

Living in each other's pockets: Nucleotide variation inside a genomic island harboring *Pan I* and its neighbors in Atlantic cod

Ubaldo Benitez Hernandez, Einar Árnason

The *Pan I* locus in Atlantic cod lies in a genomic island of divergence extending over a large genomic region. The locus has two divergent alleles, defined by a single *DraI* restriction site, that have been related to behavioral differences of habitat selection by depth and temperature. The *Pan I* locus is known to be under an unusual mix of balancing selection and selective sweeps within the functional types. Here we study nucleotide variation in a 12.5 kb region inside the genomic island harboring *Pan I* and neighboring loci for sortilin 1 (*Sort1*) and ataxin 7-like 2 (*Atxn7/2*) which we partially covered. Variation of the 31 gene copies throughout the region falls into two divergent haplogroups that correlate with the 25 copies of *A* and six copies of *B* alleles of *Pan I*. The unfolded site frequency spectrum for the part with Pacific cod used as the outgroup is trimodal with a mode at singletons and two high frequency modes at 6/31 and 25/31 representing the two genealogical lineages. The folded site frequency spectrum for the entire region similarly has a high frequency mode of mutations that have accumulated on the two lineages. The high frequency of singletons is accounted for by multiple merger coalescent models. Parameter estimates using these models indicate sweepstakes reproduction. The high frequency modes of the spectrum is evidence for balancing selection. Analysis of non-synonymous changes shows that *Pan I* is at least one focus of selection within the genomic island. There may be multiple sites of selection and epistatic interactions. There is extensive linkage disequilibrium throughout the region. We suggest that the genomic island of divergence is a supergene of co-adapted complexes possibly locked together by structural variation.

Living in each other's pockets: Nucleotide variation inside a genomic island harboring *Pan I* and its neighbors in Atlantic cod

Ubaldo Benitez Hernandez¹ and Einar Árnason²

¹Institute of Life and Environmental Sciences, University of Iceland, Reykjavik, Iceland

²Institute of Life and Environmental Sciences, University of Iceland, Reykjavik, Iceland

ABSTRACT

The *Pan I* locus in Atlantic cod lies in a genomic island of divergence extending over a large genomic region. The locus has two divergent alleles, defined by a single *Dral* restriction site, that have been related to behavioral differences of habitat selection by depth and temperature. The *Pan I* locus is known to be under an unusual mix of balancing selection and selective sweeps within the functional types. Here we study nucleotide variation in a 12.5 kb region inside the genomic island harboring *Pan I* and neighboring loci for sortilin 1 (*Sort1*) and ataxin 7-like 2 (*Atxn7l2*) which we partially covered. Variation of the 31 gene copies throughout the region falls into two divergent haplogroups that correlate with the 25 copies of *A* and six copies of *B* alleles of *Pan I*. The unfolded site frequency spectrum for the part with Pacific cod used as the outgroup is trimodal with a mode at singletons and two high frequency modes at 6/31 and 25/31 representing the two genealogical lineages. The folded site frequency spectrum for the entire region similarly has a high frequency mode of mutations that have accumulated on the two lineages. The high frequency of singletons is accounted for by multiple merger coalescent models. Parameter estimates using these models indicate sweepstakes reproduction. The high frequency modes of the spectrum is evidence for balancing selection. Analysis of non-synonymous changes shows that *Pan I* is at least one focus of selection within the genomic island. There may be multiple sites of selection and epistatic interactions. There is extensive linkage disequilibrium throughout the region. We suggest that the genomic island of divergence is a supergene of co-adapted complexes possibly locked together by structural variation.

Keywords: Balancing selection, Linkage Disequilibrium, Pantophysin, Pan I, Genomic Island, Sortilin, Ataxin-7 like, Atlantic cod

INTRODUCTION

Selection, the differential fecundity and mortality of genotypes, is a most powerful evolutionary force. Organisms exploit finite resources with differential efficiency leading to fitness differences. They thereby pass on their alleles to future generations with differential efficiency and thus are selected for (Lewontin, 1974). The selective forces arise from diverse physical or biotic factors and can exist in different combinations, resulting also in diverse patterns of polymorphism. Various modes of selection exist. Under negative or purifying selection, detrimental mutations are purged. Under positive or advantageous selection variants sweep to fixation. Under balancing selection, however, alternative allelic forms exist at intermediate frequencies due to a tug of selective forces that may ensue because of spatially or temporally varying selective pressures favoring allelic forms differently, due to heterozygous advantage or overdominance, or due to inverse frequency-dependent selection.

The genealogical relationships among alleles are captured by the coalescent (Kingman, 1982) that is a retrospective model of the assignment of alleles of a sample to ancestors. Going back in time, alleles reach a point of common ancestry. Different present-day lineages of descent, in hind-sight fuse or coalesce in an ancestral form they originate from. Eventually all alleles of a sample merge in a single most recent common ancestor (Wakeley, 2009). The coalescent yields the theoretical site frequency spectrum of a population at neutrality that can be contrasted with the observed site frequency spectrum thus allowing for the detection of deviations from the theoretical expectations at neutrality. The classic coalescent model (Kingman, 1982), derived from Fisher/Wright model of reproduction among organisms

29 with non-skewed low fecundity offspring distributions, considers bifurcating coalescent events or binary
30 mergers. The more generalized models such as the Beta ($2 - \alpha, \alpha$) (Schweinsberg, 2003) and point-mass
31 coalescent (Eldon and Wakeley, 2006) that model multiple merger coalescent events, are a major research
32 focus in new developments of coalescent theory (Wakeley, 2013). These models are appropriate for
33 organisms exhibiting both high fecundity and highly-skewed heavy-tailed offspring distributions, and
34 may be a better null model for many species (Eldon and Wakeley, 2006), such as the Atlantic cod (*Gadus*
35 *morhua*, Linnaeus, 1758) a marine organism with high-fecundity and opportunity for very high variance
36 of offspring numbers due to its life history traits (Árnason, 2004).

37 The *Pan I* locus in Atlantic cod, which has been used widely as a marker for population genetics
38 analysis, shows strong differentiation among populations at different geographic scales (Fevolden and
39 Pogson, 1997; Pogson et al., 2001; Pogson and Fevolden, 2003; Karlsson and Mork, 2003; Árnason et al.,
40 2009). The locus is acknowledged to be under selection or to be linked to selected loci (Fevolden and
41 Pogson, 1997; Case et al., 2005). The continued use of this selected locus with the aim of identifying
42 local populations at the very least begs advancing the knowledge on the evolutionary forces at work at the
43 locus. Thus Nielsen et al. (2007) recommend that it not be used at microgeographical scale until more
44 knowledge on evolutionary drivers is attained.

45 Pogson (2001) characterized a 1.85 kilo base (kb) region of the *Pan I* locus in Atlantic cod, and showed
46 that it harbors an ancient polymorphic *DraI* restriction site that through absence or presence defines
47 alternative alleles or haplogroups, *A* and *B* respectively. The two alleles are maintained by balancing
48 selection and are highly divergent at the nucleotide and amino acid level. Both alleles also show signs of
49 selective sweeps within the functional types (Pogson, 2001). The two types differ by four amino acids
50 representing six amino acid replacement mutations fixed between the lineages (three in each lineage). The
51 amino acids differences reside in the first intra-vesicular loop domain (IV1) of the protein. The *DraI* site
52 defining the alleles and various other restriction sites in the region show high linkage disequilibrium (LD).
53 There is LD with sites defining a 5.7 kb restriction fragment around the *Pan I* locus (Pogson, 2001).

54 There is a strong correlation between *Pan I* allele frequency and environmental settings (coastal vs
55 offshore) of different depth. The *A* allele is found in higher proportions at coastal/shallow-water locales,
56 and the *B* alleles at offshore/deep-water locales (Pogson and Fevolden, 2003; Case et al., 2005; Árnason
57 et al., 2009). A very steep gradient is found of allele frequency, a change of 0.4% per meter of depth
58 down to about 200 m (Árnason et al., 2009). This bears on the association found between the *Pan I* locus
59 and Atlantic cod behavioral ecotypes defined using data storage tags (Pálsson and Thorsteinsson, 2003;
60 Pampoulie et al., 2008). The ecotypes exhibit either a shallow-water behavior characterized by seasonal
61 temperature trends (stationary cod), or deep-water behavior characterized by frequent vertical migrations
62 and steep temperature changes possibly representing foraging at thermal fronts (migratory cod) (Pálsson
63 and Thorsteinsson, 2003; Pampoulie et al., 2008; Thorsteinsson et al., 2012). *Pan I* genotypes of Atlantic
64 cod in relation to depth show that *AA* individuals have a shallow water behavior and *BB* a deep water
65 behavior, however, seeking shallower waters during spawning. *AB* individuals show a mixed behavior
66 (Pampoulie et al., 2008) somewhat intermediate between the homozygotes. The different cod ecotypes
67 and their associated *Pan I* genotypes share their depth range during spawning (Pampoulie et al., 2008;
68 Árnason et al., 2009). However, possible segregation may occur by the behavioral differences of the
69 ecotypes in spawning habitat-selection (Grabowski et al., 2011). Thus the markedly divergent *A* and *B*
70 alleles can be roughly classified as shallow and deep-water adapted types respectively. This presents a
71 strong parallel between a genomic island and ecological divergence associated to *Pan I*.

72 Strong heterogeneity in levels of differentiation in different parts of the genome have revealed genomic
73 islands of divergence (Wu, 2001; Renaut et al., 2013; Ruegg et al., 2014; Cruickshank and Hahn, 2014) in
74 Atlantic cod. The *Pan I* locus is located within one of such genomic islands of divergence (Bradbury et al.,
75 2013; Hemmer-Hansen et al., 2013; Karlsten et al., 2013). The genomic islands represent a non-random
76 distribution of levels of divergence in the Atlantic cod genome, and are constituted by clusters of genomic
77 regions of elevated divergence running within several linkage groups and with co-occurrence of loci
78 most likely implicated in selective processes (Bradbury et al., 2013). The genomic island containing
79 *Pan I* is in linkage group (LG) 1 and has been found to be linked to the aforementioned *Pan I* ecotypes
80 (Hemmer-Hansen et al., 2013). Thus *Pan I* and loci co-occurring at the same genomic island are likely to
81 be functionally related to the capability of the organism to thrive in different environments with complex
82 differences based upon a discrete region of genomic divergence composed by multiple linked loci.

83 Different environments entail multidimensional differences that must be met by organisms inhabiting

84 those environments. When alternative forms of an organism inhabit different environments, divergent
85 selection may be involved in building supergenes or switch-genes (see e.g. Thompson and Jiggins, 2014).
86 Supergenes are genomic architectures of multiple, functional, co-adapted loci in tight linkage and little
87 recombination by a variety of mechanisms. The polymorphic variants segregate together in particular
88 combinations of alleles as if they were a single locus. Often those variants are kept at intermediate
89 frequencies due to a balance of selective vectors (Thompson and Jiggins, 2014) as in butterfly mimicry.
90 The latter particular combinations of alleles at co-adapted loci, which are reflected on complex phenotypes,
91 allow alternative forms of an organism to meet the multidimensional challenges of particular habitats, with
92 no maladaptive intermediate combinations thanks to little recombination and tight linkage among loci
93 (Thompson and Jiggins, 2014). *Pan I* has shown highly divergent alleles likely maintained by balancing
94 selection and the region shows considerable LD (Pogson, 2001). The question arises how far those
95 influences extend into the neighboring loci of *Pan I*? How tight-knit is the LD within *Pan I* and its
96 surrounding loci? High divergence, tight-knit LD among multilocus variants, the implication of balancing
97 selection, and suggested functional correlation among loci would present conditions for the build up
98 of a supergene structure. If the 20 cM genomic Island of divergence at LG1 (Hemmer-Hansen et al.,
99 2013) represents a supergene of co-adapted complexes there may of course be multiple sites of epistatic
100 interactions throughout this genomic island.

