size; incomplete lin-
3736 eage sorting is more likely when either of the two are high. However, as
3737 described in chapter 2, there are problems with this approach which leave
3738 room for future improvement.

3739 For example, recombination blocks can become fragmented by ac-
3740 cumulation of subsequent mutations following the recombination event.
3741 Consequently, older recombination blocks tend to be smaller, when they

3742 are actually larger. Thus, not all mutations are accounted for, resulting in
3743 an underestimate of the divergence time particularly for old recombination
3744 events/regions of incomplete lineage sorting.

3745 Furthermore, some methods of resolving introgression from incomplete
3746 lineage sorting require knowledge of branching orders, and sometimes
3747 these are unknown, and sometimes this is even because of the influence
3748 of introgression or incomplete lineage sorting. To solve this issue for the
3749 malaria parasite, Fontaine et al. 2015 obtained the correct species branch-
3750 ing order of the *An. gambiae* complex and two *Pyretophorus* out-group
3751 species. To do this in the face of introgression and incomplete lineage sort-
3752 ing they used 50kb non-overlapping windows across a genome alignment
3753 and computed phylogenies for each window. At least 85 tree topologies
3754 were observed. When these were sorted according to chromosome arm
3755 and their relative frequency, the most commonly observed topology for
3756 the X chromosome was highly discordant with the most commonly ob-
3757 served topology for the autosomes. They then grouped these phylogenetic
3758 topologies, into three distinct topology categories based on the relative phy-
3759 logenetic positions of two species: *An. arabiensis* and *An. quadriannulatus*,
3760 and they observed the topology category most commonly observed on
3761 the X chromosome, was not the same as for the autosomes. Dating the
3762 internal nodes of phylogenies for each topology category allowed them to
3763 distinguish which category of topology best represented the true branching
3764 order, and which represented topologies that were caused by introgression.
3765 Given that almost all of the autosome was represented by a topology cate-
3766 gory that is affected by introgression and linkage disequilibrium, traditional
3767 phylogenetic methods for resolving a species level topology, which typically
3768 invoke some majority rule, would certainly have resulted in the incorrect
3769 answer.

3770 The method utilized in their work will be of great benefit to researchers

3771 studying complicated genomes where introgression, and incomplete lineage
3772 sorting, are prevalent. A likely future direction for the development of
3773 HybridCheck will be to take these methodological ideas and implement
3774 tools that make it trivial for researchers to decompose the gene trees
3775 computed across a genome, identify topological categories from those
3776 trees, and organize them, before analyzing the divergence times of the
3777 phylogenies in each topological category. In the future HybridCheck should
3778 make it simple to perform such an analysis along with other methods such
3779 as Patterson's D , f_d , and tests to distinguish introgression from incomplete
3780 lineage sorting. It should make it trivial to compile such multiple lines of
3781 evidence into a more complete picture of introgression, incomplete lineage
3782 sorting, and linkage, across genomes.

3783 **5.0.2.3 *F. cylindrus***

3784 The study of *F. cylindrus* is in preparation to be submitted to the journal
3785 Nature this year. As such it is not possible to describe the impact in terms
3786 of a number of citations, or who has cited it and why at this time. However,
3787 as stated in discussion of chapter 4, reviewer comments led to further
3788 sequencing with PacBio technology, which resulted in confirmation that we
3789 had obtained strong evidence of diverged alleles. Furthermore, it is known
3790 that at time of writing, that unpublished data and correspondence from a
3791 colleague and co-author of the paper, Chris Bowler (perscom), that similar
3792 evidence of diverged alleles and differential expression has been found in
3793 another diatom species that his group study. Therefore, it could be that
3794 the data presented in this thesis and in the paper, are the first evidence
3795 of a common phenomenon and mechanism of adaptation in this group of
3796 organisms. Future work on this topic has already been described in the
3797 discussion of chapter 4: Imminent future work will show how the sequences
3798 of allelic pairs previously identified align to the new FALCON assembly.

3799 This will reveal which pairs align to different haplotypes of a FALCON
3800 'bubble' (true allelic pairs), and which pairs align to the same haplotype
3801 of a FALCON 'bubble' (potentially gene duplicates). Currently, multiple
3802 population samples of *F. cylindrus* are not available, and so analyses
3803 presented here used sequences from cultures, and so further population
3804 genetic analyses should be conducted in the future as more data becomes
3805 available, for example to assess the population structure of *F. cylindrus* and
3806 investigate if gene flow is occurring between subpopulations of *F. cylindrus*.

3807 In conclusion, detecting and understanding how recombination is affect-
3808 ing the genomes is critical to understanding how species of interest evolve
3809 and adapt to dynamic environments, this thesis has demonstrated how
3810 recombination appears to have influenced the evolution and adaptation of
3811 two different eukaryotic micro-organisms. Future work will expand on the
3812 bioinformatics methodological techniques implemented in this thesis, as
3813 more and more data becomes available for these two species.

- 3815 Abbott, R, D Albach, S Ansell, J. W. Arntzen, S. J. E. Baird, N Bierne,
3816 J Boughman, A Brelsford, C. A. Buerkle, R Buggs, R. K. Butlin, U
3817 Dieckmann, F Eroukhmanoff, A Grill, S. H. Cahan, J. S. Hermansen,
3818 G Hewitt, A. G. Hudson, C Jiggins, J Jones, B Keller, T Marczewski,
3819 J Mallet, P Martinez-Rodriguez, M Möst, S Mullen, R Nichols, A. W.
3820 Nolte, C Parisod, K Pfennig, A. M. Rice, M. G. Ritchie, B Seifert, C. M.
3821 Smadja, R Stelkens, J. M. Szymura, R Väinölä, J. B. W. Wolf, and D
3822 Zinner (2013). “Hybridization and speciation.” In: *Journal of Evolutionary*
3823 *Biology* 26(2), pp. 229–46.
- 3824 Adams, K. L. and J. F. Wendel (2005). “Allele-specific, bidirectional silencing
3825 of an alcohol dehydrogenase gene in different organs of interspecific
3826 diploid cotton hybrids”. In: *Genetics* 171(4), pp. 2139–2142.
- 3827 Agler, M. T., J. Ruhe, S. Kroll, C. Morhenn, S. T. Kim, D. Weigel, and E. M.
3828 Kemen (2016). “Microbial hub taxa link host and abiotic factors to plant
3829 microbiome variation”. In: *PLoS Biology* 14(1), pp. 1–31.
- 3830 Agrawal, A. F. (2001). “Sexual selection and the maintenance of sexual
3831 reproduction.” In: *Nature* 411(6838), pp. 692–695.
- 3832 Alberts, Johnson, and Lewis (2002). “Recombination”. In: *Molecular Biology*
3833 *of the Cell*. 4th editio. New York: Garland Science. Chap. General Re.
- 3834 Ardlie, K, S. N. Liu-Cordero, M. A. Eberle, M Daly, J Barrett, E Winchester,
3835 E. S. Lander, and L Kruglyak (2001). “Lower-than-expected linkage
3836 disequilibrium between tightly linked markers in humans suggests a role

- 3837 for gene conversion.” In: *American Journal of Human Genetics* 69(3),
3838 pp. 582–9.
- 3839 Armbrust, E. V. (2009). “The life of diatoms in the world’s oceans.” In: *Nature*
3840 459, pp. 185–192.
- 3841 Arnheim, N., P. Calabrese, and M. Nordborg (2003). “Hot and cold spots of
3842 recombination in the human genome: the reason we should find them
3843 and how this can be achieved.” In: *American Journal of Human Genetics*
3844 73(1), pp. 5–16.
- 3845 Arrigo, K. R., D. K. Perovich, R. S. Pickart, Z. W. Brown, G. L. V. Dijken,
3846 K. E. Lowry, M. M. Mills, M. A. Palmer, W. M. Balch, F. Bahr, N. R.
3847 Bates, C. Benitez-nelson, B. Bowler, E. Brownlee, J. K. Ehn, K. E. Frey,
3848 R. Garley, S. R. Laney, L. Lubelczyk, J. Mathis, A. Matsuoka, B. G.
3849 Mitchell, G. W. K. Moore, E. Ortega-retuerta, S. Pal, C. M. Polashenski,
3850 R. A. Reynolds, B. Schieber, H. M. Sosik, M. Stephens, and J. H. Swift
3851 (2012). “Massive phytoplankton blooms under arctic sea ice”. In: *Science*
3852 336(6087), p. 2012.
- 3853 Assmy, P., V. Smetacek, M. Montresor, C. Klaas, J. Henjes, V. H. Strass,
3854 J. M. Arrieta, U. Bathmann, G. M. Berg, E. Breitbarth, B. Cisewski, L.
3855 Friedrichs, N. Fuchs, G. J. Herndl, S. Jansen, S. Kragefsky, M. Latasa,
3856 I. Peeken, R. Rottgers, R. Scharek, S. E. Schuller, S. Steigenberger,
3857 A. Webb, and D. Wolf-Gladrow (2013). “Thick-shelled, grazer-protected
3858 diatoms decouple ocean carbon and silicon cycles in the iron-limited
3859 Antarctic Circumpolar Current”. In: *Proceedings of the National Academy*
3860 *of Sciences* 110(51), pp. 20633–20638.
- 3861 Auton, A. and G. A. T. McVean (2007). “Recombination rate estimation in
3862 the presence of hotspots”. In: *Genome Research* 17, pp. 1219–1227.
- 3863 Baack, E. J. and L. H. Rieseberg (2007). “A genomic view of introgression
3864 and hybrid speciation.” In: *Current Opinion in Genetics and Development*
3865 17(6), pp. 513–8.

- 3866 Baack, E. J., K. D. Whitney, and L. H. Rieseberg (2005). "Hybridization and
3867 genome size evolution: Timing and magnitude of nuclear DNA content
3868 increases in *Helianthus* homoploid hybrid species". In: *New Phytologist*
3869 167(2), pp. 623–630.
- 3870 Baker, M. (2012). "De novo genome assembly: what every biologist should
3871 know". In: *Nature Methods* 9(4), pp. 333–337.
- 3872 Baranwal, V. K., V. Mikkilineni, U. B. Zehr, A. K. Tyagi, and S. Kapoor
3873 (2012). "Heterosis: emerging ideas about hybrid vigour". In: *Journal of*
3874 *Experimental Botany* 63(2), pp. 695–709.
- 3875 Barrick, J. E. and R. E. Lenski (2013). "Genome dynamics during experi-
3876 mental evolution." In: *Nature Reviews Genetics* 14(12), pp. 827–39.
- 3877 Barrick, J. E., D. S. Yu, S. H. Yoon, H. Jeong, T. K. Oh, D. Schneider,
3878 R. E. Lenski, and J. F. Kim (2009). "Genome evolution and adaptation
3879 in a long-term experiment with *Escherichia coli*." In: *Nature* 461(7268),
3880 pp. 1243–7.
- 3881 Barton, N. H., D. E. G Briggs, J. A. Eisen, D. B. Goldstein, and N. H. Patel
3882 (2007). *Evolution*. Cold Springs Harbor, NY: Cold Springs Harbor Press,
3883 p. 833.
- 3884 Bayer-Giraldi, M., I. Weikusat, H. Besir, and G. Dieckmann (2011). "Char-
3885 acterization of an antifreeze protein from the polar diatom *Fragilariopsis*
3886 *cylindrus* and its relevance in sea ice". In: *Cryobiology* 63(3), pp. 210–
3887 219.
- 3888 Bazykin, A. D. (1969). "Hypothetical mechanism of speciation". In: *Evolution*
3889 23(4), pp. 685–687.
- 3890 Beck, S. and J. Trowsdale (2000). "The human Major Histocompatibility
3891 Complex: Lessons from the DNA sequence". In: *Annual Review of*
3892 *Genomics and Human Genetics* 1, pp. 117–137.
- 3893 Belhaj, K, L. M. Cano, C Prince D, A Kemen, K Yoshida, Y. F. Dagdas, G. J.
3894 Ehterington, H Schoonbeek, H. P. van Esse, J. D. G. Jones, S Kamoun,