101 The *Pan I* locus has been linked to selective forces such as temperature and salinity (Case et al., 2005)
102 and fisheries (Árnason et al., 2009). Although we do not have evidence of epistatic effects, the function
103 of the proteins coded by *Pan I* and flanking loci suggests correlation at functional level. Knowledge of
104 function is essential to understand the working of selection. The *Pan I* gene codes for pantophysin, a
105 microvesicle membrane protein involved in transport events (Haass et al., 1996), specifically in the traffick-
106 ing of the insulin-regulated glucose transporter GLUT4 (see reviews in Bradley et al., 2001; Larance et al.,
107 2008). The loci on either side of *Pan I* are *Sort1* and *Atxn7l2* (Star et al., 2011). *Sort1* codes for sortilin, a
108 protein that also is a major component of Glut4-containing microvesicles and that might be involved in
109 the translocation or biogenesis of the Glut4-containing vesicles (Lin et al., 1997). Sortilin participates in
110 trafficking processes at the Golgi apparatus and plasma membrane (Strong et al., 2012). *Atxn7l2* codes for
111 ataxin 7-like 2, a protein that contains SCA7, a zinc-binding domain that binds with TFTC/STAGA sub-
112 units (Marchler-Bauer et al., 2012). TFTC/STAGA are histone acetyltransferase-containing coactivator
113 complexes (Helmlinger et al., 2006) which are implicated in chromatin remodeling (Zhao et al., 2008).
114 Chromatin remodeling is a response mechanism to environmental stressors such as temperature or salinity
115 (de Nadal et al., 2011). From these considerations we speculate that the proteins of these genes may be
116 important for burst energy metabolism (anaerobic glycolysis). Thus the nearest neighbors of *Pan I* may
117 be functionally related to *Pan I*, with *Sort1* products acting in the same transport vesicles and *Atxn7l2*
118 products responding to similar environmental pressures.

119 If the *Pan I* locus is influenced by the dynamics of a selected linked locus the neighboring loci *Sort1*
120 and *Atxn7l2* are a reasonable location to consider for analysis of LD and selective effects. Yet, even if
121 strong LD and signatures of selection were to be found at those loci, the question still remains what is
122 the target of selection. A simple metric for elucidating the action of selective forces at a locus is the
123 proportion of codons with multiple non-synonymous substitutions. Due to the nature of the genetic code,
124 some amino acid substitutions are not accessible for certain codons. Instead a change from one amino
125 acid to another may have to go through intermediate states, accumulating non-synonymous substitutions
126 at the codons in the process. Environmental changes can create bursts of non-synonymous substitutions
127 or evolutionary bursts (Gillespie, 1984). DNA sequence studies have pointed to areas suspected of being
128 under such selective bursts (Bazykin et al., 2004, 2006). We apply this metric by comparing the proportion
129 of codons hit by multiple non-synonymous mutations at *Pan I* and its flanking regions, segments of the
130 (*Sort1*) and *Atxn7l2* loci.

131 In this paper we study detailed nucleotide variation of a 12.6 kb region within the genomic island
132 containing the *Sort1*, *Pan I*, and *Atxn7l2* loci. The main questions we address are: Do the signals of
133 selective effects already known at *Pan I* extend to its neighboring genes or is there a peak signal at *Pan I*?
134 How tightly knit is the LD in the region and what role does it have concerning selective effects? Can we
135 detect a focal point of selection within this region?

MATERIALS AND METHODS

Sampling

We selected 31 individual Atlantic cod genomic DNA samples for genotyping, cloning and sequencing, by stratified random sampling on geographic regions of the species distribution spanning east to west and north to south. The samples come from our large sample laboratory database of greater than 20,000 individuals. All localities are represented with at least 100 individuals (except the White Sea with 24 individuals). The geographic regions were the waters of Newfoundland (New), Iceland (Ice), Faroe Islands (Far), Norway (Nor), the North Sea (Nse) and the Celtic Sea (Cel) (6, 5, 5, 6, 6, and 3 samples, respectively). We isolated DNA from gill tissue using a Chelex/Proteinase K digestion (Walsh et al., 1991). We included a specimen of the closely related Pacific cod *G. macrocephalus* (Tilesius, 1810) as an outgroup.

The Icelandic Committee for Welfare of Experimental Animals, Chief Veterinary Office at the Ministry of Agriculture, Reykjavik, Iceland has determined that the research conducted here is not subject to the laws concerning the Welfare of Experimental Animals (The Icelandic Law on Animal Protection, Law 15/1994, last updated with Law 157/2012). DNA was isolated from tissue taken from dead fish on board research vessels. Fish were collected during the yearly surveys of the Icelandic Marine Research Institute. All research plans and sampling of fish, including the ones for the current project, have been evaluated and approved by the Marine Research Institute Board of Directors. Samples were also obtained from dead fish from marine research institutes in Norway, the Netherlands, Canada and the US that were similarly approved by the respective ethics boards. The samples from the US used in this study have been described in Cunningham et al. (2009) and the samples from Norway in Árnason and Pálsson (1996a). The samples from Canada consisted of DNA isolated from the samples described in Pogson (2001). The samples from the Netherlands were obtained from the Beam-Trawl-Survey (<http://www.wageningenur.nl/en/Expertise-Services/Research-Institutes/imares/Weblogs/Beam-Trawl-Survey.htm>) of the Institute for Marine Resources & Ecosystem Studies (IMARES), Wageningen University, the Netherlands, which is approved by the IMARES Animal Care Committee and IMARES Board of Directors.

Molecular analysis

We genotyped the individuals at the *Pan* I locus in the manner described by Árnason et al. (2009). The locus has two alleles or haplogroups, *A* and *B*, corresponding to the absence or presence of a *Dra*I restriction site (Fevolden and Pogson, 1997). This site is in LD with several other sites of the *Pan* I locus (Pogson, 2001). If the individual was heterozygous *A/B* for the *Pan* I we chose one allele at random from that individual for this study.

We amplified two fragments 4.3 kb and 8.7 kb in length (fragment I and fragment II respectively) from genomic DNA using Long PCR Enzyme Mix (Fermentas) in a Tetrad2 (MJ Research). The fragments had a 489 base pair (bp) overlap that contains the polymorphic *Dra*I restriction site defining the *A* and *B* alleles of *Pan* I. The merged sequence of both fragments resulted in 12.5 kb sequence (Figure S1). We had previously 454-sequenced BAC clones with about a 150 kb insert which contained the *Pan* I gene as well as a number of other genes (and see Árnason and Halldórsdóttir, 2015). With this sequence at hand, we designed primers to PCR amplify and sequence the 12.5 kb fragment containing the full *Pan* I gene. The neighboring genes were sortilin 1 (*Sort1*) and ataxin 7-like 2 (*Atxn7l2*) which we partially covered. Subsequently, the cod genome was released (Star et al., 2011) and confirmed our BAC clone sequence. We have used the genomic sequence (www.ensembl.org) and comparative data from various species ranging from humans to fish to determine the exon/intron structure of our 12.5 kb fragment.

We used primer 3 (Pogson and Mesa, 2004) and primer sc343pr66398 to amplify fragment I, and primer 20 (Pogson and Mesa, 2004) and primer sc343pr79421 to amplify fragment II (Figure S1 and Table S1). The PCR was an initial denaturation step of 2 minutes at 94°C; followed by 10 cycles of 15 seconds denaturation at 94°C, 30 seconds annealing at 50.4°C, 9 minutes 18 seconds (for fragment I) and 4 minutes 36 seconds (for fragment II) elongation at 68°C. This was followed by 25 additional cycles, increasing the elongation time by 10 seconds every cycle.

We purified fragments for cloning by agarose gel purification. We loaded 40 μ l of the PCR amplified products in a 0.8% agarose gel (1 \times TAE buffer) with crystal violet (1.6 μ g/ml) and electrophoresed at 80 volts/cm for 56 minutes using a 1 kb DNA ladder (Fermentas) as reference. We excised the gel

190 pieces under visible light with the DNA bands of interest, froze and thawed, and used the resulting DNA
191 suspension directly for cloning.

192 We TOPO-TA cloned fragments into vector pCR-XL-TOPO (Invitrogen) following the manufacturer's
193 instructions except we used 1/7 of recommended amount of vector. After cloning, we isolated plasmid
194 DNA using alkaline lysis minipreps (Birnboim and Doly, 1979). We miniprepmed five clones per
195 individual. We confirmed that the clones contained fragments of the size of interest, by EcoRI digestion
196 and agarose gel electrophoresis.

197 We genotyped the clones for the *Pan I* A and B alleles. For each individual, we took three clones from
198 the same allele for sequencing, thus sequencing 93 clones \times 2 (three clones for both fragments per each
199 of the 31 individuals). We sequenced with overlapping primers (Figure S1 and Table S1) using BigDye
200 Terminator kit (ABI) and manufacturer's protocol except we used 1/16 of recommended amount of TRR.
201 We purified reaction products with EtOH precipitation, resuspended in HiDi Formamide and ran on an
202 ABI 3500xL genetic analyzer (ABI). The area we sequenced comprised 3 loci: the entire *Pan I* locus, and
203 partial segments of *Atxn7l2* and *Sort1*. *Pan I* is located medially; *Atxn7l2* is located downstream and on
204 the same strand as *Pan I*; and *Sort1* is located upstream, and on the opposite strand of *Pan I* (Figure 1).

205 For comparative purposes, we sequenced three clones of a Pacific cod (*Gadus macrocephalus*)
206 individual, covering the same region as fragment I (4.3kb) in Atlantic cod (partial segments of *Pan I* and
207 *Atxn7l*). We used the same methods as with Atlantic cod. We could not get fragment II to amplify for
208 Pacific cod.

209 Data analysis

210 We base-called, assembled and visually inspected 1836 DNA reads into 93 sequences 12.5 kb in length
211 (a 12.5 kb sequence per clone) using the software suite Phred-Phrap-Consed (Ewing et al., 1998;
212 Gordon et al., 1998; Green, 1994). For each of the 31 individuals, we aligned its three clone sequences
213 using MUSCLE (Edgar, 2004) to edit and build a consensus sequence for each individual. We used the R
214 language and environment (R Development Core Team, 2008) with the APE (Paradis et al., 2004), Pegas
215 (Paradis, 2010b) and SeqinR (Charif and Lobry, 2007) packages and in-house functions to manipulate
216 the clone sequence alignments and build 31 individual consensus sequences. These are phased haplotypic
217 data. We applied the same procedures on 18 DNA reads from three clones of a Pacific cod individual to
218 obtain a consensus sequence 4.3 kb in length.

219 PCR errors inevitably are found in clones by cloning PCR products. We consider that taking three
220 clones is sufficient to eliminate PCR errors among the clones. We assume that the three clones do not
221 share a PCR error site (Árnason and Halldórsdóttir, 2015). Two of the clones from each individual will
222 be of the same chromosome. The third clone will be of the same chromosome with probability 1/2 and
223 of the alternative chromosome with probability 1/2 (Árnason and Halldórsdóttir, 2015). The consensus
224 sequence will eliminate PCR errors except in the rare cases in which PCR errors in one of the two clones
225 from the same chromosome has hit a site which is polymorphic in the population and found in the third
226 clone derived from the alternative chromosome. A small bias may be introduced by this. However, this
227 will be seen as recombination and we would in such cases err on the conservative side in interpretation
228 based on LD.

229 Since we sequenced two fragments for each individual and merged them (Figure S1), we had to
230 address the possibility of inadvertently forming chimeras where the consensus sequence of one fragment
231 corresponds to the alternative allele relative to the other fragment. To investigate if this was an issue
232 we aligned and contrasted the 489 bp overlap in both fragments and checked that the polymorphisms in
233 both fragments produced the same consensus sequence. Doing this we detected five possible chimeric
234 individual sequences in our first overview of the data. For each of them, we replaced the 4.3 kb consensus
235 sequence from fragment I with the sequence of the single clone in fragment I that was in phase with
236 fragment II. This means that for these five individual sequences fragment I would have PCR errors in
237 these cases that would appear as excess of singleton variable sites. We randomly eliminated the excess of
238 singletons (corresponding to PCR errors) in these five clone sequences to obtain the same average number
239 of singletons as in the 4.3 kb region of non-chimeric individuals. Therefore, the remaining singletons
240 were scattered at random, in agreement with the nature of mutations.

241 We aligned the 31 individual consensus sequences with the Fast Statistical Alignment (FSA) program
242 (Bradley et al., 2009), and visually inspected this alignment with Seaview (Galtier et al., 1996). We
243 manually edited a few indel sites where FSA had made obvious errors. We used the alignment with

244 SNIplay (Dereeper et al., 2011) for SNP detection. We used the R `ade4`, `adegenet` `LDheatmap`,
 245 and `popgen` packages (Jombart and Ahmed, 2011; Shin et al., 2006; Marchini, 2013; Dray and Dufour,
 246 2007) and various functions written by us for managing, analyzing, and plotting the data. We used
 247 `Genomicus/Phyloview` (Louis et al., 2013) to produce a multi-species comparative display in
 248 genomic context of our sequenced region and surrounding loci, showing ortholog and paralog genes
 249 (Figure S2). We used the `snpposi` functions of the `adegenet` package (Jombart and Ahmed, 2011) to
 250 plot and test the density of SNPs over the fragment. We used `Arlequin v. 3.5` (Excoffier and Lischer, 2010)
 251 for analysis of molecular variance, AMOVA, and analysis of population differentiation with pairwise F_{ST} .

252 We generated folded and unfolded site frequency spectra using R with the package `pegas` Paradis
 253 (2010a) and with in-house functions. We compared observed spectra to expectation of Kingman coalescent
 254 (θ/i) (Kingman, 1982) and to multiple merger Beta ($2 - \alpha, \alpha$) (Schweinsberg, 2003) and point-mass coa-
 255 lescent (Eldon and Wakeley, 2006) using inference methods developed by Birkner et al. (2013). We used
 256 software from Bjarki Eldon (<http://page.math.tu-berlin.de/~eldon/programs.html>)
 257 to estimate various parameters of the multiple merger coalescents (and see Árnason and Halldórsdóttir,
 258 2015). We carried out Tajima's D , F_u and L_i 's, and McDonald-Kreitman neutrality tests, analysis of
 259 recombination, and computed statistics of polymorphism and divergence with `DNAsp` (Rozas et al., 2003).

260 We performed Hudson-Kreitman-Aguadé (HKA) test (Hudson et al., 1987) in direct mode with
 261 `DNAsp` and also with a sliding window using `DNA_Slider` (McDonald, 1998). For the HKA tests, we
 262 trimmed the ends of fragment I sequence alignment (resulting in 4.2 kb) of Atlantic and Pacific cod due to
 263 low phred scores of the Pacific cod sequences at the ends of the fragment. For the HKA test in direct mode
 264 we contrasted DNA sequences (31 and one, respectively) of Atlantic and Pacific cod at the mitochondrial
 265 cytochrome *b* locus (data from Árnason, 2004) with the 4.2 kb region (fragment I) covering the *Pan I* and
 266 *Atxn712* loci partially. For the HKA test with sliding window we used the 4.2 kb alignment (fragment
 267 I) of Atlantic and Pacific cod sequences (31 and one, respectively) with silent polymorphisms and fixed
 268 differences, windows of 31 and 33 variable sites for the largest average and maximum sliding G value,
 269 and 100 replications. Also, for the HKA test we used 12.5 kb region to contrast all *Pan I^A* against *Pan I^B*
 270 as outgroup, and vice versa, using all polymorphisms and fixed differences, and 100 replications. We
 271 also performed the maximum likelihood HKA test using the `MLHKA` program (Wright and Charlesworth,
 272 2004). For this test we used the data on Hemoglobin $\alpha 2$ *Hba2* and Myoglobin *Myg* loci of Árnason and
 273 Halldórsdóttir (2015) with the 4.2 kb alignment (fragment I) of Atlantic and Pacific cod sequences for
 274 *Pan I*.

275 GenBank accession numbers for sequences reported in this paper are KR011783–KR011814.

276 RESULTS

277 Nucleotide variability

278 Extensive LD exists between various sites of the *A* and *B* alleles as already observed by Pogson (2001).
 279 We, therefore, decided to analyze and present our data with reference to the *DraI* site defining the *A* (25
 280 sequences) and *B* (six sequences) alleles of *Pan I*. Examination of the data showed that the *DraI* site
 281 defining the *A* and *B* *Pan I* difference was tied to differences at a larger scale.

282 We found maximum haplotype diversity with every sequence representing a different haplotype.
 283 Haplotype diversity was thus not informative about differences at this level of sampling. The levels of
 284 polymorphism were notably higher for *Pan I^A* sequence variants compared to *Pan I^B* (Table 1). $\hat{\pi}$ and
 285 $\hat{\theta}$ values of *Pan I^A* sequence variants were roughly 3 times and 5 times larger, respectively, than the
 286 same statistics for *Pan I^B*. Combined the sequence variants showed much higher levels of polymorphism
 287 than each group of sequence variant separately (Table 1). The levels of divergence (\hat{K}) between the
 288 *Pan I* sequence variants were higher than the levels of nucleotide diversity ($\hat{\pi}$) within both (Figure 2),
 289 throughout the region.

290 Heterozygosity per site among sequences classified according to the *Pan I A* alleles, sequences
 291 classified according to the *B* alleles, and for all sequences combined are shown in Figure S3. High
 292 heterozygosity was found throughout the region for both the *A* and *B* alleles. Nevertheless there were
 293 concentrations of high heterozygosity sites in some parts (e.g. around 1000 and 8500 for *A* alleles
 294 and 3200 for *B* alleles). There was high heterozygosity throughout the region for the combined data.
 295 Heterozygosity of 0.31 represented the fixed differences between sequences classified by the *Pan I A* and
 296 *B* alleles (25 vs six respectively). Although high heterozygosity was found throughout the region there was
 297 significant clustering of SNPs ($P = 0.006$, `snpposi.test`, Jombart and Ahmed, 2011) (Figure S4).

298 The maximum likelihood tree of the alleles showed two distinct lineages with the variation grouped
 299 according to the two *Pan I* allelic variants (Figure S5). The *Pan I^A* lineages showed higher sequence
 300 variability than the *Pan I^B* lineage. Figure S5 also showed that with respect to the outgroup the *Pan*
 301 *I^B* lineage had evolved further than the *Pan I^A* lineage (Figure S6).

302 Tajima's $D = -0.72989$ and Fu and Li's $D^* = -0.93532$ and $F^* = -1.01658$ were non-significant
 303 for the overall region ($P > 0.10$), as well as for each of the loci separately. With a sliding windows
 304 approach (100 bp and 25 bp window and step size, respectively), we found a region between 10914 bp and
 305 11038 bp, in *Atxn7l2*, with Fu and Li's $D^* = -2.7105$ and $F^* = -2.8126$ that deviate significantly from
 306 neutrality ($P < 0.05$). The McDonald-Kreitman test did not show a significant deviation from neutrality
 307 for the overall region or for each of the loci separately ($P > 0.10$, Fisher's exact test).

308 The HKA test with sliding windows indicated a significant deviation from neutrality only at the *Pan*
 309 *I* locus when considering the Atlantic cod *Pan I^A* sequence variants as ingroup and *Pan I^B* sequence
 310 variants as outgroup, but not at the *Atxn7l2* locus or when considering *Pan I^A* and *Pan I^B* sequence
 311 variants of Atlantic cod as ingroup compared against Pacific cod as outgroup. The HKA test in direct
 312 mode comparing Atlantic and Pacific cod at segments of *Pan I* and *Atxn7l2*, and at cytochrome *b*, did not
 313 indicate a significant deviation from neutrality. However, the maximum likelihood HKA analysis showed
 314 a significant HKA test ($P < 0.01$) with a selection parameter $k = 4.12$ indicative of balancing selection
 315 (Table S2).

316 Linkage disequilibrium, LD

317 We observed very strong LD among most of the high heterozygosity polymorphic sites (those with minor
 318 allele frequency 6/31 or more) of the three analyzed loci over the 12.56 kb region (Figure 3). Virtually
 319 the whole 12.56 kb region, that harbors the *Pan I* locus surrounded by partial segments of *Sort1* and
 320 *Atxn7l2*, is one LD block with maximum LD (measured by D') throughout the whole region. Very few
 321 polymorphic sites had LD values lower than maximum. However, there were notable exceptions. Three
 322 adjacent sites (sites number 8360, 8362, and 8364) were in full linkage equilibrium. There are three
 323 possible explanations for this phenomenon. First they might be due to sequencing error. We have gone
 324 over the data and found no evidence for error. Second, these may hypermutable sites. In that case the
 325 variants at these sites are not identical by descent. Third, this may be a recombination tract with the blocks
 326 on either side of that tract being held together in full LD by epistatic interactions.

327 Measures of LD depend on allele frequencies (Hedrick, 1987) and in general no measure is independent
 328 of allele frequencies (Lewontin, 1988). Excluding only singleton sites the LD of sites with a minor allele
 329 frequency of 2/31 or more also showed large LD blocks (Figure S7). However, another recombination
 330 tract was observed having intermediate D' LD values in the intergenic region of *Sort1* and *Pan I*.

331 Considering the *A* allele sequences only the LD was much lower with evidence of extensive recombi-
 332 nation among the *A* alleles (Figure S8, minor allele frequency set at 3/25). There was much less variation
 333 among *B* alleles and there were large blocks of LD but also recombination tracks with low LD (Figure S9,
 334 sites with minor allele frequency 2/6).

335 Population differentiation

336 Considering only variation for the 12558 bp region among the *Pan I A* alleles the AMOVA (Table S3)
 337 showed that most of the variation (84%) was within populations. On the basis of spatial patterns of
 338 variation at the *Ckma* gene Árnason and Halldórsdóttir (2015) observed a North (Canada, Iceland, and
 339 Norway) vs South (Faroe Islands, North Sea, and Celtic Sea) divide. Using this classification to group
 340 localities only 6% of the variation was among groups and 11% within groups. Only the within population
 341 variance component V_c and the associated F_{ST} fixation index was significant (Table S3). The lack of
 342 significance was probably to some extent due to small sample sizes but the size of the fragment counteracts
 343 that effect.

344 The pairwise F_{ST} of *A* allele variation between localities (Table S4) showed that Canada (Nova Scotia
 345 and Newfoundland combined) differed significantly from Norway, Faroe Islands, North Sea, and Celtic
 346 Sea. The differentiation of Canada and Iceland was a little over 1/3 that of Canada and the other localities
 347 but it was not significant. The only other significant difference was between Iceland and North Sea.
 348 There was no differentiation among any of the pairs of Norway, Faroe Islands, North Sea, and Celtic Sea
 349 with most F_{ST} having negative signs that are interpreted as null. These patterns were also evident in the
 350 maximum likelihood tree of variation (Figure S10). One clade was confined to Canada but the Canadian
 351 samples were not, however, confined to that clade. In general individuals from most localities were widely

352 dispersed on branches of the tree. Overall the $F_{ST} = 0.09 \pm 0.02$ among *A* alleles for the North vs South
353 areas defined in Árnason and Halldórsdóttir (2015). Thus differentiation at *Pan I* and peripheral regions
354 could be described as an east vs west differentiation with Iceland intermediate. It did not fit the north vs
355 south divide of Árnason and Halldórsdóttir (2015).

356 The sample contained only six *B* alleles, three from Iceland and three from Norway. All *B* carried
357 the ∇_2 indel considered a sign of a selective sweep (Pogson, 2001). There were two clades among the *B*
358 alleles that were defined by several sites in full LD (Figure S11). Both clades were found in both Iceland
359 and Norway and were thus not geographically restricted. The $F_{ST} = -0.09$ ($P = 0.66 \pm 0.04$) for Iceland
360 vs Norway comparison of *B* allele variants.

361 Allelic divergence

362 The $F_{ST} = 0.82$ ($P < 0.001$) between the *Pan I* *A* and *B* haplogroup variants. The average number of
363 pairwise differences for the 12558 bp between the *A* and *B* alleles $D_{xy} = 442.7$, the average within allele
364 difference $D_X = 97.4$, and the corrected pairwise difference $D_a = D_{xy} - (D_x + D_y)/2 = 384.2$. The net
365 differentiation between the *A* and *B* alleles was thus 0.031 per nucleotide over the 12558 bp region
366 (Table S5).

367 Considering the shorter 4194 bp fragment with *Gadus macrocephalus* as the outgroup the divergence
368 was similar (Table S5).