- 3895 and S Schornack (2015). “*Arabidopsis* late blight: Infection of a nonhost
3896 plant by *Albugo laibachii* enables full colonization by *Phytophthora*
3897 *infestans*”. In: *bioRxiv* (June), pp. 1–30.
- 3898 Bernstein, H., C. Bernstein, and R. E. Michod (2011). “Meiosis as an
3899 evolutionary adaptation for DNA repair.” In: *DNA repair*. Ed. by Inna
3900 Kruman. Intech. Chap. 19, pp. 357–382.
- 3901 Bernstein, H, F. A. Hopf, and R. E. Michod (1987). “The molecular basis of
3902 the evolution of sex”. In: *Advances in Genetics* 24, pp. 323–370.
- 3903 Bernstein, H. (1985). “Genetic damage, mutation, and the evolution of sex”.
3904 In: *Science* 229, pp. 1277–1281.
- 3905 Bernstein, H., H. C. Byerly, F. A. Hopf, and R. E. Michod (1984). “Origin of
3906 sex”. In: *Journal of Theoretical Biology* 110(3), pp. 323–351.
- 3907 Biémont, C and C. Vieira (2006). “Junk DNA as an evolutionary force”. In:
3908 *Nature Genetics* 43, pp. 521–524.
- 3909 Biga, M. L. B. (1955). “Review of the species of the genus *Albugo* based
3910 on the morphology of the conidia.” In: *Sydowia* 9, pp. 339–358.
- 3911 Bodmer, W. F. and L. L. Cavalli-Sforza (1968). “A migration matrix model
3912 for the study of random genetic drift.” In: *Genetics* 59(4), pp. 565–592.
- 3913 Boni, M. F., D. Posada, and M. W. Feldman (2007). “An exact nonparametric
3914 method for inferring mosaic structure in sequence triplets”. In: *Genetics*
3915 176(2), pp. 1035–1047.
- 3916 Bouckaert, R., J. Heled, D. Kühnert, T. Vaughan, C.-H. Wu, D. Xie, M. A.
3917 Suchard, A. Rambaut, and A. J. Drummond (2014). “BEAST 2: A soft-
3918 ware platform for bayesian evolutionary analysis”. In: *PLoS Computa-*
3919 *tional Biology* 10(4), e1003537.
- 3920 Boutemy, L. S., S. R. F. King, J. Win, R. K. Hughes, T. a. Clarke, T. M. a.
3921 Blumenschein, S. Kamoun, and M. J. Banfield (2011). “Structures of
3922 *Phytophthora* RXLR effector proteins: a conserved but adaptable fold

- 3923 underpins functional diversity.” In: *The Journal of Biological Chemistry*
3924 286(41), pp. 35834–42.
- 3925 Bowler, C., A. Vardi, and A. E. Allen (2010). “Oceanographic and biogeo-
3926 chemical insights from diatom genomes.” In: *Annual Review of Marine*
3927 *Science* 2, pp. 333–365.
- 3928 Brent, R. (1973). *Algorithms for minimization without derivatives*. Mineola,
3929 New York: Dover Publications, p. 206.
- 3930 Bromham, L. and D. Penny (2003). “The modern molecular clock.” In:
3931 *Nature Reviews Genetics* 4, pp. 216–224.
- 3932 Brown, J. K. M. (2003). “A cost of disease resistance: Paradigm or peculiar-
3933 ity?” In: *Trends in Genetics* 19(12), pp. 667–671.
- 3934 Brown, R. P. and Z. Yang (2011). “Rate variation and estimation of di-
3935 vergence times using strict and relaxed clocks”. In: *BMC Evolutionary*
3936 *Biology* 11(1), p. 271.
- 3937 Bruen, T. C., H. Philippe, and D. Bryant (2006). “A simple and robust
3938 statistical test for detecting the presence of recombination”. In: *Genetics*
3939 172(April), pp. 2665–2681.
- 3940 Bubb, K. L., D. Bovee, D. Buckley, E. Haugen, M. Kibukawa, M. Paddock, A.
3941 Palmieri, S. Subramanian, Y. Zhou, R. Kaul, P. Green, and M. V. Olson
3942 (2006). “Scan of human genome reveals no new loci under ancient
3943 balancing selection”. In: *Genetics* 173(4), pp. 2165–2177.
- 3944 Buckling, A. and P. B. Rainey (2002). “Antagonistic coevolution between a
3945 bacterium and a bacteriophage.” In: *Proceedings. Biological sciences /*
3946 *The Royal Society* 269(1494), pp. 931–6.
- 3947 Buerkle, C. A., R. J. Morris, M. A. Asmussen, and L. H. Rieseberg (2000).
3948 “The likelihood of homoploid hybrid speciation.” In: *Heredity* 84, pp. 441–
3949 451.
- 3950 Buerkle, C. A. and L. H. Rieseberg (2008). “The rate of genome stabilization
3951 in homoploid hybrid species”. In: *Evolution* 62(2), pp. 266–275.

- 3952 Burke, J. M., T. J. Voss, and M. L. Arnold (1998). "Genetic interactions and
3953 natural selection in Louisiana Iris hybrids". In: *Evolution* 52(5), pp. 1304–
3954 1310.
- 3955 Burke, J. M. and M. L. Arnold (2001). "Genetics and the fitness of hybrids".
3956 In: *Annual Review of Genetics* 35, pp. 31–52.
- 3957 Burton, R. S. (1990a). "Hybrid Breakdown in developmental time in the
3958 Copepod *Tigriopus californicus*". In: *Evolution* 44(7), pp. 1814–1822.
- 3959 — (1990b). "Hybrid breakdown in physiological response: a mechanistic
3960 approach". In: *Evolution* 44(7), pp. 1806–1813.
- 3961 Burton, R. S., P. D. Rawson, and S. Edmands (1999). "Genetic architecture
3962 of physiological phenotypes: Empirical evidence for coadapted gene
3963 complexes". In: *American Zoologist* 39(2), pp. 451–462.
- 3964 Cabot, E. L., A. W. Davis, N. A. Johnson, and C. I. Wu (1994). "Genetics
3965 of reproductive isolation in the *Drosophila simulans* clade: Complex
3966 epistasis underlying hybrid male sterility". In: *Genetics* 137(1), pp. 175–
3967 189.
- 3968 Cambareri, E. B., M. J. Singer, and E. U. Selker (1991). "Recurrence of
3969 repeat-induced point mutation (RIP) in *Neurospora crassa*". In: *Genetics*
3970 127(4), pp. 699–710.
- 3971 Casabianca, S., A. Penna, E. Pecchioli, A. Jordi, G. Basterretxea, and C.
3972 Vernesi (2012). "Population genetic structure and connectivity of the
3973 harmful dinoflagellate *Alexandrium minutum* in the Mediterranean Sea."
3974 In: *Proceedings of the Royal Society B: Biological Sciences* 279(1726),
3975 pp. 129–38.
- 3976 Cavalier-Smith, T (2002). "Origins of the machinery of recombination and
3977 sex." In: *Heredity* 88, pp. 125–141.
- 3978 Chakraborty, R. and M. Nei (1977). "Bottleneck effects on average het-
3979 erozygosity and genetic distance with the stepwise mutation model". In:
3980 *Evolution* 31(2), pp. 347–356.

- 3981 Chappell, T. M. and M. D. Rausher (2016). "Evolution of host range in
3982 *Coleosporium ipomoeae*, a plant pathogen with multiple hosts". In: *Pro-*
3983 *ceedings of the National Academy of Sciences* 113(19), p. 201522997.
- 3984 Charlesworth, B. and D. Charlesworth (2010). *Elements of Evolutionary*
3985 *Genetics*. Greenwood Village, Colorado: Roberts and Company, p. 734.
- 3986 Charlesworth, D. (2006). "Balancing selection and its effects on sequences
3987 in nearby genome regions". In: *PLoS Genetics* 2(4), pp. 379–384.
- 3988 Charlesworth, D. and J. H. Willis (2009). "The genetics of inbreeding de-
3989 pression." In: *Nature Reviews Genetics* 10(11), pp. 783–96.
- 3990 Chepurinov, V. A., D. G. Mann, K. Sabbe, and W. Vyverman (2004). "Ex-
3991 perimental studies on sexual reproduction in diatoms". In: *International*
3992 *Review of Cytology* 237, pp. 91–154.
- 3993 Choi, D and M. J. Priest (1955). "A key to the genus *Albugo*". In: *Mycotaxon*
3994 53, pp. 261–272.
- 3995 Choi, Y.-J., H.-D. Shin, and M. Thines (2009). "The host range of *Albugo*
3996 *candida* extends from *Brassicaceae* through *Cleomaceae* to *Cappa-*
3997 *raceae*". In: *Mycological Progress* 8(4), pp. 329–335.
- 3998 Choi, Y.-J., H.-D. Shin, S. Ploch, and M. Thines (2011). "Three new phy-
3999 logenetic lineages are the closest relatives of the widespread species
4000 *Albugo candida*." In: *Fungal biology* 115(7), pp. 598–607.
- 4001 Clark Cockerham, C. and Z. B. Zeng (1996). "Design III with marker loci".
4002 In: *Genetics* 143(3), pp. 1437–1456.
- 4003 Coop, G. and M. Przeworski (2007). "An evolutionary view of human recom-
4004 bination." In: *Nature Reviews Genetics* 8(1), pp. 23–34.
- 4005 Cooper, A. J., A. O. Latunde-Dada, A Woods-Tör, J Lynn, J. A. Lucas, I. R.
4006 Crute, and E. B. Holub (2008). "Basic compatibility of *Albugo candida*
4007 in *Arabidopsis thaliana* and *Brassica juncea* causes broad-spectrum
4008 suppression of innate immunity." In: *Molecular plant-microbe interactions*
4009 : *MPMI* 21(6), pp. 745–756.

- 4010 Cousyn, C, L De Meester, J. K. Colbourne, L Brendonck, D Verschuren, and
4011 F Volckaert (2001). "Rapid, local adaptation of zooplankton behavior
4012 to changes in predation pressure in the absence of neutral genetic
4013 changes." In: *Proceedings of the National Academy of Sciences of the
4014 United States of America* 98(11), pp. 6256–6260.
- 4015 Crow, J. F. (1994). "Advantages of sexual reproduction". In: *Developmental
4016 Genetics* 15(3), pp. 205–213.
- 4017 Crow, J. F. and M. Kimura (1970). *An Introduction to Population Genetics
4018 Theory*. New York: Harper and Row, p. 656.
- 4019 Cruzan, M. B. and M. L. Arnold (1994). "Assortative Mating and Natural-
4020 Selection in an Iris Hybrid Zone". In: *Evolution* 48(6), pp. 1946–1958.
- 4021 Darlington, C. D. and K. Mather (1950). *The Elements of Genetics*. New
4022 York: Macmillan, p. 462.
- 4023 Davidovich, N. A. and S. S. Bates (1998). "Patterns of sexual reproduction
4024 in diatoms". In: *Hydrobiologia* 269-270, pp. 11–20.
- 4025 Davis, B and B. Columbia (2011). "Genetic load". In: *eLS* July, pp. 1–4.
- 4026 Dawson, E., G. R. Abecasis, S. Bumpstead, Y. Chen, S. Hunt, D. M. Beare,
4027 J. Pabial, T. Dibling, E. Tinsley, S. Kirby, D. Carter, M. Papaspyridonos,
4028 S. Livingstone, R. Ganske, E. Lohmussaar, J. Zernant, N. Tonisson,
4029 M. Remm, R. Magi, T. Puurand, J. Vilo, A. Kurg, K. Rice, P. Deloukas,
4030 R. Mott, A. Metspalu, D. R. Bentley, L. R. Cardon, and I. Dunham (2002).
4031 "A first-generation linkage disequilibrium map of human chromosome
4032 22". In: *Nature* 418, pp. 544–548.
- 4033 Decaestecker, E., S. Gaba, J. A. M. Raeymaekers, R. Stoks, L. Van Kerck-
4034 hoven, D. Ebert, and L. De Meester (2007). "Host-parasite 'Red Queen'
4035 dynamics archived in pond sediment." In: *Nature* 450(7171), pp. 870–
4036 873.