369 Genomic aspects

370 The *Pan I* locus was comprised of 6 exons and 5 introns. We identified six exons in the segment of the
371 *Atxn7l2* locus, and seven exons in the segment containing the *Sort1* locus. .

372 There was a clear difference between sequences classified according to the two *Pan I* allelic variants.
373 There were 121 fixed substitutions between *Pan I* allelic variants out of a total of 349 variable sites
374 found in the entire region (Table 1 and Figure 1). Out of 121 fixed substitutions, eight were non-
375 synonymous, seven were synonymous, and 106 were in non-coding regions. Six non-synonymous and
376 three synonymous substitutions were at the *Pan I* locus, with two codons showing multiple (two) non-
377 synonymous substitutions each and two other codons with one non-synonymous substitution each. Of
378 those two codons with multiple non-synonymous substitutions (codons 61 and 64 in Table 5 of Pogson,
379 2001), one had C and A, and A and T nucleotides at the first and third position, respectively, for the *A*
380 and *B* allelic sequence variants (codons CAA and AAT respectively). At this same codon most gadoid
381 species sequenced by Pogson and Mesa (2004) including *G. macrocephalus*/*G. ogac* had A in both first
382 and third position (codon AAA). The other codon, number 64, had A and G, and C and A nucleotides at
383 the first and second codon position, respectively (codons ACC and GAC), for the *A* and *B* allelic sequence
384 variants (Figure 1). At this same codon most gadoid species sequenced by Pogson and Mesa (2004) had
385 G and A in the first and second position (codon GAC) whereas *G. ogac* (which Pogson, 2001, used as the
386 outgroup for the *A* and *B* alleles of *Pan I*) and *G. macrocephalus* had an A in both the first and second
387 position (codon AAC). Each of the *Sort1* and *Atxn7l2* loci had one non-synonymous substitution and
388 these loci had one and three synonymous substitutions, respectively (Figure 1). The fixed substitutions
389 located furthest apart were 12088 bp apart (Figure 1).

390 We looked for and analyzed the *DraI*, *BstEII*, *BstXI*, *PstI*, and *SacII* restriction sites referred to by
391 Pogson (2001). In his Figure 1 *BstEII* and *PstI* are 5.7 kb apart on either side of the *DraI* site and in strong
392 LD with the *DraI* site and with each other. The *DraI* and *BstEII* sites were fixed substitutions between
393 *Pan I* sequence variants and *PstI* was polymorphic within *Pan I*^A sequence variants. However, we did not
394 find a *BstEII* site 5' to the *DraI* site as observed by Pogson (2001). Instead, we found a *BstEII* site 3'
395 to the *DraI* site at position 11308 in our sequence. It was also 3' to the *PstI* site at position 11257 in our
396 sequence (Figure 1). This site was in perfect LD with the *A* and *B* *Pan I* alleles and thus behaves much
397 like the *BstEII* site that Pogson (2001) observed.

398 Site frequency spectra and coalescent models

399 The unfolded site frequency spectrum of the 4.2 kb region of *Pan I* and *Atxn7l2* with Pacific cod used as
400 the outgroup is in Figure 4. There were three peaks in the spectrum, singletons, sextuplets, and twenty-five
401 tuplets. The two latter peaks of 25 and 15 sites respectively represent the fixed differences between
402 sequences classified according to the *B* and *A* *Pan I* alleles respectively. The Kingman coalescent model
403 did not fit well. The Beta ($2 - \alpha, \alpha$) and the point-mass coalescent models gave a better fit, in particular

404 for the singletons. None of the coalescent models could account for the high frequency of sextuplets and
 405 twenty-five tuplets.

406 The folded site frequency spectrum of the entire region (Figure 5; folded because we did not have an
 407 outgroup for the whole region) was bimodal, with peaks at singleton sites and combined sextuplet and
 408 twenty-five tuplet in all sites (136 sites). This peak was almost as high as the singleton class. As was the
 409 case for the unfolded spectrum, the Kingman coalescent model gave the worst fit. Both the Beta ($2 - \alpha, \alpha$)
 410 and point mass coalescent models gave better fit to the data except for the high sextuplet/twenty-five tuplet
 411 peak. None of the coalescent models of neutrality predicted the high peak at intermediate frequency.

412 The site frequency spectra of the *A* and *B* alleles separately (Figure S12) were unimodal. Again the
 413 Kingman coalescent model did not fit the data well whereas the Beta ($2 - \alpha, \alpha$) and point-mass coalescent
 414 gave significantly better fits (Table S6).

415 The parameter estimates for the Beta ($2 - \alpha, \alpha$) and point mass multiple-merger coalescent models
 416 are in Figure S13. The α parameter for the *A* and *B* alleles were similar to those for the *Myg* and *Hba2*
 417 genes (Árnason and Halldórsdóttir, 2015) slightly above 1.0. However, for the combined data $\alpha = 1.475$.
 418 A similar effect was observed for the ψ parameter which was similar ($\psi = 0.245$) for the sequences
 419 classified according to *Pan I A* alleles as for the *Myg* and *Hba2* loci. Sequences classified according
 420 to the *Pan I B* alleles had an even higher $\psi = 0.325$. For the combined data was considerably lower or
 421 $\psi = 0.100$.

422 DISCUSSION

423 Function of proteins

424 The *Pan I* codes for pantophysin, a protein whose function is involved in vesicle transport pathways in
 425 adipocytes, especially in the trafficking of insulin-regulated glucose transporter GLUT4 (reviewed in
 426 Bradley et al., 2001). Thus it is likely to be involved in energy metabolism, possibly burst activity. The
 427 allelic variants of *Pan I* have been associated to behavioral profiles with the *Pan I^A* allele connected to
 428 shallow waters and seasonal temperature changes while the *Pan I^B* allele is connected to deeper waters and
 429 steep temperature fluctuations (Pampoulie et al., 2008). The differences in *Pan I^A* allelic frequencies at
 430 different geographic scales have been connected to temperature and salinity Case et al. (2005). Sortilin is
 431 a major protein component of Glut4-containing microvesicles that might be involved in the translocation
 432 or biogenesis of the GLUT4-containing vesicles (Lin et al., 1997). Sortilin is also involved in trafficking
 433 processes at the Golgi apparatus and plasma membrane (Strong et al., 2012) whose expression is connected
 434 to hepatic reduction in triglycerides and to obesity (Ai et al., 2012). Thus it also seems involved in energy
 435 metabolism. Atxn712 codes for a protein involved in chromatin remodeling activities (Marchler-Bauer
 436 et al., 2012; Zhao et al., 2008). Chromatin dynamics have been documented to act as a control of gene
 437 expression and show a response to stress episodes mediated by e.g. temperature or salinity (de Nadal
 438 et al., 2011), the very same drivers that *Pan I* has been linked to Case et al. (2005). The attributes of
 439 the proteins thus suggest on one hand co-location and shared metabolic pathways of Glut4-containing
 440 vesicles for pantophysin and sortilin, and on the other hand shared correlations to steep fluctuations in
 441 temperature or depth-related environmental vectors among pantophysin and ataxin-12.

442 Allelic divergence and spatial differentiation

443 There is a deep divergence of the *A* and *B* alleles only a little less than the divergence of Atlantic cod
 444 and Pacific cod. This is in line with results of Pogson and Mesa (2004) who found that the *A* and *B* split
 445 predated the divergence of Atlantic cod and Walleye pollock *Gadus chalcogrammus*. Using mitogenomic
 446 data Coulson et al. (2006) date the Atlantic cod vs. Pacific cod split at 4 mya and the Atlantic cod vs.
 447 Walleye pollock split at 3.8 mya. Accordingly the *A* vs *B* divergence in between those date, perhaps
 448 3.9 mya. However, these dates are based on the Kingman coalescent. Times scales under the more
 449 appropriate multiple merger coalescents considered here may be considerably shorter (Árnason and
 450 Halldórsdóttir, 2015). Furthermore, if the *A* and *B* divergence is driven by repeated selective sweeps
 451 within each haplotype Pogson (2001) and strong selection time may be shorter.

452 Árnason and Halldórsdóttir (2015) considered as one possible explanation a historical hypothesis of
 453 ancient isolation and recent admixture for the *Ckma* gene in Atlantic cod. Their evidence did not support
 454 the historical hypothesis. We can use our results to shed further light on the issue. The *Ckma* gene shows
 455 large differentiation between a region that Árnason and Halldórsdóttir (2015) called South (Faroe Islands,
 456 North Sea, Baltic Sea, Celtic Sea, and Irish Sea) and North (Canada, Greenland, Iceland, Norway, Barents

457 Sea, and White Sea) with highly significant pairwise $F_{ST} \approx 0.8$ between North and South localities and
458 no differentiation between localities within each region. Pogson and Fevolden (2003) devised a test of the
459 historical vs selection hypothesis (Árnason and Pálsson, 1996b) of coastal vs North East Arctic cod in
460 northern Norway. They stated that patterns of neutral variation within the *A* allelic class of *Pan I* would
461 be a sensitive indicator of the historical hypothesis. Isolation and admixture are part of the breeding
462 structure of a population with genome-wide effects (Wright, 1931). Different genomic regions should be
463 concordant in their behavior (Bernardi et al., 1993) both neutral genes under random drift and selective
464 genes. However, supposedly neutral variation within the *A* haplogroup of *Pan I* and neighboring loci is
465 not congruent with the North vs South divide considered by Árnason and Halldórsdóttir (2015). Instead
466 the differentiation is more east vs west. Thus spatial differentiation in Atlantic cod probably is primarily
467 driven by selection (c.f. Bradbury et al., 2010) and not by history.

468 **Balancing selection**

469 Our evidence strongly suggests selective effects at the *Pan I* locus and its peripheral regions, partial
470 segments of the *Sort1* and *Atxn7l2* loci. Our evidence also points to the *Pan I* locus as one target of
471 selection.

472 The patterns of distribution of polymorphism at both site frequency spectra clearly indicate departure
473 from neutrality and the action of balancing selection. The patterns that we find in these 4.2 kb and 12.5
474 kb regions are in agreement with the findings of Pogson (2001) at a 1.85 Kb region of *Pan I*. The *Pan I*
475 locus contains an ancient polymorphism undergoing a mixture of directional and balancing selection that
476 has maintained two highly differentiated alleles (Pogson, 2001). In the unfolded site frequency spectrum
477 covering 4.2 kb, the signature of balancing selection is in the form of high frequency sixtuplets and
478 twenty-five tuplets that do not fit any of the theoretical expectations of neutrality. The high frequency
479 peaks of the spectrum are at opposite frequencies (6/31 and 25/31) and correspond to the differentiation
480 between the six *Pan I^A* and the 25 *Pan I^B* sequence variants (i.e. six and 25 4.2 kb sequences classified
481 according to the long-lived polymorphism maintained by balancing selection as proposed by Pogson,
482 2001). For the whole 12.5 kb region, the folded site frequency spectrum also exhibits the signature of
483 balancing selection as the conflated sixtuples/twenty-five tuples peak clearly surpasses all theoretical
484 expectations. It is almost as common as the singleton class. The unfolded and folded spectra accord
485 well with each other. In the unfolded spectrum, the sixtuplets and twenty-five tuplets represent 40 fixed
486 differences between *Pan I^B* and *Pan I^A* sequence variants (25 and 15 differences, respectively) over a
487 4.2 kb region in 31 sequences, i.e. 0.0095 fixed differences per site. In the folded one, the conflated
488 sixtuples/twenty-five tuplets have a frequency of 136 differences over a 12.5 kb region in 31 sequences,
489 i.e. 0.0108 fixed differences per site. The site frequency spectra also match the phylogeny of alleles. The
490 frequency of 25 and 15 sites of the sixtuplets and twenty-five tuplets respectively mean that the *Pan I^B*
491 sequence variants are 10 sites further from the Pacific cod outgroup than are the *Pan I^A* sequence variants.
492 Thus the *Pan I^B* have evolved further from the outgroup as seen in the phylogeny.