- 4037 Defaveri, J. and J. Meril (2014). “Local adaptation to salinity in the three-
4038 spined stickleback?” In: *Journal of Evolutionary Biology* 27(2), pp. 290–
4039 302.
- 4040 Delaneau, O., J. Marchini, and J.-F. Zagury (2012). “A linear complexity
4041 phasing method for thousands of genomes.” In: *Nature Methods* 9(2),
4042 pp. 179–81.
- 4043 Delaneau, O., J.-F. Zagury, and J. Marchini (2013). “Improved whole-
4044 chromosome phasing for disease and population genetic studies.” In:
4045 *Nature Methods* 10(1), pp. 5–6.
- 4046 Delaneau, O., B. Howie, A. J. Cox, J.-F. Zagury, and J. Marchini (2013).
4047 “Haplotype estimation using sequencing reads”. In: *American Journal of*
4048 *Human Genetics* 93(4), pp. 687–696.
- 4049 Depotter, J. R. L., M. F. Seidl, T. A. Wood, and B. P.H. J. Thomma (2016).
4050 “Interspecific hybridization impacts host range and pathogenicity of fil-
4051 amentous microbes”. In: *Current Opinion in Microbiology* 32, pp. 7–
4052 13.
- 4053 Devlin, B. and N. Risch (1995). “Linkage disequilibrium measures for fine-
4054 scale mapping: a comparison”. In: *Genomics* 29, pp. 311–322.
- 4055 Diogo, D., C. Bouchier, C. D’Enfert, and M. E. Bougnoux (2009). “Loss
4056 of heterozygosity in commensal isolates of the asexual diploid yeast
4057 *Candida albicans*”. In: *Fungal Genetics and Biology* 46(2), pp. 159–168.
- 4058 Dobzhansky, T (1936). “Studies on hybrid sterility. II. Localization of steril-
4059 ity factors in *Drosophila Pseudoobscura* Hybrids.” In: *Genetics* 21(2),
4060 pp. 113–135.
- 4061 — (1970). *Genetics of the Evolutionary Process*. A Columbia Paperback.
4062 New York and London: Columbia University Press, p. 506.
- 4063 Dong, S., S. Raffaele, and S. Kamoun (2015). “The two-speed genomes
4064 of filamentous pathogens: Waltz with plants”. In: *Current Opinion in*
4065 *Genetics and Development* 35, pp. 57–65.

- 4066 Dong, S., R. Stam, L. M. Cano, J. Song, J. Sklenar, K. Yoshida, T. O.
4067 Bozkurt, R. Oliva, Z. Liu, M. Tian, J. Win, M. J. Banfield, A. M. E. Jones,
4068 R. A. L. Van der Hoorn, and S. Kamoun (2014). "Effector Specializa-
4069 tion in a Lineage of the Irish Potato Famine Pathogen". In: *Science*
4070 343(6170), pp. 552–555.
- 4071 Dres, M. and J. Mallet (2002). "Host races in plant-feeding insects and their
4072 importance in sympatric speciation". In: *Philosophical Transactions of*
4073 *the Royal Society B: Biological Sciences* 357(1420), pp. 471–492.
- 4074 Drummond, A. and A. Rambaut (2007). "BEAST: Bayesian evolutionary
4075 analysis by sampling trees". In: *BMC Evolutionary Biology* 7(1), p. 214.
- 4076 Drummond, A. J. and M. A. Suchard (2010). "Bayesian random local clocks,
4077 or one rate to rule them all." In: *BMC biology* 8, p. 114.
- 4078 Drummond, A. J., M. A. Suchard, D. Xie, and A. Rambaut (2012). "Bayesian
4079 phylogenetics with BEAUti and the BEAST 1.7". In: *Molecular Biology*
4080 *and Evolution* 29(8), pp. 1969–1973.
- 4081 Dujon, B. (2010). "Yeast evolutionary genomics". In: *Nature Reviews Ge-*
4082 *netics* 11(7), pp. 512–524.
- 4083 Durand, E. Y., N. Patterson, D. Reich, and M. Slatkin (2011). "Testing for
4084 ancient admixture between closely related populations." In: *Molecular*
4085 *biology and evolution* 28(8), pp. 2239–52.
- 4086 Edmands, S. (1999). "Heterosis and outbreeding depression in interpopula-
4087 tion crosses spanning a wide range of divergence". In: *Evolution* 53(6),
4088 pp. 1757–1768.
- 4089 Edwards, S. V. and P. W. Hedrick (1998). "Evolution and ecology of MHC
4090 molecules: From genomics to sexual selection". In: *Trends in Ecology*
4091 *and Evolution* 13(8), pp. 305–311.
- 4092 Eisen, J. A. (2000). "Horizontal gene transfer among microbial genomes:
4093 new insights from complete genome analysis." In: *Current Opinion in*
4094 *Genetics and Development* 10(6), pp. 606–11.

- 4095 Endler, J. A. (1977). "Geographic variation, speciation, and clines." In:
4096 *Monographs in Population Biology* 10, pp. 1–246.
- 4097 Endler, J. A. (1982). "Problems in distinguishing historical from ecological
4098 factors in biogeography". In: *American Zoologist* 22(2), pp. 441–452.
- 4099 Ersek, T, J. T. English, and J. E. Schoelz (1995). "Creation of species
4100 hybrids of *Phytophthora* with modified host ranges by zoospore fusion".
4101 In: *Phytopathology* 85(11), pp. 1343–1347.
- 4102 Felsenstein, J. and S. Yokoyama (1976). "The evolutionary advantage of
4103 recombination. II. Individual selection for recombination". In: *Genetics*
4104 83(4), pp. 845–859.
- 4105 Felsenstein, J. (1981). "Evolutionary trees from DNA sequences: A maxi-
4106 mum likelihood approach". English. In: *Journal of Molecular Evolution*
4107 17(6), pp. 368–376.
- 4108 Fisher, R. A. (1930). *The Genetical Theory of Natural Selection*. Oxford: At
4109 The Clarendon Press, p. 308.
- 4110 Fontaine, M. C., J. B. Pease, A. Steele, R. M. Waterhouse, D. E. Neafsey,
4111 I. V. Sharakhov, X. Jiang, A. B. Hall, F. Catteruccia, E. Kakani, S. N.
4112 Mitchell, Y. Wu, H. A. Smith, R. R. Love, M. K. N. Lawniczak, M. A.
4113 Slotman, S. J. Emrich, M. W. Hahn, and N. J. Besansky (2015). "Ex-
4114 tensive introgression in a malaria vector species complex revealed by
4115 phylogenomics". In: *Science* 347(6217), pp. 12585241–12585246.
- 4116 Ford Doolittle, W. (1999). "Lateral genomics". In: *Trends in Biochemical*
4117 *Sciences* 24(12), pp. 5–8.
- 4118 Frank, S (1993). "Coevolutionary genetics of plant pathogens". In: *Evolu-*
4119 *tionay Ecology* 7, pp. 45–75.
- 4120 Frankham, R (1996). "Relationship of genetic variation to population size in
4121 wildlife". In: *Conservation Biology* 10(6), pp. 1500–1508.
- 4122 Gallagher, J. (1983). "Cell Enlargement in *Skeletonema Costatum* (*Bacillar-*
4123 *iophyceae*)". In: *Journal of phycology* 19(4), pp. 539–542.

- 4124 Galtier, N. and V. Daubin (2008). “Dealing with incongruence in phyloge-
4125 nomic analyses.” In: *Philosophical Transactions of the Royal Society of*
4126 *London. Series B, Biological Sciences* 363(1512), pp. 4023–4029.
- 4127 Gardner, K., C. A. Buerkle, J. Whitton, and L. H. Rieseberg (2000). “Infer-
4128 ring Epistasis in Wild Sunflower Hybrid Zones”. In: *Epistasis and the*
4129 *Evolutionary Process*. Ed. by J. B. Wolf, E. D. Brodie, and M. J. Wade.
4130 Oxford: Oxford University Press. Chap. 16, pp. 264–279.
- 4131 Gerrish, P. J. and R. E. Lenski (1998). “The fate of competing beneficial
4132 mutations in an asexual population.” In: *Genetica* 102-103(1-6), pp. 127–
4133 44.
- 4134 Gijzen, M., C. Ishmael, and S. D. Shrestha (2014). “Epigenetic control of
4135 effectors in plant pathogens”. In: *Frontiers in Plant Science* 5, pp. 1–4.
- 4136 Giraud, A, M Radman, I Matic, and F Taddei (2001). “The rise and fall of
4137 mutator bacteria.” In: *Current Opinion in Microbiology* 4(5), pp. 582–5.
- 4138 Goodwin, S, L. S. Sujkowski, A. T. Dyer, B. A. Fry, and W. E. Fry (1995).
4139 “Direct detection of gene flow and probably sexual reproduction of *Phy-*
4140 *tophthora infestans* in northern North America”. In: *Phytopathology* 85,
4141 pp. 473–479.
- 4142 Goodwin, S. B., B. A. Cohen, and W. E. Fry (1994). “Panglobal distribution of
4143 a single clonal lineage of the Irish potato famine fungus.” In: *Proceedings*
4144 *of the National Academy of Sciences of the United States of America*
4145 91(24), pp. 11591–5.
- 4146 Goodwin, S. B., L. J. Spielman, J. M. Matuszak, S. M. Bergeron, and W. E.
4147 Fry (1992). “Clonal diversity and genetic differentiation of *Phytophthora*
4148 *infestans* populations in northern and central Mexico”. In: *Phytopathol-*
4149 *ogy* 82, pp. 955–961.
- 4150 Goodwin, S., J. D. McPherson, and W. R. McCombie (2016). “Coming of
4151 age: ten years of next-generation sequencing technologies”. In: *Nature*
4152 *Reviews Genetics* 17(6), pp. 333–351.

- 4153 Grauer, D. and W.-H. Li (2000). *Fundamentals of Molecular Evolution*.
4154 Second Edi. Sunderland: Sinauer Associates. Chap. Genes, gen, p. 481.
- 4155 Greaves, J. H., R Redfern, P. B. Ayres, and J. E. Gill (1977). "Warfarin
4156 resistance: a balanced polymorphism in the Norway rat". In: *Genetical
4157 Research* 30(3), pp. 257–63.
- 4158 Green, R. E. et al. (2010). "A draft sequence of the Neandertal genome".
4159 In: *Science* 328(5979), pp. 710–722.
- 4160 Hacquard, S., B. Kracher, T. Maekawa, S. Vernaldi, P. Schulze-Lefert, and
4161 E. Ver Loren van Themaat (2013). "Mosaic genome structure of the
4162 barley powdery mildew pathogen and conservation of transcriptional
4163 programs in divergent hosts." In: *Proceedings of the National Academy
4164 of Sciences of the United States of America* 110(24), pp. 2219–28.
- 4165 Haldane, J. B. S. (1948). "The theory of a cline". In: *Journal of Genetics*
4166 48(3), pp. 277–284.
- 4167 Haldane, J. B. S. and S. D. Jayakar (1963). "Polymorphism due to selection
4168 of varying direction". In: *Journal of Genetics* 58(2), pp. 237–242.
- 4169 Hanski, I. (1998). "Metapopulation dynamics". In: *Nature* 396, pp. 41–49.
- 4170 Hasegawa, M., H. Kishino, and T.-a. Yano (1985). "Dating of the human-
4171 ape splitting by a molecular clock of mitochondrial DNA". In: *Journal of
4172 Molecular Evolution* 22(2), pp. 160–174.
- 4173 Heath, T. A., J. P. Huelsenbeck, and T. Stadler (2014). "The fossilized birth-
4174 death process for coherent calibration of divergence-time estimates". In:
4175 *Proceedings of the National Academy of Sciences* 111, pp. 2957–2966.
- 4176 Hedrick, P. W. (1980). "Hitchhiking: a comparison of linkage and partial
4177 selfing." In: *Genetics* 94(3), pp. 791–808.
- 4178 — (1987). "Gametic disequilibrium measures: proceed with caution." In:
4179 *Genetics* 117(2), pp. 331–341.

- 4180 Hedrick, P. W. and M. E. Gilpin (1997). "Genetic effective size of a metapopulation". In: *Metapopulation Dynamics: Ecology, Genetics, and Evolution*.
4181 Ed. by I. Hanski and M. Gilpin. New York: Academic Press, pp. 166–182.
4182 Hedrick, P. W. (2010). *The genetics of populations*. 4th. Sudbury: Jones
4183 and Bartlett Publishers, p. 675.
4184 — (2013). "Adaptive introgression in animals: examples and comparison to
4185 new mutation and standing variation as sources of adaptive variation".
4186 In: *Molecular Ecology* 22(18), pp. 4606–4618.
4187 Hein, J., M. Schierup, and C. Wiuf (2004). *Gene Genealogies, Variation and
4188 Evolution: A primer in coalescent theory*. Oxford and New York: Oxford
4189 University Press, p. 296.
4190 Hewitt, G. M. (1985). "Analysis of hybrid zones". In: *Annual Review of
4191 Ecology, Evolution, and Systematics* 16(1985), pp. 113–148.
4192 Hill, W. G. and A. Robertson (1966). "The effect of linkage on artificial
4193 selection." In: *Genetics Research* 8, pp. 269–294.
4194 Hill, W. G. and A. Robertson (1968). "Linkage disequilibrium in finite popula-
4195 tions." In: *Theoretical and applied genetics*. 38(6), pp. 226–31.
4196 Hoekstra, H. E. and J. A. Coyne (2007). "The locus of evolution: Evo devo
4197 and the genetics of adaptation". In: *Evolution* 61(5), pp. 995–1016.
4198 Hollocher, H. and C. I. Wu (1996). "The genetics of reproductive isolation
4199 in the *Drosophila simulans* clade: X vs. autosomal effects and male vs.
4200 female effects". In: *Genetics* 143(3), pp. 1243–1255.
4201 Holub, E. B., E. Brose, M. Tör, C. Clay, I. R. Crute, J. L. Beynon, M. Tor, C
4202 Clay, I. R. Crute, and J. L. Beynon (1995). "Phenotypic and genotypic
4203 variation in the interaction between *Arabidopsis thaliana* and *Albugo
4204 candida*." In: *Molecular Plant-Microbe Interactions : MPMI* 8(6), pp. 916–
4205 928.
4206 Howie, B., J. Marchini, and M. Stephens (2011). "Genotype imputation with
4207 thousands of genomes". In: *G3* 1(6), pp. 457–470.
4208

- 4209 Hudson, R. R. and N. L. Kaplan (1988). "The coalescent process in models
4210 with selection and recombination." In: *Genetics* 120(3), pp. 831–40.
- 4211 Hudson, R. R. (2001). "Two-locus sampling distributions and their applica-
4212 tion". In: *Genetics* 159, pp. 1805–1817.
- 4213 Huson, D. H. (1998). "SplitsTree: analyzing and visualizing evolutionary
4214 data." In: *Bioinformatics* 14(1), pp. 68–73.
- 4215 Huson, D. H. and D. Bryant (2006). "Application of phylogenetic networks in
4216 evolutionary studies." In: *Molecular Biology and Evolution* 23(2), pp. 254–
4217 67.
- 4218 Ingvarsson, P. K. and M. C. Whitlock (2000). "Heterosis increases the
4219 effective migration rate." In: *Proceedings. Biological sciences / The
4220 Royal Society* 267(1450), pp. 1321–1326.
- 4221 Jayaraman, R (2011). "Phase variation and adaptation in bacteria: A Red
4222 Queen arms race". In: *Current Science* 100(8), pp. 1163–1171.
- 4223 Jex, A. R., M. A. Schneider, and T. H. Cribb (2006). "The importance of
4224 host ecology in thelastomatoid (Nematoda: Oxyurida) host specificity."
4225 In: *Parasitology International* 55(3), pp. 169–74.
- 4226 Johansen-Morris, A. D. and R. G. Latta (2006). "Fitness consequences of
4227 hybridization between ecotypes of *Avena barbata*: hybrid breakdown, hy-
4228 brid vigor, and transgressive segregation." In: *Evolution* 60(8), pp. 1585–
4229 1595.
- 4230 Jouet, A., M. McMullan, and C. van Oosterhout (2015). "The effects of
4231 recombination, mutation and selection on the evolution of the Rp1 re-
4232 sistance genes in grasses". In: *Molecular Ecology* 24(12), pp. 3077–
4233 3092.
- 4234 Judson, O. P. and B. B. Normark (1996). "Ancient asexual scandals". In:
4235 *Trends in Ecology and Evolution* 11(2), pp. 41–46.

- 4236 Jukes, T. and C. Cantor (1969). "Evolution of Protein Molecules". In: *Mam-*
4237 *malian Protein Metabolism III*. Ed. by H. N. Munro. New York: Academic
4238 Press, pp. 21–132.
- 4239 Kang, S. H. and G. A. Fryxell (1992). "*Fragilariopsis cylindrus* (Grunow)
4240 Krieger: The most abundant diatom in water column assemblages of
4241 Antarctic marginal ice-edge zones". In: *Polar Biology* 12(6-7), pp. 609–
4242 627.
- 4243 Kaplan, N. L., R. R. Hudson, and C. H. Langley (1989). "The 'hitchhiking
4244 effect' revisited". In: *Genetics* 123(4), pp. 887–899.
- 4245 Karrenberg, S., C. Lexer, and L. H. Rieseberg (2007). "Reconstructing
4246 the history of selection during homoploid hybrid speciation." In: *The*
4247 *American Naturalist* 169(6), pp. 725–37.
- 4248 Kauppi, L., A. J. Jeffreys, and S. Keeney (2004). "Where the crossovers are:
4249 recombination distributions in mammals." In: *Nature Reviews Genetics*
4250 5(6), pp. 413–424.
- 4251 Kemen, E. and J. D. G. Jones (2012). "Obligate biotroph parasitism: can
4252 we link genomes to lifestyles?" In: *Trends in Plant Science*, pp. 1–10.
- 4253 Kemen, E., A. Gardiner, T. Schultz-Larsen, A. C. Kemen, A. L. Balmuth,
4254 A. Robert-Seilaniantz, K. Bailey, E. Holub, D. J. Studholme, D. MacLean,
4255 and J. D. G. Jones (2011). "Gene gain and loss during evolution of
4256 obligate parasitism in the white rust pathogen of *Arabidopsis thaliana*".
4257 In: *PLoS Biology* 9(7), e1001094.
- 4258 Kimura, M. (1962). "On the probability of fixation of mutant genes in a
4259 population." In: *Genetics* 47(391), pp. 713–719.
- 4260 Kimura, M (1968). "Evolutionary rate at the molecular level." In: *Nature*
4261 217(5129), pp. 624–626.
- 4262 Kimura, M. and J. F. Crow (1964). "The number of alleles that can be
4263 maintained in a finite population". In: *Genetics* 49, pp. 725–738.

- 4264 Kimura, M and T Ohta (1971). *Theoretical Aspects of Population Genetics*.
4265 Monographs in Population Biology. Princeton, New Jersey: Princeton
4266 University Press.
- 4267 Kimura, M. (1980). "A simple method for estimating evolutionary rates
4268 of base substitutions through comparative studies of nucleotide se-
4269 quences". In: *Journal of Molecular Evolution* 16(2), pp. 111–120.
- 4270 — (1983). *The Neutral Theory of Molecular Evolution*. Cambridge: Cam-
4271 bridge University Press, p. 367.
- 4272 Kimura, M. and J. F. Crow (1963). "The measurement of effective population
4273 number". In: *Evolution* 17(3), pp. 279–288.
- 4274 Koenig, R. L., R. C. Ploetz, and H. C. Kistler (1997). "*Fusarium oxysporum*
4275 f. sp. *cubense* consists of a small number of divergent and globally
4276 distributed clonal lineages." In: *Phytopathology* 87(9), pp. 915–923.
- 4277 Kovach, R. P., B. K. Hand, P. A. Hohenlohe, T. F. Cosart, M. C. Boyer, H. H.
4278 Neville, C. C. Muhlfeld, S. J. Amish, K. Carim, S. R. Narum, W. H. Lowe,
4279 F. W. Allendorf, and G. Luikart (2016). "Vive la résistance: genome-
4280 wide selection against introduced alleles in invasive hybrid zones". In:
4281 *Proceedings of the Royal Society B: Biological Sciences* 283(1843).
- 4282 Kozmin, S, G Slezak, A Reynaud-Angelin, C Elie, Y de Rycke, S Boiteux,
4283 and E Sage (2005). "UVA radiation is highly mutagenic in cells that
4284 are unable to repair 7,8-dihydro-8-oxoguanine in *Saccharomyces cere-*
4285 *visiae*". In: *Proceedings of the National Academy of Sciences of the*
4286 *United States of America* 102(38), pp. 13538–43.
- 4287 Krell, A. (2006). "Salt stress tolerance in the psychrophilic diatom *Fragilari-*
4288 *opsis cylindrus*". PhD thesis. Bremen, p. 124.
- 4289 Kuhner, M. K. (2006). "LAMARC 2.0: maximum likelihood and Bayesian
4290 estimation of population parameters." In: *Bioinformatics* 22(6), pp. 768–
4291 770.

- 4292 Lai, Z., T. Nakazato, M. Salmaso, J. M. Burke, S. Tang, S. J. Knapp, and
4293 L. H. Rieseberg (2005). "Extensive chromosomal repatterning and the
4294 evolution of sterility barriers in hybrid sunflower species". In: *Genetics*
4295 171(1), pp. 291–303.
- 4296 Lamichhaney, S., J. Berglund, M. S. Almén, K. Maqbool, M. Grabherr, A.
4297 Martinez-Barrio, M. Promerová, C.-J. Rubin, C. Wang, N. Zamani, B. R.
4298 Grant, P. R. Grant, M. T. Webster, and L. Andersson (2015). "Evolution
4299 of Darwins finches and their beaks revealed by genome sequencing".
4300 In: *Nature* 518(7539), pp. 371–375.
- 4301 Lamour, A. K. H., J. Mudge, D. Gobena, O. P. Hurtado-gonzales, K Bharti,
4302 R. S. Donahoo, S. Finley, E. Huitema, J. Hulvey, and D. Platt (2010).
4303 "Genome sequencing and mapping reveal loss of heterozygosity as a
4304 mechanism for rapid adaptation in the vegetable pathogen *Phytophthora*
4305 *capsici*." In: *Molecular Plant-Microbe Interactions* 25(10), pp. 1350–60.
- 4306 Lamour, K. H. and S. Kamoun (2009). *Oomycete Genetics and Genomics:*
4307 *Diversity, Interactions and Research Tools*. Chichester: Wiley Blackwell,
4308 p. 582.
- 4309 Lamour, K. H., R. Stam, J. Jupe, E. Huitema, E. Road, and D. Dd (2012).
4310 "Pathogen profile The oomycete broad-host-range pathogen *Phytoph-*
4311 *thora capsici*". In: *Molecular Plant Pathology* 13(4), pp. 329–337.
- 4312 Lande, R. (1994). "Risk of population extinction from fixation of new deleterious
4313 mutations". In: *Evolution* 48(5), pp. 1460–1469.
- 4314 Le Quesne, W. J. (1969). "A method of selection of characters in numerical
4315 taxonomy". In: *Systematic Biology* 18(2), pp. 201–205.
- 4316 Lemey, P, M Salemi, and A. Vandamme (2009). *The Phylogenetic Hand-*
4317 *book*. Cambridge: Cambridge University Press, p. 723.
- 4318 Lemey, P., M. Lott, D. P. Martin, and V. Moulton (2009). "Identifying re-
4319 combinants in human and primate immunodeficiency virus sequence
4320 alignments using quartet scanning." In: *BMC bioinformatics* 10, p. 126.

- 4321 Lenormand, T. (2002). "Gene flow and the limits to natural selection". In:
4322 *Trends in Ecology and Evolution* 17(4), pp. 183–189.
- 4323 Levins, R. (1969). "Some demographic and genetic consequences of envi-
4324 ronmental heterogeneity for biological control". In: *Bull. Entomol. Soc.*
4325 *Am.* 15, pp. 237–240.
- 4326 Lewontin, R. C. (1988). "On measures of gametic disequilibrium". In: *Ge-*
4327 *netics* 120(3), pp. 849–852.
- 4328 Lewontin, R. C. and K Kojima (1960). "The evolutionary dynamics of com-
4329 plex polymorphisms". In: *Evolution* 14(4), pp. 458–472.
- 4330 Li, C. C. (1976). *First Course in Population Genetics*. Pacific Grove, CA:
4331 Boxwood Press, p. 631.
- 4332 Li, H. (2011). "Improving SNP discovery by base alignment quality". In:
4333 *Bioinformatics* 27(8), pp. 1157–1158.
- 4334 Li, H. and R. Durbin (2009). "Fast and accurate short read alignment with
4335 Burrows-Wheeler transform". In: *Bioinformatics* 25(14), pp. 1754–1760.
- 4336 Li, H., B. Handsaker, A. Wysoker, T. Fennell, J. Ruan, N. Homer, G. Marth,
4337 G. Abecasis, and R. Durbin (2009). "The sequence alignment/map
4338 format and SAMtools". In: *Bioinformatics* 25(16), pp. 2078–2079.
- 4339 Li, W.-h. (2008). "Molecular Clocks". In: *eLS*.
- 4340 Lippman, Z. B. and D. Zamir (2007). "Heterosis: revisiting the magic". In:
4341 *Trends in Genetics* 23(2), pp. 60–66.
- 4342 Little, T. J. and P. D. Hebert (1996). "Ancient asexuals: scandal or artifact?"
4343 In: *Trends in Ecology & Evolution* 11(7), p. 296.
- 4344 Löytynoja, A. and N. Goldman (2005). "An algorithm for progressive mul-
4345 tiple alignment of sequences with insertions." In: *Proceedings of the*
4346 *National Academy of Sciences of the United States of America* 102(30),
4347 pp. 10557–62.