493 From the coalescent models, the better fit to the site frequency spectra that we observe with the Beta
494 ($2 - \alpha, \alpha$) and point-mass coalescent than with the Kingman coalescent is in agreement with observations
495 by Birkner et al. (2013) in mitochondrial DNA data (Árnason et al., 2000; Sigurgíslason and Árnason,
496 2003; Árnason, 2004) of Atlantic cod. The better fit to the site frequency spectra by the more generalized
497 coalescent models than the Kingman coalescent is most likely because the last considers bifurcated
498 coalescent events only, while the generalized models allow for multiple merger coalescent events and
499 accommodate large variance in the number of offspring (Eldon and Wakeley, 2006). Thus, the better fit of
500 the generalized coalescent models to the observed site frequency spectra, especially capturing the high
501 frequency of singletons, may indicate large variance in offspring numbers in Atlantic cod and sweepstakes
502 reproduction (Árnason and Halldórsdóttir, 2015). Such large variance is likely due to then the occurrence
503 of frequent small and infrequent large offspring reproductive events as otherwise there would be no
504 genetic variation (Árnason, 2004). The parameter estimate α represents the probability of large offspring
505 events (i.e. large families), which is most likely as it approaches $\alpha = 1$. As α approaches $\alpha = 2$ the
506 model behaves like the Kingman coalescent; the parameter ψ points to the proportion of reproduction that
507 can be ascribed to an individual, where the model behaves like the Kingman coalescent when $\psi = 0$ and
508 multiple merger coalescent events are predominant when $\psi = 1$ (Birkner et al., 2013). An indicator of
509 balancing selection is that α is larger when we consider both allelic sequence variants *A* and *B* combined
510 than when we consider each allelic variant separately; i.e. when considering both *A* and *B* combined α

511 corresponds to a coalescent behavior that tends to longer coalescent times (with accumulated mutations)
 512 between alleles. In contrast, when we consider each allelic variant separately α has a lower value tending
 513 to faster coalescent times between alleles. From ψ , we also conclude that it indicates balancing selection
 514 as the combined data (both *A* and *B*) show ψ values that are lower and indicate longer coalescent times
 515 than when we consider each allelic sequence variant separately. These observations have a parallel on
 516 nucleotide variability in that the combined data have larger $\hat{\pi}$ and $\hat{\theta}$ than each allelic variant by itself.
 517 This is also an indication of balancing selection.

518 Concerning nucleotide variability, the signature of balancing selection is recognized in the patterns
 519 of silent divergence and polymorphism at *Pan I* and its neighboring genes. Balancing selection has
 520 been detected by looking at peaks of silent diversity among alleles that have experienced amino acid
 521 substitutions underlying adaptation to different environmental niches (Storz et al., 2007). Here, we see that
 522 the highly diverged sequence variants *A* and *B* (classified according to *Pan I*) have undergone amino acid
 523 replacements (six at *Pan I* and one at each of the neighboring regions) which strongly suggests functional
 524 differences in protein products. We have a series of peaks of silent divergence among sequence variants *A*
 525 and *B* larger than the polymorphism within each variant class. Thus the time to coalescence is longer and
 526 there are accumulated silent mutations among *A* and *B* sequence variants, and there is a shorter time to
 527 coalescence and less neutral mutations within each sequence variant class. This scenario fits in an iterative
 528 fashion that seen by Storz et al. (2007), who observed that the silent diversity among functional variants
 529 underlying altitude adaptation in deer mice was older and larger than the silent diversity segregating
 530 within such variants, signaling the works of diversifying selection (Storz et al., 2007). In this study, the
 531 larger levels of divergence among sequence variants *A* and *B* than the levels of polymorphism within each
 532 sequence variant separately is also evident in the larger measurements of $\hat{\theta}$ and $\hat{\pi}$ of the combined data
 533 than for each haplogroup separately. The same effect becomes apparent in the heterozygosity which is
 534 twofold and sixfold larger for the data combined than for the sequence variants classified according to
 535 the *Pan I^A* and *Pan I^B* alleles. Although these patterns point to selective effects at the *Pan I* gene and its
 536 neighboring genes, we can not rule out admixture (Bernardi et al., 1993).

537 The high LD observed throughout the entire 12.5 kb region is in agreement with the high LD detected
 538 by Therkildsen et al. (2013) in vast genomic tracts using genome scans with sets of gene-associated SNP
 539 genotyping. Our genomic horizon is narrower but at a higher resolution, revealing a very fine-grained
 540 LD at *Pan I* and flanking loci. The multiple peaks of divergence among sequence variants are a signal
 541 of balancing selection (Storz et al., 2007), and the iterative nature of the signal is strengthened by the
 542 fine-grained and high LD levels. As nearly all the variation is so tightly connected at *Pan I* and flanking
 543 loci, a signature of selection as peaks of divergence over relatively lower levels of polymorphism is seen
 544 repeatedly throughout the 12.5 kb region.

545 The summary statistics in general were non-significant, but some summary statistics for neutrality
 546 tests such as Tajima's *D* and Fu and Li's *D** and *F** are sensitive to low sample sizes (Guinand et al., 2004,
 547 revised in). Pogson (2001) has demonstrated that strong departures from neutrality are not necessarily
 548 reflected in test statistics.

549 The evidence indicates that there may have been bursts of non-synonymous substitutions at *Pan I*
 550 locus. This implies that *Pan I* locus is at least one focus of selection. The sites that contribute to the
 551 differentiation of the sequence variants (i.e. fixed differences among sequence variants) are more or less
 552 spread throughout the whole 12.5 kb region. However, a majority of the amino acid replacements sites
 553 that are fixed between the *A* and *B* haplogroups are located in the *Pan I* gene. They represent radical
 554 amino acid changes and probably lead to functional differences of the corresponding proteins (Pogson,
 555 2001). The *A* and *B* *Pan I* alleles each have a codon that must have experienced multiple non-synonymous
 556 substitutions relative to the outgroup as already noted by Pogson (2001). Contrasting this codon to other
 557 gadoids sequenced by Pogson and Mesa (2004) suggests that for codon 61 (the codon numbers refer to
 558 Pogson and Mesa, 2004, sequence) *G. ogac* is the same as the ancestral allele (as Pogson, 2001, assumed)
 559 and that the *A* and *B* variants of *Pan I* in *G. morhua* both carry derived alleles. This implies adaptive
 560 substitutions in both the *A* and *B* *Pan I* Atlantic cod lineages. It also suggests that in codon 64 the *B*
 561 variant of *Pan I* in *G. morhua* carries the ancestral allele shared by most other gadoids while the *A* variant
 562 of *Pan I* and *G. ogac* both carry derived alleles at the first codon position, implying an older adaptive
 563 substitution, prior to the separation of *G. morhua* and the *G. macrocephalus/G. ogac* lineage (Coulson
 564 et al., 2006). However, the *Pan I A* has an additional adaptive change in the second position of the
 565 codon. Thus *G. ogac* is an appropriate outgroup in codon 61 but not as a distant ancestor for codon 64.

566 These codons with multiple non-synonymous substitutions at *Pan I* constitute bursts of non-synonymous
567 substitutions in the same lineage (Atlantic cod) and show that the differences between the *A* and *B*
568 haplogroups occur as adaptive changes on both lineages. There is build-up of differences over time.
569 This is an ancient balanced polymorphism (Charlesworth, 2006) and not simply a partial selective sweep
570 bringing a particular chromosomal region to a high frequency (for example as seen in human G6PD and
571 β globin polymorphisms Verrelli et al., 2002; Currat et al., 2002). This in conjunction with the number of
572 other non-synonymous changes in the same region of the *Pan I* gene implies that the *Pan I* gene is most
573 likely a focal point of selection. Due to the nature of the genetic code certain amino acid substitutions
574 cannot occur without going through intervening amino acid states. Thus if selection favors an amino
575 acid change of that type, there will be accumulation of non-synonymous mutations within respective
576 codons. Thus codons with multiple non-synonymous mutations in the same lineage (what is referred
577 to as evolutionary bursts Gillespie, 1984) are signals of focal points of selection (Bazykin et al., 2004,
578 2006). The intervening state for the codons we observe with multiple non-synonymous mutations are the
579 corresponding codons observed in *G. ogac* by Pogson (2001), thus signaling the operation of balancing
580 selection at *Pan I* in Atlantic cod. Hughes (2007) has criticized studies that look for the operation of
581 selection under the criteria of concentration of amino acid replacements within a limited region. However,
582 our main focus is not the concentration of amino acid replacements at a gene, but rather the occurrence
583 of multiple codons with multiple non-synonymous substitutions at divergent haplogroups observed at a
584 particular gene.

585 The entire 12.5 kb region of *Pan I*, *Sort1*, and *Atxn7l2*) is located within a much larger region
586 of a genomic island of divergence (Bradbury et al., 2013; Hemmer-Hansen et al., 2013). This genomic
587 island of divergence has been connected (Hemmer-Hansen et al., 2007) to behavioral differences related
588 to different habitat use with respect to temperature and depth regimes (Pampoulie et al., 2008). Recently,
589 *Pan I* ecotypes in cod have been associated to polymorphisms at the Rhodopsin (*Rh1*) locus also located in
590 LG1 as *Pan I*, with potential involvement of behavioral differences and visual capabilities as rhodopsin is
591 a pigment involved in dim-light vision (Pampoulie et al., 2015). There is evidence for the build-up of the
592 two haplogroups, two functionally balanced types, by selection as already stated. Neutral variation will
593 accumulate on the two genealogical lineages (Charlesworth et al., 2003). Utilization of different habitats
594 with complex multidimensional differences may entail complex phenotypic differences with bearings on
595 genomic structures known as supergenes (Thompson and Jiggins, 2014). The implication of balancing
596 selection, the prevalence of divergence and of high levels of fine-grained LD, and a possible correlation in
597 function suggested by protein function at the loci in the region, together hint the build up of a supergene
598 inclusive of the region where *Pan I* and flanking segments of *Sort1* and *Atxn7l2* are located. The tightly
599 knit LD throughout the 12.5 kb region is likely a product of the selective effects detected throughout the
600 whole region, and seemingly little recombination among the *A* and *B* sequence variants classified by the
601 *Pan I* alleles.

602 The effects of balancing selection at a single locus will extend only short distances from the selected
603 sites with free recombination (Wiuf and Hein, 1999). Signs of a long standing balanced polymorphism
604 therefore are the result of a build-up of co-adapted complexes of epistatic interactions among multiple
605 sites or due to suppression of recombination (Wiuf and Hein, 1999). The very high LD observed here and
606 the peculiar site frequency spectra with peaks at exactly opposite frequencies and no variation around
607 the peaks (c.f. Árnason and Halldórsdóttir, 2015) imply very little recombination. We suggest that
608 the genomic island of divergence is a supergene of co-adapted complexes possibly locked together by
609 structural variation (Joron et al., 2011; Thompson and Jiggins, 2014). There may well be multiple selective
610 sites within the genomic island. *Pan I* is very likely one such site.