- 4348 Lunt, D. H. (2008). "Genetic tests of ancient asexuality in root knot ne-
4349 matodes reveal recent hybrid origins." In: *BMC evolutionary biology* 8,
4350 p. 194.
- 4351 Lynch, M., J. Conery, and R. Burger (1995). "Mutational meltdowns in
4352 sexual populations". In: *Evolution* 49(6), pp. 1067–1080.
- 4353 Lynch, M. and W. Gabriel (1990). "Mutation load and the survival of small
4354 populations". In: *Evolution* 44(7), pp. 1725–1737.
- 4355 Macholán, M., P. Munclinger, M. Šugerková, P. Dufková, B. Bímová, E.
4356 Božíková, J. Zima, and J. Piálek (2007). "Genetic analysis of autosomal
4357 and X-linked markers across a mouse hybrid zone". In: *Evolution* 61(4),
4358 pp. 746–771.
- 4359 Madigan, M., J. Martinko, D. Stahl, and D. Clark (2012). "Genetics of
4360 Bacteria and Archaea". In: *Biology of Microorganisms*. Thirteenth. New
4361 York: Pearson. Chap. 10, pp. 291–316.
- 4362 Mäkinen, V, D Belazzougui, F Cunial, and A. I. Tomescu (2015). *Genome-
4363 Scale Algorithm Design: Biological Sequence Analysis in the Era of
4364 High-Throughput Sequencing*. Cambridge: Cambridge University Press.
- 4365 Mallet, J. (2005). "Hybridization as an invasion of the genome". In: *Trends
4366 in Ecology and Evolution* 20(5), pp. 229–237.
- 4367 Mann, D. G., S. M. McDonald, M. M. Bayer, S. J. M. Droop, V. A. Chepurnov,
4368 R. E. Loke, A Ciobanu, and J. M. H. Du Buf (2004). "The *Sellaphora
4369 pupula* species complex (Bacillariophyceae): morphometric analysis,
4370 ultrastructure and mating data provide evidence for five new species".
4371 In: *Phycologia* 43(4), pp. 459–482.
- 4372 Mann, D. G. (1989). "The species concept in diatoms: Evidence for mor-
4373 phologically distinct, sympatric gamodemes in four epipellic species". In:
4374 *Plant Systematics and Evolution* 164(1-4), pp. 215–237.
- 4375 Marchini, J. and B. Howie (2010). "Genotype imputation for genome-wide
4376 association studies". In: *Nature Reviews Genetics* 11(7), pp. 499–511.

- 4377 Marchini, J., B. Howie, S. Myers, G. McVean, and P. Donnelly (2007). “A new
4378 multipoint method for genome-wide association studies by imputation of
4379 genotypes.” In: *Nature Genetics* 39(7), pp. 906–13.
- 4380 Martin, D and E Rybicki (2000). “RDP: detection of recombination amongst
4381 aligned sequences.” In: *Bioinformatics* 16(6), pp. 562–563.
- 4382 Martin, D. P., B. Murrell, M. Golden, A. Khoosal, and B. Muhire (2015).
4383 “RDP4: Detection and analysis of recombination patterns in virus genomes”.
4384 In: *Virus Evolution* 1(1), pp. 1–5.
- 4385 Martin, F. and S. Kamoun (2012). *Effectors in plant-microbe interactions*.
4386 Chichester: Wiley Blackwell, p. 444.
- 4387 Martin, S. H., J. W. Davey, and C. D. Jiggins (2014). “Evaluating the use of
4388 ABBA-BABA statistics to locate introgressed loci.” In: *Molecular Biology
4389 and Evolution* 32(1), pp. 244–257.
- 4390 Martin, S. H., K. K. Dasmahapatra, N. J. Nadeau, C. Salazar, J. R. Walters,
4391 F. Simpson, M. Blaxter, A. Manica, J. Mallet, and C. D. Jiggins (2013).
4392 “Genome-wide evidence for speciation with gene flow in *Heliconius*
4393 butterflies”. In: *Genome Research* 23, pp. 1817–1828.
- 4394 Martinsen, G. D., T. G. Whitham, R. J. Turek, and P. Keim (2001). “Hybrid
4395 populations selectively filter gene introgression between species.” In:
4396 *Evolution* 55(7), pp. 1325–1335.
- 4397 Maruyama, T. (1970). “On the rate of decrease of heterozygosity in circu-
4398 lar stepping stone models of populations”. In: *Theoretical Population
4399 Biology* 1(1), pp. 101–119.
- 4400 May, R. M. and R. M. Anderson (1983). “Epidemiology and genetics in
4401 the coevolution of parasites and hosts.” In: *Proceedings of the Royal
4402 Society of London. Series B* 219(1216), pp. 281–313.
- 4403 Maynard-Smith, J and J Haigh (1974). “Hitch-hiking effect of a favorable
4404 gene”. In: *Genetics Research* 23(1), pp. 23–35.

- 4405 Maynard Smith, J and N. H. Smith (1998). "Detecting recombination from
4406 gene trees." In: *Molecular Biology and Evolution* 15(2), pp. 590–599.
- 4407 Mayr, E (1942). *Systematics and the Origin of Species from the Viewpoint*
4408 *of a Zoologist*. New York: Coloumbia University Press, p. 334.
- 4409 McMahon, M. S., C. W. Sham, and D. K. Bishop (2007). "Synthesis-dependent
4410 strand annealing in meiosis." In: *PLoS biology* 5(11), e299.
- 4411 McMullan, M., A. Gardiner, K. Bailey, E. Kemen, B. B. J. Ward, V. Cevik,
4412 A. Robert-Seilaniantz, T. Schultz-Larsen, A. Balmuth, E. Holub, C. van
4413 Oosterhout, and J. D. J. Jones (2015). "Evidence for suppression of
4414 immunity as a driver for genomic introgressions and host range expan-
4415 sion in races of *Albugo candida*, a generalist parasite." In: *eLife* 4(4),
4416 pp. 1–24.
- 4417 McVean, G., P. Awadalla, and P. Fearnhead (2002). "A coalescent-based
4418 method for detecting and estimating recombination from gene sequences".
4419 In: *Genetics* 160(March), pp. 1231–1241.
- 4420 Meena, P., C Chattopadhyay, F. Singh, B. Singh, and A. Gupta (2002).
4421 "Yield loss in indian mustard due to white rust and effect of some cultural
4422 practices on Alternaria blight and white rust severity". In: *Brassica* 4(1 &
4423 2), pp. 18–24.
- 4424 Mehrabi, R., A. H. Bahkali, K. A. Abd-Elsalam, M. Moslem, S. Ben M'Barek,
4425 A. M. Gohari, M. K. Jashni, I. Stergiopoulos, G. H. J. Kema, and P. J.G. M.
4426 De Wit (2011). "Horizontal gene and chromosome transfer in plant
4427 pathogenic fungi affecting host range". In: *FEMS Microbiology Reviews*
4428 35(3), pp. 542–554.
- 4429 Meselson, M. and D. M. Welch (2007). "Stable heterozygosity?" In: *Science*
4430 318(October), pp. 202–204.
- 4431 Metzgar, D., T. Scripps, and L. Jolla (2007). "Mutation Rates: Evolution". In:
4432 *eLS*, pp. 1–3.

- 4433 Mielczarek, M. and J. Szyda (2016). "Review of alignment and SNP calling
4434 algorithms for next-generation sequencing data". In: *Journal of Applied*
4435 *Genetics* 57(1), pp. 71–79.
- 4436 Mimitou, E. P. and L. S. Symington (2009). "Nucleases and helicases take
4437 center stage in homologous recombination". In: *Trends in Biochemical*
4438 *Sciences* 34(5), pp. 264–272.
- 4439 Mock, T., R. P. Otiillar, J. Strauss, A. E. Allen, C. L. Dupont, S. Frickenhaus,
4440 F. Maumus, M. McMullan, A. Salamov, R. Sanges, S. Schmutz, A.
4441 Toseland, A. Veluchamy, B. J. Ward, T. Wu, K. W. Barry, A. Falciatore,
4442 M. I. Ferrante, A. E. Fortunato, G. Glockner, A. Gruber, R. Hipkin, M. G.
4443 Janech, P. Kroth, F. Leese, E. A. Lindquist, B. R. Lyon, J. Martin, P.
4444 Mayer, M. Parker, H. Quesneville, J. A. Raymond, C. Uhlig, R. E. Valas,
4445 K. U. Valentin, A. Z. Worden, E. V. Armbrust, C. Bowler, B. R. Green,
4446 C. van Oosterhout, and I. V. Grigoriev (0). "Extensive genetic diversity
4447 and differential bi-allelic expression in a Southern Ocean diatom". In:
4448 *Nature*.
- 4449 Mock, T., R. P. Otiillar, J. Strauss, M. McMullan, P. Paajanen, J. Schmutz,
4450 A. Salamov, R. Sanges, A. Toseland, B. J. Ward, A. E. Allen, C. L.
4451 Dupont, S. Frickenhaus, F. Maumus, A. Veluchamy, T. Wu, K. W. Barry,
4452 A. Falciatore, M. I. Ferrante, A. E. Fortunato, G. Glöckner, A. Gruber,
4453 R. Hipkin, M. G. Janech, P. G. Kroth, F. Leese, E. A. Lindquist, B. R.
4454 Lyon, J. Martin, C. Mayer, M. Parker, H. Quesneville, J. A. Raymond,
4455 C. Uhlig, R. E. Valas, K. U. Valentin, A. Z. Worden, E. V. Armbrust, M. D.
4456 Clark, C. Bowler, B. R. Green, V. Moulton, C. van Oosterhout, and I. V.
4457 Grigoriev (2017). "Evolutionary genomics of the cold-adapted diatom
4458 *Fragilariopsis cylindrus*". In: *Nature* 541(7638), pp. 536–540.
- 4459 Morgan, W. and S. Kamoun (2007). "RXLR effectors of plant pathogenic
4460 oomycetes." In: *Current Opinion in Microbiology* 10(4), pp. 332–8.

- 4461 Morrell, P. L., T. D. Williams-Coplin, A. L. Lattu, J. E. Bowers, J. M. Chandler,
4462 and A. H. Paterson (2005). "Crop-to-weed introgression has impacted
4463 allelic composition of johnsongrass populations with and without recent
4464 exposure to cultivated sorghum". In: *Molecular Ecology* 14(7), pp. 2143–
4465 2154.
- 4466 Muller, H. J. (1932). "Some genetic aspects of sex". In: *American Naturalist*
4467 66, pp. 118–138.
- 4468 Muller, H. J. (1942). "Isolating mechanisms, evolution, and temperature". In:
4469 *Biology Symposium* 6, pp. 71–125.
- 4470 Myerowitz, R. (1997). "Tay-Sachs disease-causing mutations and neutral
4471 polymorphisms in the Hex A gene". In: *Human Mutation* 9(3), pp. 195–
4472 208.
- 4473 Nadeau, N. J., A Whibley, R. T. Jones, J. W. Davey, K. K. Dasmahapatra,
4474 S. W. Baxter, M. A. Quail, M Joron, R. H. Ffrench-Constant, M. L. Blaxter,
4475 J Mallet, and C. D. Jiggins (2012). "Genomic islands of divergence
4476 in hybridizing *Heliconius butterflies* identified by large-scale targeted
4477 sequencing." In: *Philosophical Transactions of the Royal Society B:*
4478 *Biological Sciences*, pp. 343–353.
- 4479 Nader, J (0). "A Genetic Basis for Anthroponosis by Cryptosporidium". In:
4480 *Unpublished*.
- 4481 Nagai, S., Y. Hori, T. Manabe, and I. Imai (1995). "Restoration of cell size
4482 by vegetative cell enlargement in *Coscinodiscus wailesii* (Bacillario-
4483 phyceae)". In: *Phycologia* 34(6), pp. 533–535.
- 4484 Needleman, S. B. and C. D. Wunsch (1970). "A general method applicable
4485 to the search for similarities in the amino acid sequence of two proteins".
4486 In: *Journal of Molecular Biology* 48(3), pp. 443–453.
- 4487 Nei, M (1987). *Molecular Evolutionary Genetics*. New York: Columbia Uni-
4488 versity Press, p. 512.

- 4489 Nei, M and W. H. Li (1973). "Linkage disequilibrium in subdivided popula-
4490 tions." In: *Genetics* 75(1), pp. 213–219.
- 4491 Nei, M. (2005). "Selectionism and neutralism in molecular evolution". In:
4492 *Molecular Biology and Evolution* 22(12), pp. 2318–2342.
- 4493 Nei, M., T. Maruyama, and R. Chakraborty (1975). "The bottleneck effect
4494 and genetic variability in populations". In: *Evolution* 29(1), pp. 1–10.
- 4495 Nielsen, R., J. S. Paul, A. Albrechtsen, and Y. S. Song (2011). "Genotype
4496 and SNP calling from next-generation sequencing data." In: *Nature*
4497 *Reviews Genetics* 12(6), pp. 443–51.
- 4498 Nunney, L. (1999). "The effective size of a hierarchically structured popula-
4499 tion". In: *Evolution* 53(1), pp. 1–10.
- 4500 Ochman, H, J. G. Lawrence, and E. A. Groisman (2000). "Lateral gene
4501 transfer and the nature of bacterial innovation." In: *Nature* 405(6784),
4502 pp. 299–304.
- 4503 O'Connell, J., D. Gurdasani, O. Delaneau, N. Pirastu, S. Ulivi, M. Cocca, M.
4504 Traglia, J. Huang, J. E. Huffman, I. Rudan, R. McQuillan, R. M. Fraser,
4505 H. Campbell, O. Polasek, G. Asiki, K. Ekoru, C. Hayward, A. F. Wright,
4506 V. Vitart, P. Navarro, J. F. Zagury, J. F. Wilson, D. Toniolo, P. Gasparini,
4507 N. Soranzo, M. S. Sandhu, and J. Marchini (2014). "A general approach
4508 for haplotype phasing across the full spectrum of relatedness". In: *PLoS*
4509 *Genetics* 10(4), e1004234.
- 4510 Ohta, T (1973). "Slightly deleterious mutant substitutions in evolution." In:
4511 *Nature* 246(5428), pp. 96–98.
- 4512 Ohta, T. (1971). "Associative overdominance caused by linked detrimental
4513 mutations". In: *Genetics Research* 18, pp. 277–286.
- 4514 Ohta, T. and M. Kimura (1969). "Linkage disequilibrium due to random
4515 genetic drift". In: *Genetics Research* 13(691), pp. 47–55.

- 4516 Oosterhout, C. van (2009). "A new theory of MHC evolution: beyond selec-
4517 tion on the immune genes." In: *Proceedings. Biological Sciences / The*
4518 *Royal Society* 276(1657), pp. 657–65.
- 4519 Padhukasahasram, B., J. D. Wall, P. Marjoram, and M. Nordborg (2006).
4520 "Estimating recombination rates from single-nucleotide polymorphisms
4521 using summary statistics". In: *Genetics* 174(3), pp. 1517–1528.
- 4522 Padidam, M, S Sawyer, and C. M. Fauquet (1999). "Possible emergence
4523 of new geminiviruses by frequent recombination." In: *Virology* 265(2),
4524 pp. 218–225.
- 4525 Palopoli, M. F. and C. I. Wu (1994). "Genetics of hybrid male sterility be-
4526 tween *Drosophila* sibling species: A complex web of epistasis is revealed
4527 in interspecific studies". In: *Genetics* 138, pp. 329–341.
- 4528 Paraskevis, D., K. Deforche, P. Lemey, G. Magiorkinis, A. Hatzakis, and
4529 A. M. Vandamme (2005). "SlidingBayes: Exploring recombination using
4530 a sliding window approach based on Bayesian phylogenetic inference".
4531 In: *Bioinformatics* 21(7), pp. 1274–1275.
- 4532 Paten, B., D. Earl, and N. Nguyen (2011). "Cactus : Algorithms for genome
4533 multiple sequence alignment". In: *Genome Research* 21, pp. 1512–
4534 1528.
- 4535 Paterson, S., T. Vogwill, A. Buckling, R. Benmayor, A. J. Spiers, N. R.
4536 Thomson, M. Quail, F. Smith, D. Walker, B. Libberton, A. Fenton, N.
4537 Hall, and M. a. Brockhurst (2010). "Antagonistic coevolution accelerates
4538 molecular evolution." In: *Nature* 464(7286), pp. 275–278.
- 4539 Payseur, B. A., M. Place, and J. L. Weber (2008). "Linkage Disequilibrium
4540 between STRPs and SNPs across the Human Genome". In: *American*
4541 *Journal of Human Genetics* 82(5), pp. 1039–1050.
- 4542 Pedersen, A. B., S. Altizer, M. Poss, A. A. Cunningham, and C. L. Nunn
4543 (2005). "Patterns of host specificity and transmission among parasites of

- 4544 wild primates". In: *International Journal for Parasitology* 35(6), pp. 647–
4545 657.
- 4546 Peng, B. and M. Kimmel (2005). "simuPOP: A forward-time population
4547 genetics simulation environment". In: *Bioinformatics* 21(18), pp. 3686–
4548 3687.
- 4549 Ploch, S., Y. J. Choi, C. Rost, H. D. Shin, E. Schilling, and M. Thines
4550 (2010). "Evolution of diversity in *Albugo* is driven by high host specificity
4551 and multiple speciation events on closely related Brassicaceae". In:
4552 *Molecular Phylogenetics and Evolution* 57(2), pp. 812–820.
- 4553 Policy, S. D. and D. J. Conway (2001). "Strong diversifying selection on
4554 domains of the *Plasmodium falciparum* apical membrane antigen 1
4555 gene". In: *Genetics* 158(4), pp. 1505–1512.
- 4556 Pondaven, P., O. Ragueneau, P. Tréguer, A. Hauvespre, L. Dezileau, and
4557 J. L. Reyss (2000). "Resolving the opal paradox in the Southern Ocean".
4558 In: *Nature* 405(6783), 168172.
- 4559 Posada, D and K. A. Crandall (2001). "Evaluation of methods for detecting
4560 recombination from DNA sequences: computer simulations." In: *Pro-
4561 ceedings of the National Academy of Sciences of the United States of
4562 America* 98(24), pp. 13757–62.
- 4563 Posada, D., K. A. Crandall, and E. C. Holmes (2002). "Recombination in
4564 evolutionary genomics." In: *Annual Review of Genetics* 36, pp. 75–97.
- 4565 Pouchkina-Stantcheva, N. N., B. M. McGee, C. Boschetti, D. Tolleter, S.
4566 Chakrabortee, A. V. Popova, F. Meersman, D. Macherel, D. K. Hinch, and
4567 A. Tunnacliffe (2007). "Functional divergence of former alleles in an
4568 ancient asexual invertebrate." In: *Science* 318(5848), pp. 268–71.
- 4569 Poullickova, A. (2008). "Morphology, cytology and sexual reproduction in
4570 the aerophytic cave diatom *Luticola dismutica* (Bacillariophyceae)". In:
4571 *Preslia* 80(1), pp. 87–99.

- 4572 Poulin, R and D Mouillot (2003). "Parasite specialization from a phyloge-
4573 netic perspective: a new index of host specificity". In: *Parasitology* 126,
4574 pp. 473–480.
- 4575 — (2005). "Combining phylogenetic and ecological information into a new
4576 index of host specificity." In: *The Journal of Parasitology* 91(3), pp. 511–
4577 514.
- 4578 Poulin, R. (2011). *Evolutionary Ecology of Parasites*. Second Edi. New
4579 Jersey: Princeton University Press, p. 360.
- 4580 Poulin, R. and D. B. Keeney (2008). "Host specificity under molecular and
4581 experimental scrutiny". In: *Trends in Parasitology* 24(1), pp. 24–28.
- 4582 Pritchard, J. K. and M Przeworski (2001). "Linkage disequilibrium in humans:
4583 models and data." In: *American Journal of Human Genetics* 69(1), pp. 1–
4584 14.
- 4585 Proulx, S. R. and P. C. Phillips (2006). "Allelic divergence precedes and
4586 promotes gene duplication." In: *Evolution* 60(2003), pp. 881–892.
- 4587 Ptak, S. E., K. Voelpel, and M. Przeworski (2004). "Insights into recombina-
4588 tion from patterns of linkage disequilibrium in humans". In: *Genetics*
4589 167(1), pp. 387–397.
- 4590 Raffaele, S. and S. Kamoun (2012). "Genome evolution in filamentous plant
4591 pathogens: why bigger can be better". In: *Nature Reviews Microbiology*
4592 (May), pp. 1–14.
- 4593 Raffaele, S., J. Win, L. M. Cano, and S. Kamoun (2010). "Analyses of
4594 genome architecture and gene expression reveal novel candidate vir-
4595 ulence factors in the secretome of *Phytophthora infestans*". In: *BMC*
4596 *Genomics* 11(1), p. 637.
- 4597 Rainey, P. B., A Buckling, R Kassen, and M Travisano (2000). "The emer-
4598 gence and maintenance of diversity: insights from experimental bacterial
4599 populations." In: *Trends in Ecology & Evolution* 15(6), pp. 243–247.

- 4600 Raymond, J. A. and H. J. Kim (2012). "Possible role of horizontal gene
4601 transfer in the colonization of sea ice by Algae". In: *PLoS ONE* 7(5).
- 4602 Ridley, M. (2004). *Evolution*. 3rd. Malden: Blackwell Publishing, p. 751.
- 4603 Rieseberg, L. H., M. A. Archer, and R. K. Wayne (1999). "Transgressive
4604 segregation, adaptation and speciation." In: *Heredity* 83, pp. 363–372.
- 4605 Rieseberg, L. H., B. Sinervo, C. R. Linder, M. C. Ungerer, and D. M. Arias
4606 (1996). "Role of gene interactions in hybrid speciation: Evidence from
4607 ancient and experimental hybrids". In: *Science* 272(5262), pp. 741–745.
- 4608 Rieseberg, L. H. and S. E. Carney (1998). "Plant hybridization". In: *New
4609 Phytologist* 140(4), pp. 599–624.
- 4610 Rieseberg, L. H., J. Whitton, and K. Gardner (1999). "Hybrid zones and
4611 the genetic architecture of a barrier to gene flow between two sunflower
4612 species". In: *Genetics* 152(2), pp. 713–727.
- 4613 Robbins, T. R., L. E. Walker, K. D. Gorospe, S. A. Karl, A. W. Schrey, E. D.
4614 McCoy, and H. R. Mushinsky (2014). "Rise and fall of a hybrid zone:
4615 Implications for the roles of aggression, mate choice, and secondary
4616 succession". In: *Journal of Heredity* 105(2), pp. 226–236.
- 4617 Robertson, A. (1962). "Selection for heterozygotes in small populations." In:
4618 *Genetics* 47(1953), pp. 1291–1300.
- 4619 Rogers, J. and R. A. Gibbs (2014). "Comparative primate genomics:emerging
4620 patterns of genome content and dynamics". In: *Nature Reviews Genetics*
4621 15(5), pp. 347–359.
- 4622 Rose, L. E., P. D. Bittner-Eddy, C. H. Langley, E. B. Holub, R. W. Michelmore,
4623 and J. L. Beynon (2004). "The maintenance of extreme amino acid
4624 diversity at the disease resistance gene, RPP13, in *Arabidopsis thaliana*".
4625 In: *Genetics* 166(3), pp. 1517–1527.
- 4626 Rousset, F. (1997). "Genetic differentiation and estimation of gene flow
4627 from f statistics under isolation by distance". In: *Genetics* 145, pp. 1219–
4628 1228.