611 REFERENCES

- 612 Ai, D., Baez, J. M., Jiang, H., Conlon, D. M., Hernandez-Ono, A., Frank-Kamenetsky, M., Milstein, S.,
613 Fitzgerald, K., Murphy, A. J., Woo, C. W., Strong, A., Ginsberg, H. N., Tabas, I., Rader, D. J., and Tall,
614 A. R. (2012). Activation of ER stress and mTORC1 suppresses hepatic sortilin-1 levels in obese mice.
615 *Journal of Clinical Investigation*, 122(5):1677–1687.
- 616 Árnason, E. (2004). Mitochondrial cytochrome *b* DNA variation in the high fecundity Atlantic cod:
617 Trans-Atlantic clines and shallow gene-genealogy. *Genetics*, 166:1871–1885.
- 618 Árnason, E. and Halldórsdóttir, K. (2015). Nucleotide variation and balancing selection at the *Ckma* gene
619 in Atlantic cod: analysis with multiple merger coalescent models. *PeerJ*, 3:e786.
- 620 Árnason, E., Hernandez, U. B., and Kristinsson, K. (2009). Intense habitat-specific fisheries-induced
621 selection at the molecular *Pan I* locus predicts imminent collapse of a major cod fishery. *PLoS ONE*,
622 4(5):e5529.
- 623 Árnason, E. and Pálsson, S. (1996a). Mitochondrial cytochrome *b* DNA sequence variation of Atlantic
624 cod, *Gadus morhua*, from Norway. *Molecular Ecology*, 5:715–724.
- 625 Árnason, E. and Pálsson, S. (1996b). Mitochondrial cytochrome *b* DNA sequence variation of Atlantic
626 cod *Gadus morhua*, from Norway. *Molecular Ecology*, 5(6):715–724.
- 627 Árnason, E., Petersen, P. H., Kristinsson, K., Sigurgíslason, H., and Pálsson, S. (2000). Mitochondrial
628 cytochrome *b* DNA sequence variation of Atlantic cod from Iceland and Greenland. *Journal of Fish*
629 *Biology*, 56(2):409–430.
- 630 Bazykin, G. A., Dushoff, J., Levin, S. A., and Kondrashov, A. S. (2006). Bursts of nonsynonymous
631 substitutions in HIV-1 evolution reveal instances of positive selection at conservative protein sites.
632 *Proceedings of the National Academy of Sciences of the United States of America*, 103(51):19396–
633 19401.
- 634 Bazykin, G. A., Kondrashov, F. A., Ogurtsov, A. Y., Sunyaev, S., and Kondrashov, A. S. (2004). Pos-
635 itive selection at sites of multiple amino acid replacements since rat–mouse divergence. *Nature*,
636 429(6991):558–562.
- 637 Bernardi, G., Sordino, P., and Powers, D. A. (1993). Concordant mitochondrial and nuclear DNA
638 phylogenies for populations of the teleost fish *Fundulus heteroclitus*. *Proceedings of the National*
639 *Academy of Sciences of the United States of America*, 90(20):9271–9274.
- 640 Birkner, M., Blath, J., and Eldon, B. (2013). Statistical properties of the site-frequency spectrum associated
641 with Λ -coalescents. *Genetics*, 195:1037–1053.
- 642 Birnboim, H. and Doly, J. (1979). A rapid alkaline extraction procedure for screening recombinant
643 plasmid DNA. *Nucleic Acids Research*, 7(6):1513–1523.
- 644 Bradbury, I. R., Hubert, S., Higgins, B., Borza, T., Bowman, S., Paterson, I. G., Snelgrove, P. V. R.,
645 Morris, C. J., Gregory, R. S., Hardie, D. C., Hutchings, J. A., Ruzzante, D. E., Taggart, C. T., and
646 Bentzen, P. (2010). Parallel adaptive evolution of Atlantic cod on both sides of the Atlantic Ocean in
647 response to temperature. *Proceedings of the Royal Society, B*, 277:3725–3734.
- 648 Bradbury, I. R., Hubert, S., Higgins, B., Bowman, S., Borza, T., Paterson, I. G., Snelgrove, P. V. R.,
649 Morris, C. J., Gregory, R. S., Hardie, D., Hutchings, J. A., Ruzzante, D. E., Taggart, C. T., and Bentzen,
650 P. (2013). Genomic islands of divergence and their consequences for the resolution of spatial structure
651 in an exploited marine fish. *Evolutionary Applications*, 6:450–461.
- 652 Bradley, R., Cleveland, K., and Cheatham, B. (2001). The adipocyte as a secretory organ: mechanisms of
653 vesicle transport and secretory pathways. *Recent progress in hormone research*, 56:329–358.
- 654 Bradley, R. K., Roberts, A., Smoot, M., Juvekar, S., Do, J., Dewey, C., Holmes, I., and Pachter, L. (2009).
655 Fast statistical alignment. *PLoS Computational Biology*, 5(5):e1000392+.
- 656 Case, R. A. J., Hutchinson, W. F., Hauser, A. L., Van Oosterhout, A. C., and Carvalho, A. G. R. (2005).
657 Macro- and micro-geographic variation in pantophysin (*pan i*) allele frequencies in NE Atlantic cod
658 *Gadus morhua*. *Marine Ecology Progress Series*, 301:267–278.
- 659 Charif, D. and Lobry, J. R. (2007). SeqinR 1.0-2: a contributed package to the R project for statistical
660 computing devoted to biological sequences retrieval and analysis. In Bastolla, U., Porto, M., Roman,
661 H., and Vendruscolo, M., editors, *Structural approaches to sequence evolution: Molecules, networks,*
662 *populations*, Biological and Medical Physics, Biomedical Engineering, pages 207–232. Springer Verlag,
663 New York. ISBN : 978-3-540-35305-8.
- 664 Charlesworth, B., Charlesworth, D., and Barton, N. H. (2003). The effects of genetic and geographic
665 structure on neutral variation. *Annual Review of Ecology, Evolution and Systematics*, 34:99–125.

- 666 Charlesworth, D. (2006). Balancing selection and its effects on sequences in nearby genome regions.
667 *PLoS Genetics*, 2:e64.
- 668 Coulson, M. W., Marshall, H. D., Pepin, P., and Carr, S. M. (2006). Mitochondrial genomics of gadine
669 fishes: Implications for taxonomy and biogeographic origins from whole-genome data sets. *Genome*,
670 49:1115–1130.
- 671 Cruickshank, T. E. and Hahn, M. W. (2014). Reanalysis suggests that genomic islands of speciation are
672 due to reduced diversity, not reduced gene flow. *Molecular Ecology*, 23:3133–3157.
- 673 Cunningham, K. M., Canino, M. F., Spies, I. B., and Hauser, L. (2009). Genetic isolation by distance and
674 localized fjord population structure in Pacific cod (*Gadus macrocephalus*): Limited effective dispersal
675 in the northeastern Pacific Ocean. *Canadian Journal of Fisheries and Aquatic Sciences*, 66:153–166.
- 676 Currat, M., Trabuchet, G., Rees, D., Perrin, P., Harding, R. M., Clegg, J. B., Langaney, A., and Excoffier,
677 L. (2002). Molecular analysis of the β^S -globin gene cluster in the Niokholo Mandenka population
678 reveals a recent origin of the β^S Senegal mutation. *American Journal of Human Genetics*, 70:207–223.
- 679 de Nadal, E., Ammerer, G., and Posas, F. (2011). Controlling gene expression in response to stress.
680 *Nature Review Genetics*.
- 681 Dereeper, A., Nicolas, S., Le Cunff, L., Bacilieri, R., Doligez, A., Peros, J. P., Ruiz, M., and This, P.
682 (2011). SNIPlay: a web-based tool for detection, management and analysis of SNPs. Application to
683 grapevine diversity projects. *BMC Bioinformatics*, 12(1):134.
- 684 Dray, S. and Dufour, A. (2007). The Ade4 package: Implementing the duality diagram for ecologists.
685 *Journal of Statistical Software*, 22:1–20.
- 686 Edgar, R. C. (2004). MUSCLE: multiple sequence alignment with high accuracy and high throughput.
687 *Nucleic Acids Research*, 32(5):1792–1797.
- 688 Eldon, B. and Wakeley, J. (2006). Coalescent processes when the distribution of offspring number among
689 individuals is highly skewed. *Genetics*, 172:2621–2633.
- 690 Ewing, B., Hillier, L., Wendl, M. C., and Green, P. (1998). Base-calling of automated sequencer traces
691 using Phred. I. Accuracy assessment. *Genome Research*, 8(3):175–185.
- 692 Excoffier, L. and Lischer, H. E. L. (2010). Arlequin Suite Ver 3.5: a new series of programs to perform
693 population genetics analyses under Linux and Windows. *Molecular Ecology Resources*, 10(3):564–567.
- 694 Fevolden, S. E. and Pogson, G. H. (1997). Genetic divergence at the Synaptophysin (*Syp I*) locus
695 among Norwegian coastal and North-east Arctic populations of Atlantic cod. *Journal of Fish Biology*,
696 51:895–908.
- 697 Flicek, P., Amode, M. R., Barrell, D., Beal, K., Billis, K., Brent, S., Carvalho-Silva, D., Clapham, P.,
698 Coates, G., Fitzgerald, S., Gil, L., Girón, C. G., Gordon, L., Hourlier, T., Hunt, S., Johnson, N.,
699 Juettemann, T., Kähäri, A. K., Keenan, S., Kulesha, E., Martín, F. J., Maurel, T., McLaren, W. M.,
700 Murphy, D. N., Nag, R., Overduin, B., Pignatelli, M., Pritchard, B., Pritchard, E., Riat, H. S., Ruffier,
701 M., Sheppard, D., Taylor, K., Thormann, A., Trevanion, S. J., Vullo, A., Wilder, S. P., Wilson, M.,
702 Zadissa, A., Aken, B. L., Birney, E., Cunningham, F., Harrow, J., Herrero, J., Hubbard, T. J., Kinsella,
703 R., Muffato, M., Parker, A., Spudich, G., Yates, A., Zerbino, D. R., and Searle, S. M. (2014). Ensembl
704 2014. *Nucleic Acids Research*, 42(D1):D749–D755.
- 705 Galtier, N., Gouy, M., and Gautier, C. (1996). SEAVIEW and PHYLOWIN: two graphic tools for sequence
706 alignment and molecular phylogeny. *Computer Applications in the Biosciences*, 12(6):543–548.
- 707 Gillespie, J. H. (1984). Molecular evolution over the mutational landscape. *Evolution*, 38(5):1116.
- 708 Gordon, D., Abajian, C., and Green, P. (1998). Consed: a graphical tool for sequence finishing. *Genome
709 Research*, 8(3):195–202.
- 710 Grabowski, T. B., Thorsteinsson, V., McAdam, B. J., and Marteinsdóttir, G. (2011). Evidence of segregated
711 spawning in a single marine fish stock: Sympatric divergence of ecotypes in Icelandic cod? *PLoS ONE*,
712 6(3):e17528.
- 713 Green, P. (1994). Documentation for Phrap. [http://bozeman.mbt.washington.edu
714 /phrap.docs /phrap.html](http://bozeman.mbt.washington.edu/phrap.docs/phrap.html).
- 715 Guinand, B., Lemaire, C., and Bonhomme, F. (2004). How to detect polymorphisms undergoing selection
716 in marine fishes? A review of methods and case studies, including flatfishes. *Journal of Sea Research*,
717 51:167–182.
- 718 Haass, N. K., Kartenbeck, J., and Leube, R. E. (1996). Pantophysin is a ubiquitously expressed Synapto-
719 physin homologue and defines constitutive transport vesicles. *Journal of Cell Biology*, 134:731–746.
- 720 Hedrick, P. W. (1987). Genetic disequilibrium measures: Proceed with caution. *Genetics*, 117:331–341.

- 721 Helmlinger, D., Hardy, S., Abou-Sleymane, G., Eberlin, A., Bowman, A. B., Gansmüller, A., Picaud,
722 S., Zoghbi, H. Y., Trottier, Y., Tora, L., and Devys, D. (2006). Glutamine-expanded Ataxin-7 alters
723 TFTC/STAGA recruitment and chromatin structure leading to photoreceptor dysfunction. *PLoS Biology*,
724 4(3):e67. 16494529[pmid].
- 725 Hemmer-Hansen, J., Nielsen, E. E., Frydenberg, J., and Loeschcke, V. (2007). Adaptive divergence
726 in a high gene flow environment: *Hsc 70* variation in the European flounder (*Platichthys flesus* L.).
727 *Heredity*, 99:592–600.
- 728 Hemmer-Hansen, J., Nielsen, E. E., Therkildsen, N. O., Taylor, M. I., Ogden, R., Geffen, A. J., Bekkevold,
729 D., Helyar, S., Pampoulie, C., Johansen, T., Consortium, F., and Carvalho, G. R. (2013). A genomic
730 island linked to ecotype divergence in Atlantic cod. *Molecular Ecology*, 22:2653–2667.
- 731 Hudson, R. R., Kreitman, M., and Aguadé, M. (1987). A test of neutral molecular evolution based on
732 nucleotide data. *Genetics*, 116(1):153–159.
- 733 Hughes, A. L. (2007). Looking for Darwin in all the wrong places: the misguided quest for positive
734 selection at the nucleotide sequence level. *Heredity*, 99(4):364–373.
- 735 Jombart, T. and Ahmed, I. (2011). Adegnet 1.3-1: New tools for the analysis of genome-wide SNP data.
736 *Bioinformatics*, 27:3070–3071.
- 737 Joron, M., Frezal, L., Jones, R. T., Chamberlain, N. L., Lee, S. F., Haag, C. R., Whibley, A., Becuwe,
738 M., Baxter, S. W., Ferguson, L., Wilkinson, P. A., Salazar, C., Davidson, C., Clark, R., Quail, M. A.,
739 Beasley, H., Glithero, R., Lloyd, C., Sims, S., Jones, M. C., Rogers, J., Jiggins, C. D., and ffrrench
740 Constant, R. H. (2011). Chromosomal rearrangements maintain a polymorphic supergene controlling
741 butterfly mimicry. *Nature*, 477:203–206.
- 742 Karlsen, B. O., Klingan, K., Emblem, Å., Jørgensen, T. E., Jueterbock, A., Furmanek, T., Hoarau, G.,
743 Johansen, S. D., Nordeide, J. T., and Moum, T. (2013). Genomic divergence between the migratory
744 and stationary ecotypes of Atlantic cod. *Molecular Ecology*, 22(20):5098–5111.
- 745 Karlsson, S. and Mork, J. (2003). Selection-induced variation at the Pantophysin locus (*Pan I*) in a
746 Norwegian fjord population of cod (*Gadus morhua*). *Molecular Ecology*, 12:3265–3274.
- 747 Kingman, J. F. C. (1982). The coalescent. *Stochastic Processes and their Applications*, 13:235–248.
- 748 Larance, M., Ramm, G., and James, D. E. (2008). The GLUT4 code. *Molecular Endocrinology*,
749 22(2):226–233.
- 750 Lewontin, R. (1974). *The Genetic Basis of Evolutionary Change*. Columbia University Press, New York,
751 New York.
- 752 Lewontin, R. (1988). On measures of gametic disequilibrium. *Genetics*, 120:849–852.
- 753 Lin, B.-Z., Pilch, P. F., and Kandror, K. V. (1997). Sortilin is a major protein component of Glut4-
754 containing vesicles. *Journal of Biological Chemistry*, 272(39):24145–24147.
- 755 Louis, A., Muffato, M., and Roest Crolius, H. (2013). Genomicus: five genome browsers for comparative
756 genomics in eukaryota. *Nucleic Acids Research*, 41(D1):D700–D705.
- 757 Marchini, J. L. (2013). *popgen: Statistical and Population Genetics*. R package version 1.0-3.
- 758 Marchler-Bauer, A., Zheng, C., Chitsaz, F., Derbyshire, M. K., Geer, L. Y., Geer, R. C., Gonzales, N. R.,
759 Gwadz, M., Hurwitz, D. I., Lanczycki, C. J., Lu, F., Lu, S., Marchler, G. H., Song, J. S., Thanki,
760 N., Yamashita, R. A., Zhang, D., and Bryant, S. H. (2012). CDD: conserved domains and protein
761 three-dimensional structure. *Nucleic Acids Research*, 41(D1):D348–D352.
- 762 McDonald, J. H. (1998). Improved tests for heterogeneity across a region of DNA sequence in the ratio of
763 polymorphism to divergence. *Molecular Biology and Evolution*, 15(4):377–384.
- 764 Muffato, M., Louis, A., Poisnel, C.-E., and Crolius, H. R. (2010). Genomicus: a database and a browser
765 to study gene synteny in modern and ancestral genomes. *Bioinformatics*, 26(8):1119–1121.
- 766 Nielsen, E. E., MacKenzie, B. R., Magnussen, E., and Meldrup, D. (2007). Historical analysis of *Pan i* in
767 Atlantic cod (*Gadus morhua*): temporal stability of allele frequencies in the southeastern part of the
768 species distribution. *Canadian Journal of Fisheries and Aquatic Sciences*, 64(10):1448–1455.
- 769 Pálsson, Ó. K. and Thorsteinsson, V. (2003). Migration patterns, ambient temperature, and growth of
770 Icelandic cod (*Gadus morhua*): evidence from storage tag data. *Canadian Journal of Fisheries and
771 Aquatic Sciences*, 60:1409–1423.
- 772 Pampoulie, C., Jakobsdóttir, K. B., Marteinsdóttir, G., and Thorsteinsson, V. (2008). Are vertical behaviour
773 patterns related to the Pantophysin locus in the Atlantic cod (*Gadus morhua* L.)? *Behavior Genetics*,
774 38:76–81.
- 775 Pampoulie, C., Skirnisdóttir, S., Star, B., Jentoft, S., Jónsdóttir, I. G., Hjórleifsson, E., Thorsteinsson,

- 776 V., Pálsson, Ó. K., Berg, P. R., Andersen, O., Magnúsdóttir, S., Helyar, S. J., and Daníelsdóttir, A. K.
777 (2015). Rhodopsin gene polymorphism associated with divergent light environments in atlantic cod.
778 *Behavioral Genetics*, 45(2):236–244.
- 779 Paradis, E. (2010a). Pegas: an R package for population genetics with an integrated–modular approach.
780 *Bioinformatics*, 26:419–420.
- 781 Paradis, E. (2010b). pegas: an R package for population genetics with an integrated–modular approach.
782 *Bioinformatics*, 26(3):419–420.
- 783 Paradis, E., Claude, J., and Strimmer, K. (2004). APE: analyses of phylogenetics and evolution in R
784 language. *Bioinformatics*, 20:289–290.
- 785 Pogson, G. H. (2001). Nucleotide polymorphism and natural selection at the Pantophysin (*Pan I*) locus in
786 the Atlantic cod, *Gadus morhua* (L.). *Genetics*, 157:317–330.
- 787 Pogson, G. H. and Fevolden, S. E. (2003). Natural selection and the genetic differentiation of coastal
788 and Arctic populations of the Atlantic cod in northern Norway: a test involving nucleotide sequence
789 variation at the Pantophysin (*Pan I*) locus. *Molecular Ecology*, 12(1):63–74.
- 790 Pogson, G. H. and Mesa, K. A. (2004). Positive Darwinian selection at the Pantophysin (*Pan I*) locus in
791 marine gadid fishes. *Molecular Biology and Evolution*, 21(1):65–75.
- 792 Pogson, G. H., Taggart, C. T., Mesa, K. A., and Boutilier, R. G. (2001). Isolation by distance in the
793 Atlantic cod, *Gadus morhua*, at large and small geographic scales. *Evolution*, 55(1):131–146.
- 794 R Development Core Team (2008). *R: A Language and Environment for Statistical Computing*. R
795 Foundation for Statistical Computing, Vienna, Austria.
- 796 Renaut, S., Grassa, C. J., Yeaman, S., Moyers, B. T., Lai, Z., Kane, N. C., Bowers, J. E., Burke, J. M., and
797 Rieseberg, L. H. (2013). Genomic islands of divergence are not affected by geography of speciation in
798 sunflowers. *Nature Communications*, 4:1827.
- 799 Rozas, J., Sánchez-DelBarrio, J. C., Messeguer, X., and Rozas, R. (2003). DnaSP, DNA polymorphism
800 analyses by the coalescent and other methods. *Bioinformatics*, 19:2496–2497.
- 801 Ruegg, K., Anderson, E. C., Boone, J., Pouls, J., and Smith, T. B. (2014). A role for migration-linked
802 genes and genomic islands in divergence of a songbird. *Molecular Ecology*, 23(19):4757–4769.
- 803 Schweinsberg, J. (2003). Coalescent processes obtained from supercritical Galton-Watson processes.
804 *Stochastic Processes and their Applications*, 106:107–139.
- 805 Shin, J.-H., Blay, S., McNeney, B., and Graham, J. (2006). LDheatmap: An R function for graphical dis-
806 play of pairwise linkage disequilibria between single nucleotide polymorphisms. *Journal of Statistical*
807 *Software*, 16(c03):1–9.
- 808 Sigurgíslason, H. and Árnason, E. (2003). Extent of mitochondrial DNA sequence variation in Atlantic
809 cod from the Faroe islands: a resolution of gene genealogy. *Heredity*, 91(6):557–564.
- 810 Star, B., Nederbragt, A. J., Jentoft, S., Grimholt, U., Malmstrom, M., Gregers, T. F., Rounge, T. B.,
811 Paulsen, J., Solbakken, M. H., Sharma, A., Wetten, O. F., Lanzen, A., Winer, R., Knight, J., Vogel,
812 J.-H., Aken, B., Andersen, O., Lagesen, K., Tooming-Klunderud, A., Edvardsen, R. B., Tina, K. G.,
813 Espelund, M., Nepal, C., Previti, C., Karlsen, B. O., Moum, T., Skage, M., Berg, P. R., Gjoen, T.,
814 Kuhl, H., Thorsen, J., Malde, K., Reinhardt, R., Du, L., Johansen, S. D., Searle, S., Lien, S., Nilsen,
815 F., Jonassen, I., Omholt, S. W., Stenseth, N. C., and Jakobsen, K. S. (2011). The genome sequence of
816 Atlantic cod reveals a unique immune system. *Nature*, 477(7363):207–210.
- 817 Storz, J. F., Sabatino, S. J., Hoffmann, F. G., Gering, E. J., Moriyama, H., Ferrand, N., Monteiro, B., and
818 Nachman, M. W. (2007). The molecular basis of high-altitude adaptation in deer mice. *PLoS Genetics*,
819 3(3):e45.
- 820 Strong, A., Ding, Q., Edmondson, A. C., Millar, J. S., Sachs, K. V., Li, X., Kumaravel, A., Wang,
821 M. Y., Ai, D., Guo, L., Alexander, E. T., Nguyen, D., Lund-Katz, S., Phillips, M. C., Morales, C. R.,
822 Tall, A. R., Kathiresan, S., Fisher, E. A., Musunuru, K., and Rader, D. J. (2012). Hepatic sortilin
823 regulates both apolipoprotein B secretion and LDL catabolism. *The Journal of Clinical Investigation*,
824 122(8):2807–2816.
- 825 Therkildsen, N. O., Hemmer-Hansen, J., Als, T. D., Swain, D. P., Morgan, M. J., Trippel, E. A., Palumbi,
826 S. R., Meldrup, D., and Nielsen, E. E. (2013). Microevolution in time and space: SNP analysis of
827 historical DNA reveals dynamic signatures of selection in Atlantic cod. *Molecular Ecology*, 22(9):2424–
828 2440.
- 829 Thompson, M. J. and Jiggins, C. D. (2014). Supergenes and their role in evolution. *Heredity*, 113:1–8.
- 830 Thorsteinsson, V., Pálsson, Ó. K., Tómasson, G. G., Jónsdóttir, I. G., and Pampoulie, C. (2012). Consis-

831 tency in the behaviour types of the atlantic cod: repeatability, timing of migration and geo-location.
832 *Marine Ecology Progress Series*, 462:251–260.

833 Verrelli, B. C., McDonald, J. H., Argyropoulos, G., Destro-Bisol, G., Froment, A., Drousiotou, A.,
834 Lefranc, G., Helal, A. N., Loiselet, J., and Tishkoff, S. A. (2002). Evidence for balancing selection from
835 nucleotide sequence analyses of human *G6PD*. *American Journal of Human Genetics*, 71:1112–1128.

836 Wakeley, J. (2009). *Coalescent Theory: An Introduction*. Roberts & Company Publishers, Greenwood
837 Village, Colorado, USA.

838 Wakeley, J. (2013). Coalescent theory has many new branches. *Theoretical Population Biology*, 87:1–4.

839 Walsh, P. S., Metzger, D. A., and Higuchi, R. (1991). Chelex 100 as a medium for simple extraction of
840 DNA for PCR-based typing from forensic material. *BioTechniques*, 10:506–513.

841 Wiuf, C. and Hein, J. (1999). Recombination as a point process along sequences. *Theoretical Population
842 Biology*, 55:248–259.

843 Wright, S. (1931). Evolution in Mendelian populations. *Genetics*, 16:97–159.

844 Wright, S. I. and Charlesworth, B. (2004). The HKA test revisited: A maximum-likelihood-ratio test of
845 the standard neutral model. *Genetics*, 168(2):1071–1076.

846 Wu, C.-I. (2001). The genic view of the process of speciation. *Journal of Evolutionary Biology*,
847 14(6):851–865.

848 Zhao, Y., Lang, G., Ito, S., Bonnet, J., Metzger, E., Sawatsubashi, S., Suzuki, E., Guezennec, X. L.,
849 Stunnenberg, H. G., Krasnov, A., Georgieva, S. G., Schüle, R., Takeyama, K.-I., Kato, S., LászlóTora,
850 and Devys, D. (2008). A TFTC/STAGA module mediates histone H2A and H2B deubiquitination,
851 coactivates nuclear receptors, and counteracts heterochromatin silencing. *Molecular Cell*, 29(1):92–101.