- 4629 Ruderfer, D. M., S. C. Pratt, H. S. Seidel, and L. Kruglyak (2006). "Population
4630 genomic analysis of outcrossing and recombination in yeast." In: *Nature*
4631 *Genetics* 38(9), pp. 1077–1081.
- 4632 Ruhe, J., M. Agler, A. Placzek, K. Kramer, I. Finkemeier, and E. Kemen
4633 (2016). "Obligate biotroph pathogens of the genus *Albugo* are better
4634 adapted to active host defense compared to niche competitors". In:
4635 *Frontiers in Plant Science* 7(820), pp. 1–17.
- 4636 Russell, N (1988). "Oswald Avery and the origin of molecular biology." In:
4637 *British Journal for the History of Science* 21, pp. 393–400.
- 4638 Sabbe, K., V. A. Chepurinov, W. Vyverman, and D. G. Mann (2004). "Apomixis
4639 in *Achnanthes* (Bacillariophyceae); development of a model system for
4640 diatom reproductive biology". In: *European Journal of Phycology* 39(3),
4641 pp. 327–341.
- 4642 Saharan, G. S. and P. R. Verma (1992). *White rusts: a review of eco-*
4643 *nomically important species*. Tech. rep. Ottawa, Ontario: International
4644 Development Research Centre.
- 4645 Saharan, G. S., P. R. Verma, P. D. Meena, and A Kumar (2014). *White Rust*
4646 *of Crucifers: Biology, Ecology and Management*. New Delhi: Springer
4647 India.
- 4648 Sanford, E and M. W. Kelly (2011). "Local adaptation in marine inverte-
4649 brates". In: *Annual Review of Marine Science* 3, pp. 509–535.
- 4650 Sanger, F. and A. R. Coulson (1975). "A rapid method for determining
4651 sequences in DNA by primed synthesis with DNA polymerase". In:
4652 *Journal of Molecular Biology* 94(3), pp. 441–448.
- 4653 Sanger, F, S Nicklen, and A. R. Coulson (1977). "DNA sequencing with
4654 chain-terminating inhibitors." In: *Proceedings of the National Academy*
4655 *of Sciences of the United States of America* 74(12), pp. 5463–7.
- 4656 Scally, A., J. Y. Dutheil, L. W. Hillier, G. E. Jordan, I. Goodhead, J. Herrero,
4657 A. Hobolth, T. Lappalainen, T. Mailund, T. Marques, S. Mccarthy, S. H.

- 4658 Montgomery, P. C. Schwalie, Y. A. Tang, M. C. Ward, Y. Xue, B. Yng-
4659 vadottir, C. Alkan, L. N. Andersen, S. Fitzgerald, T. a. Graves, Y. Gu,
4660 P. Heath, and A. Heger (2012). "Insights into hominid evolution from the
4661 gorilla genome sequence". In: *Nature* 483(7388), pp. 169–175.
- 4662 Schilthuizen, M., R. F. Hoekstra, and E. Gittenberger (1999). "Selective
4663 increase of a rare haplotype in a land snail hybrid zone". In: *Proceedings*
4664 *of the Royal Society B: Biological Sciences* 266(1434), p. 2181.
- 4665 Schlupp, I., A. Taebel-Hellwig, and M. Tobler (2010). "Equal fecundity in
4666 asexual and sexual mollies (*Poecilia*)". In: *Environmental Biology of*
4667 *Fishes* 88(2), pp. 201–206.
- 4668 Schurko, A. M., M. Neiman, and J. M. Logsdon (2009). "Signs of sex: what
4669 we know and how we know it". In: *Trends in Ecology and Evolution* 24,
4670 pp. 208–217.
- 4671 Scotti-Saintagne, C., S. Mariette, I. Porth, P. G. Goicoechea, T. Barreneche,
4672 C. Bodénès, K. Burg, and A. Kremer (2004). "Genome scanning for
4673 interspecific differentiation between two closely related oak species
4674 (*Quercus robur* L. and *Q. petraea* (Matt.) Liebl.)" In: *Genetics* 168(3),
4675 pp. 1615–1626.
- 4676 Seger, J and J Antonovics (1988). "Dynamics of some simple host-parasite
4677 models with more than two in each species genotypes". In: *Philosophical*
4678 *Transactions of the Royal Society of London. Series B , Biological*
4679 319(1196), pp. 541–555.
- 4680 Shee, C., J. L. Gibson, and S. M. Rosenberg (2012). "Two mechanisms
4681 produce mutation hotspots at DNA breaks in *Escherichia coli*". In: *Cell*
4682 *Reports* 2(4), pp. 714–721.
- 4683 Shriner, D., D. Nickle, M. Jensen, and J. Mullins (2003). "Potential impact
4684 of recombination on sitewise approaches for detecting positive natural
4685 selection". In: *Genetical Research* 81(2), pp. 115–121.

- 4686 Sildever, S., J. Seibom, I. Lips, and A. Godhe (2016). "Competitive advantage and higher fitness in native populations of genetically structured
4687 planktonic diatoms." In: *Environmental Microbiology* 18(12), pp. 4403–
4688 4411.
- 4690 Slatkin, M. (1995). "A measure of population subdivision based on mi-
4691 crosatellite allele frequencies". In: *Genetics* 139(1), pp. 457–462.
- 4692 Slatkin, M. (1977). "Gene flow and genetic drift in a species subject to
4693 frequent local extinctions". In: *Theoretical Population Biology* 12(3),
4694 pp. 253–262.
- 4695 Slattery, M., H. N. Kamel, S. Ankisetty, D. J. Gochfeld, C. A. Hoover, and
4696 R. W. Thacker (2008). "Hybrid vigor in a tropical pacific soft-coral com-
4697 munity". In: *Ecological Monographs* 78(3), pp. 423–443.
- 4698 Smith, G. R. (2012). "How RecBCD enzyme and Chi promote DNA break
4699 repair and recombination: A molecular biologist's view". In: *Microbiology
4700 and Molecular Biology Reviews* 76(2), pp. 217–228.
- 4701 Smith, J. M. (1992). "Analyzing the mosaic structure of genes." In: *Journal
4702 of Molecular Evolution* 34(2), pp. 126–129.
- 4703 Smith, T. and M. Waterman (1981). "Identification of common molecular
4704 subsequences". In: *Journal of Molecular Biology* 147(1), pp. 195–197.
- 4705 Spurgin, L. G. and D. S. Richardson (2010). "How pathogens drive genetic
4706 diversity: MHC, mechanisms and misunderstandings." In: *Proceedings.
4707 Biological Sciences / The Royal Society* 277, pp. 979–988.
- 4708 Spurgin, L. G., C. van Oosterhout, J. C. Illera, S. Bridgett, K. Gharbi, B. C.
4709 Emerson, and D. S. Richardson (2011). "Gene conversion rapidly gener-
4710 ates major histocompatibility complex diversity in recently founded bird
4711 populations". In: *Molecular Ecology* 20, pp. 5213–5225.

- 4712 Steele, K. A., E. Humphreys, C. R. Wellings, and M. J. Dickinson (2001).
4713 “Support for a stepwise mutation model for pathogen evolution in aus-
4714 tralasian *Puccinia striiformis* f.sp. *tritici* by use of molecular markers”. In:
4715 *Plant Pathology* 50(2), pp. 174–180.
- 4716 Stefansson, T. S., B. A. McDonald, and Y. Willi (2014). “The influence of
4717 genetic drift and selection on quantitative traits in a plant pathogenic
4718 fungus”. In: *PLoS ONE* 9(11).
- 4719 Stephens, M. and P. Donnelly (2003). “A comparison of Bayesian methods
4720 for haplotype reconstruction from population genotype data”. In: *Am. J.*
4721 *Hum. Genet* 73(2002), pp. 1162–1169.
- 4722 Stukenbrock, E. H. (2013). “Evolution, selection and isolation: a genomic
4723 view of speciation in fungal plant pathogens”. In: *New Phytologist* 199(4),
4724 pp. 895–907.
- 4725 Stukenbrock, E. H. (2016). “The role of hybridization in the evolution and
4726 emergence of new fungal plant pathogens.” In: *Phytopathology* 106(2),
4727 pp. 104–112.
- 4728 Stukenbrock, E. H. and T. Bataillon (2012). “A population genomics per-
4729 spective on the emergence and adaptation of new plant pathogens in
4730 agro-ecosystems”. In: *PLoS Pathogens* 8(9), pp. 1–4.
- 4731 Stukenbrock, E. H., F. B. Christiansen, T. T. Hansen, J. Y. Duthiel, and
4732 M. H. Schierup (2012). “Fusion of two divergent fungal individuals led to
4733 the recent emergence of a unique widespread pathogen species.” In:
4734 *Proceedings of the National Academy of Sciences of the United States*
4735 *of America* 109(27), pp. 1–6.
- 4736 Stumpf, M. P. H. and G. A. T. McVean (2003). “Estimating recombination
4737 rates from population-genetic data.” In: *Nature Reviews Genetics* 4(12),
4738 pp. 959–68.

- 4739 Suchard, M. A., R. E. Weiss, K. S. Dorman, and J. S. Sinsheimer (2002).
4740 “Oh brother, where art thou? A Bayes factor test for recombination with
4741 uncertain heritage.” In: *Systematic Biology* 51(5), pp. 715–728.
- 4742 Sung, P. and H. Klein (2006). “Mechanism of homologous recombination:
4743 mediators and helicases take on regulatory functions”. In: *Nat Rev Mol*
4744 *Cell Biol* 7(10), pp. 739–750.
- 4745 Svensson, O., A. Smith, J. García-alonso, and C. van Oosterhout (2016).
4746 “Hybridization generates a hopeful monster : a hermaphroditic selfing
4747 cichlid”. In: *Royal Society Open Science* 3, p. 150684.
- 4748 Swanson-Wagner, R. A., Y. Jia, R. DeCook, L. a. Borsuk, D. Nettleton, and
4749 P. S. Schnable (2006). “All possible modes of gene action are observed
4750 in a global comparison of gene expression in a maize F1 hybrid and its
4751 inbred parents.” In: *Proceedings of the National Academy of Sciences*
4752 *of the United States of America* 103(18), pp. 6805–6810.
- 4753 Taddei, F and M Radman (1997). “Role of mutator alleles in adaptive
4754 evolution”. In: *Nature* 387, pp. 700–702.
- 4755 Tanaka, M. M., C. T. Bergstrom, and B. R. Levin (2003). “The evolution
4756 of mutator genes in bacterial populations: the roles of environmental
4757 change and timing.” In: *Genetics* 164(3), pp. 843–54.
- 4758 Tavaré, S (1986). “Some probabilistic and statistical problems in the analysis
4759 of DNA sequences”. In: *American Mathematical Society: Lectures on*
4760 *Mathematics in the Life Sciences* 17, pp. 57–86.
- 4761 Tenailon, O., B. Toupance, L. Nagard, and B. Godelle (1999). “Mutators,
4762 Population Size, Adaptive Landscape and the Adaptation of Asexual
4763 Populations of Bacteria”. In: *Genetics* 152, pp. 485–493.
- 4764 Thines, M., Y. J. Choi, E. Kemen, S. Ploch, E. Holub, H.-D. Shin, and J.
4765 Jones (2009). “A new species of *Albugo* parasitic to *Arabidopsis thaliana*
4766 reveals new evolutionary patterns in white blister rusts (*Albuginaceae*)”.

- 4767 In: *Persoonia - Molecular Phylogeny and Evolution of Fungi* 22(1),
4768 pp. 123–128.
- 4769 Thines, M. (2014). “Phylogeny and evolution of plant pathogenic oomycetes
4770 - a global overview”. In: *European Journal of Plant Pathology* 138(3),
4771 pp. 431–447.
- 4772 Thomas, D. N. and G. S. Dieckmann (2002). “Antarctic Sea ice - a habitat
4773 for extremophiles.” In: *Science* 295(5555), pp. 641–644.
- 4774 Thomson, G. (1977). “The effect of a selected locus on linked neutral loci”.
4775 In: *Genetics* 85(4), pp. 753–788.
- 4776 Thorne, J. L., H Kishino, and I. S. Painter (1998). “Estimating the rate of
4777 evolution of the rate of molecular evolution.” In: *Molecular Biology and*
4778 *Evolution* 15(12), pp. 1647–57.
- 4779 True, J. R., B. S. Weir, and C. C. Laurie (1996). “A genome-wide survey of
4780 hybrid incompatibility factors by the introgression of marked segments
4781 of *Drosophila mauritiana* chromosomes into *Drosophila simulans*”. In:
4782 *Genetics* 142(3), pp. 819–837.
- 4783 Turner, T. L., M. W. Hahn, and S. V. Nuzhdin (2005). “Genomic islands of
4784 speciation in *Anopheles gambiae*”. In: *PLoS Biology* 3(9), pp. 1572–
4785 1578.
- 4786 Ungerer, M. C., S. J. Baird, J Pan, and L. H. Rieseberg (1998). “Rapid hybrid
4787 speciation in wild sunflowers.” In: *Proceedings of the National Academy*
4788 *of Sciences of the United States of America* 95(20), pp. 11757–11762.
- 4789 Ungerer, M. C., S. C. Strakosh, and Y. Zhen (2006). “Genome expansion
4790 in three hybrid sunflower species is associated with retrotransposon
4791 proliferation”. In: *Current Biology* 16(20), R872–R873.
- 4792 Vancoppenolle, M., K. M. Meiners, C. Michel, L. Bopp, F. Brabant, G. Carnat,
4793 B. Delille, D. Lannuzel, G. Madec, S. Moreau, J. L. Tison, and P. van
4794 der Merwe (2013). “Role of sea ice in global biogeochemical cycles:

- 4795 Emerging views and challenges". In: *Quaternary Science Reviews* 79,
4796 pp. 207–230.
- 4797 Vilenchik, M. M. and A. G. Knudson (2003). "Endogenous DNA double-
4798 strand breaks: production, fidelity of repair, and induction of cancer." In:
4799 *Proceedings of the National Academy of Sciences* 100(22), pp. 12871–
4800 12876.
- 4801 Visser, J. A.G. M. de (2002). "The fate of microbial mutators." In: *Microbiol-*
4802 *ogy* 148, pp. 1247–52.
- 4803 Visser, J. A.G. M. de and S. F. Elena (2007). "The evolution of sex: empirical
4804 insights into the roles of epistasis and drift." In: *Nature Reviews Genetics*
4805 8(2), pp. 139–149.
- 4806 Voordeckers, K., C. A. Brown, K. Vanneste, E. van der Zande, A. Voet,
4807 S. Maere, and K. J. Verstrepen (2012). "Reconstruction of ancestral
4808 metabolic enzymes reveals molecular mechanisms underlying evolu-
4809 tionary innovation through gene duplication". In: *PLoS Biology* 10(12),
4810 e1001446.
- 4811 Wahlund, S. (1928). "Zusammensetzung von populationen und korrelation-
4812 serscheinungen vom standpunkt der vererbungslehre aus betrachtet".
4813 In: *Hereditas* 11, pp. 65–106.
- 4814 Wakeley, J (2009). *Coalescent Theory: An Introduction*. Greenwood Village,
4815 Colorado: Roberts and Company, p. 328.
- 4816 Walker, J and M. J. Priest (2007). "A new species of *Albugo* on *Pterostylis*
4817 (*Orchidaceae*) from Australia: confirmation of the genus *Albugo* on a
4818 monocotyledonous host." In: *Australian Plant Pathology* 36, pp. 181–
4819 185.
- 4820 Wall, J. D. (2000). "A comparison of estimators of the population recombina-
4821 tion rate". In: *Molecular Biology and Evolution* 17(1), pp. 156–163.

- 4822 Wang, S., D. Bailey, K. Lindsay, J. K. Moore, and M. Holland (2014). "Impact
4823 of sea ice on the marine iron cycle and phytoplankton productivity". In:
4824 *Biogeosciences* 11(17), pp. 4713–4731.
- 4825 Waples (2002). "Definition and estimation of the effective population size
4826 in the conservation of endangered species". In: *Population Viability*
4827 *Analysis*. Ed. by S. R. Beissinger and D. R. McCullough. Chicago:
4828 Chicago University Press. Chap. 8, pp. 147–168.
- 4829 Ward, B. and C. van Oosterhout (2016). "Hybridcheck: Software for the rapid
4830 detection, visualization and dating of recombinant regions in genome
4831 sequence data". In: *Molecular Ecology Resources* 16(2), pp. 534–539.
- 4832 Watson, J. D. and F. H. C. Crick (1953). "Molecular structure of nucleic
4833 acids". In: *Nature* 171(4356), pp. 737–738.
- 4834 Watterson, G. A. (1975). "On the number of segregating sites in genetical
4835 models without recombination." In: *Theoretical population biology* 7,
4836 pp. 256–276.
- 4837 Weisse, T., U. Scheffel, P. Stadler, and W. Foissner (2007). "Local adapta-
4838 tion among geographically distant clones of the cosmopolitan freshwater
4839 ciliate *Meseres corlissi*. II. Response to pH". In: *Aquatic Microbial Ecol-*
4840 *ogy* 47(3), pp. 289–297.
- 4841 Welch, M. D. and M Meselson (2000). "Evidence for the evolution of bdelloid
4842 rotifers without sexual reproduction or genetic exchange." In: *Science*
4843 288(5469), pp. 1211–1215.
- 4844 Werner, D (1977). *The Biology of Diatoms*. Ed. by D. Werner. Oakland,
4845 California: University of California Press, p. 498.
- 4846 Whibley, A. C., N. B. Langlade, C. Andalo, A. I. Hanna, A. Bangham, C.
4847 Thébaud, and E. Coen (2006). "Evolutionary paths underlying flower
4848 color variation in *Antirrhinum*". In: *Science* 313(5789), pp. 963–966.
- 4849 Whisson, S., R. Vetukuri, A. Avrova, and C. Dixelius (2012). "Can silencing
4850 of transposons contribute to variation in effector gene expression in

- 4851 *Phytophthora infestans?*” In: *Mobile Genetic Elements* 2(2), pp. 110–
4852 114.
- 4853 White, C., K. A. Selkoe, J. Watson, D. A. Siegel, D. C. Zacherl, and R. J. Too-
4854 nen (2010). “Ocean currents help explain population genetic structure”.
4855 In: *Proceedings The Royal Society* 277(1688), pp. 1685–94.
- 4856 White, M. A. C., R. E. Blackith, R. M. Blackith, and J. Cheney (1966).
4857 “Cytogenetics of the viatica group morabine grasshoppers. I. the coastal
4858 species”. In: *Australian Journal of Zoology* 15(2), pp. 263–302.
- 4859 Whitlock, M. C. and N. H. Barton (1997). “The effective size of subdivided
4860 population”. In: *Genetics* 146(1), pp. 427–441.
- 4861 Wicker, T., F. Sabot, A. Hua-Van, J. L. Bennetzen, P. Capy, B. Chalhoub, A.
4862 Flavell, P. Leroy, M. Morgante, O. Panaud, E. Paux, P. SanMiguel, and
4863 A. H. Schulman (2007). “A unified classification system for eukaryotic
4864 transposable elements”. In: *Nature Reviews Genetics* 8(12), pp. 973–
4865 982.
- 4866 Win, J., K. V. Krasileva, S. Kamoun, K. Shirasu, B. J. Staskawicz, and M. J.
4867 Banfield (2012). “Sequence divergent RXLR effectors share a structural
4868 fold conserved across plant pathogenic oomycete species”. In: *PLoS*
4869 *Pathogens* 8(1).
- 4870 Worobey, M. (2001). “A novel approach to detecting and measuring recom-
4871 bination: New insights into evolution in viruses, bacteria, and mitochon-
4872 dria”. In: *Molecular Biology and Evolution* 18(8), pp. 1425–1434.
- 4873 Wright, S. (1932). “The roles of mutations, inbreeding, crossbreeding
4874 and selection in evolution”. In: *Proceedings of the 11th International*
4875 *Congress of Genetics* 1, pp. 356–366.
- 4876 Wright, S. (1931). “Evolution in Mendelian populations”. In: *Genetics* 16,
4877 pp. 97–159.
- 4878 — (1937). “The distribution of gene frequencies in populations.” In: *Science*
4879 23, pp. 307–320.

- 4880 — (1940). “Breeding structure of populations in relation to speciation”. In:
4881 *The American Naturalist* 74(752), pp. 232–248.
- 4882 — (1943). “Isolation by distance”. In: *Genetics* 28(2), pp. 114–138.
- 4883 Yampolsky, L. Y., T. M. M. Schaer, and D. Ebert (2014). “Adaptive phenotypic
4884 plasticity and local adaptation for temperature tolerance in freshwater
4885 zooplankton.” In: *Proceedings. Biological sciences / The Royal Society*
4886 281(1776), p. 20132744.
- 4887 Yatabe, Y., N. C. Kane, C. Scotti-Saintagne, and L. H. Rieseberg (2007).
4888 “Rampant gene exchange across a strong reproductive barrier between
4889 the annual sunflowers, *Helianthus annuus* and *H. petiolaris*.” In: *Genet-*
4890 *ics* 175(4), pp. 1883–93.
- 4891 Zhan, J., R. E. Pettway, and B. A. McDonald (2003). “The global genetic
4892 structure of the wheat pathogen *Mycosphaerella graminicola* is char-
4893 acterized by high nuclear diversity, low mitochondrial diversity, regular
4894 recombination, and gene flow”. In: *Fungal Genetics and Biology* 38(3),
4895 pp. 286–297.

Appendices

4897 **APPENDIX A**

4898 **FALCON assembly haplotype divergence**

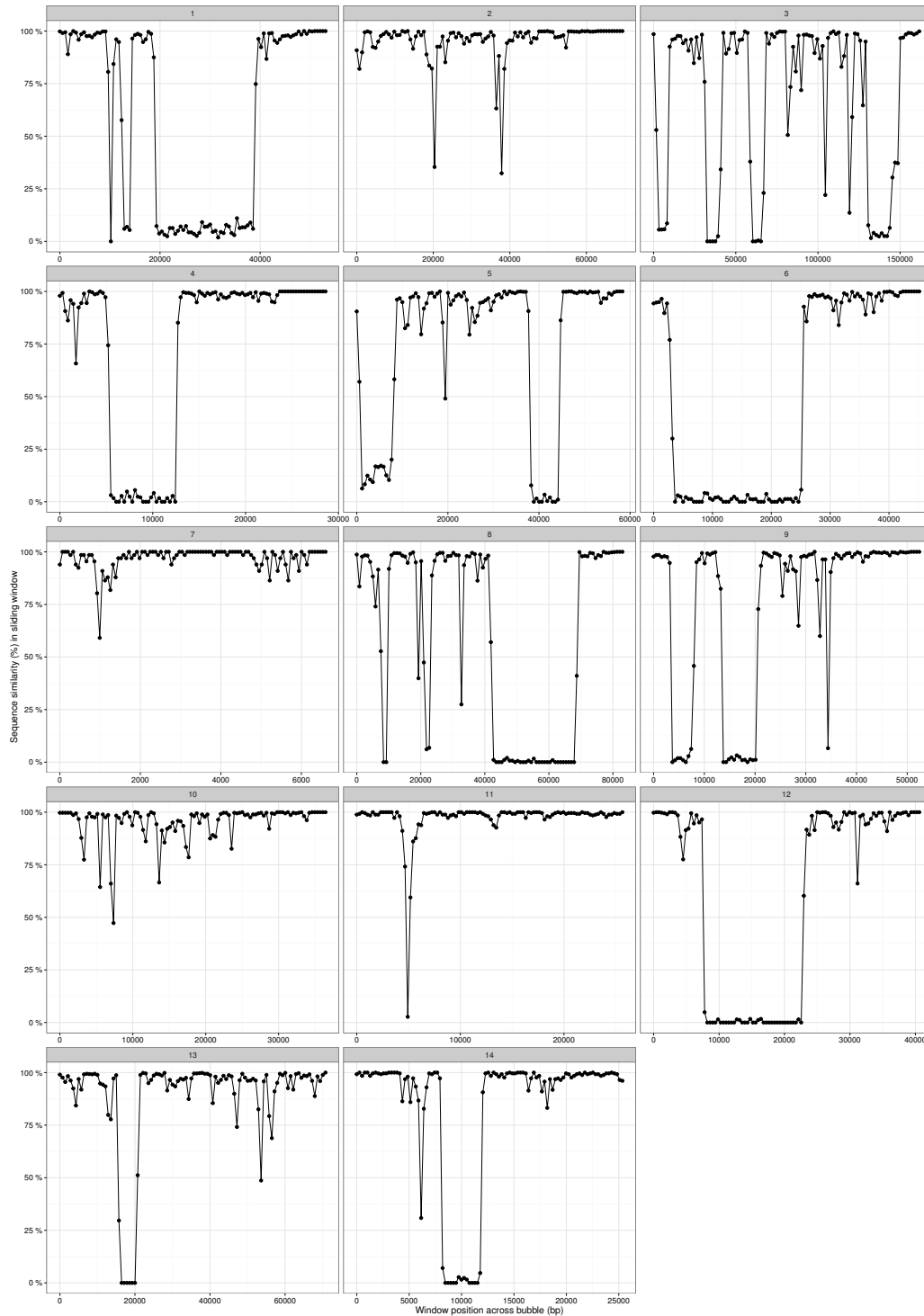


Figure A.1: Sequence similarity calculated with sliding windows across each haplotype 'bubble' in chromosome 000002F, from the *F. cylindrus* FALCON genome assembly. Regions of divergence and indels are apparent.