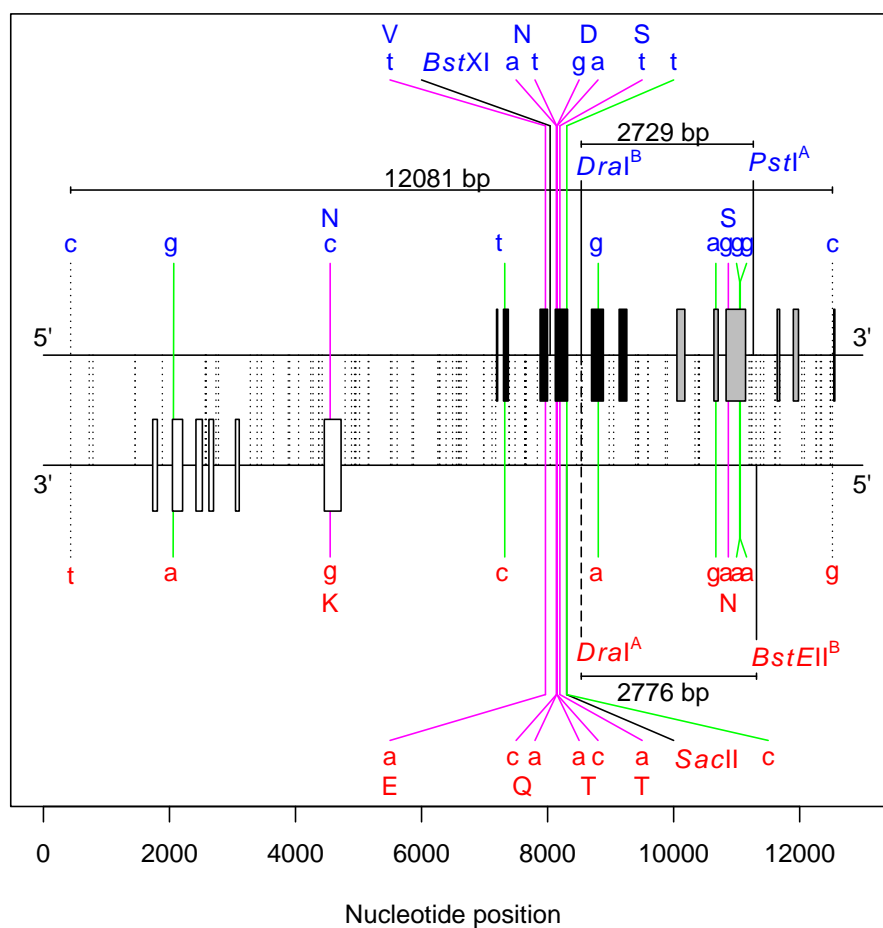


Figure 1. Map of polymorphism of the 12.5 kb region containing the *Pan I* locus and its peripheral regions, the *Sort1* and *Atxn712* loci (partial segments). Boxes represent the exons of *Sort1* (partial segment), *Pan I* and *Atxn712* (partial segment), in white, black and gray, respectively. Variation is displayed with respect to the *DraI* site defining the A and B alleles of the *Pan I* locus (Pogson, 2001). The solid black horizontal lines running through the boxes represent introns (between boxes of the same color) and intergenic space (between boxes of different color). The polymorphic *DraI* restriction site is represented with a solid and a dashed line for the *DraI*^A and *DraI*^B variants, corresponding to presence and absence of recognition site, respectively. *Pan I*^A and *Pan I*^B haplotypes are annotated in red and blue, respectively. Fixed non-synonymous and synonymous substitutions appear as solid vertical lines in magenta and green, respectively. Fixed substitutions in non-coding regions appear as vertical dotted lines between the solid horizontal lines except the outermost sites extend above and below. At the ends of the solid vertical lines, the substitution bases appear as lowercase letters, and the amino acid variants appear as uppercase letters (for non-synonymous substitutions). Restriction sites (cf. Pogson, 2001) appear as black vertical lines.

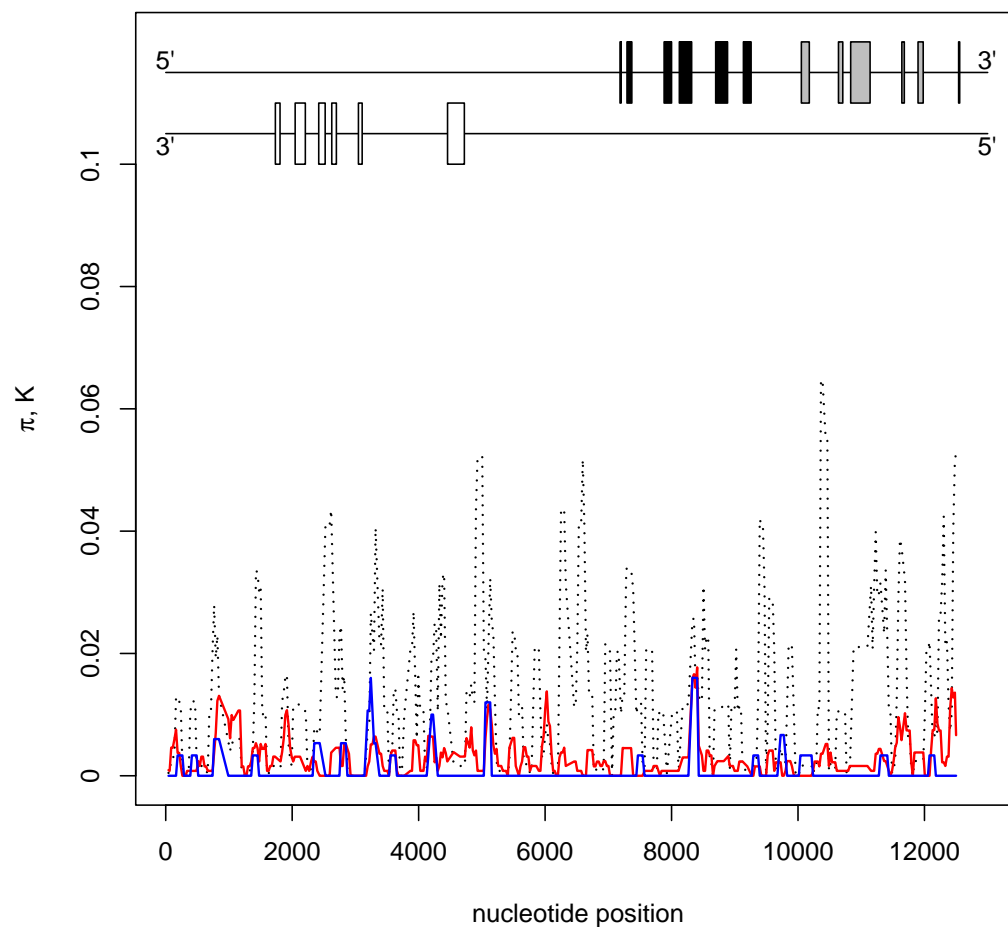


Figure 2. Polymorphism and divergence at the *Pan I* locus and its peripheral regions, the *Sort1* and *Atxn7l2* loci (partial segments). Levels of polymorphism were calculated from silent, intronic, and intergenic sites, with a sliding window size of 100 bp and step size of 25 bp. Divergence (\hat{K}) between *Pan I^A* and *Pan I^B* allelic types is represented by a dotted line. Nucleotide diversity ($\hat{\pi}$) for *Pan I^A* and *Pan I^B* allelic types shown in red and blue, respectively. Boxes represent the exons of *Sort1* (partial segment), *Pan I* and *Atxn7l2* (partial segment), in white, black and gray, respectively. The solid, black, horizontal lines running through the boxes represent introns (between boxes of the same color) and intergenic space (between boxes of different color).

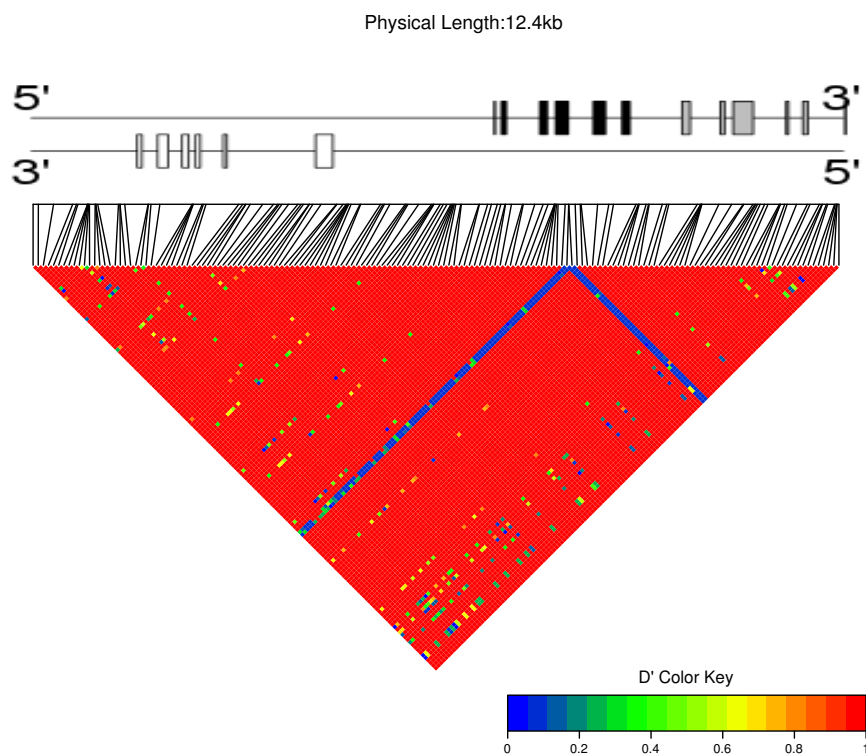


Figure 3. Linkage disequilibrium D' heatmap at high heterozygosity sites of the *Pan I* locus and its peripheral regions, the *Sort1* and *Atxn7l2* loci. Minor allele frequency set at 6/31, the frequency of the *B* alleles of the *Pan I* locus among the 31 samples.

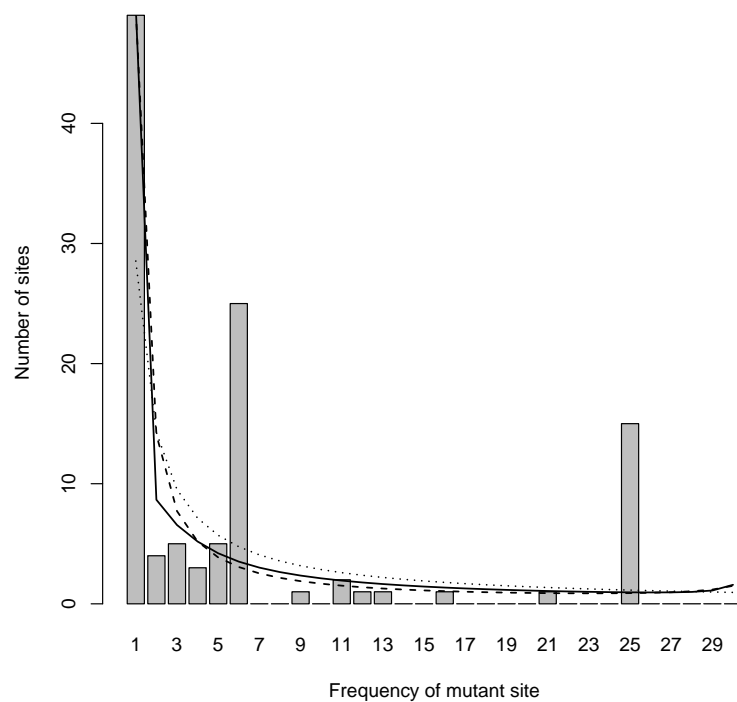


Figure 4. Unfolded site frequency spectrum of Atlantic cod *Pan I* and *Atxn7l2* genes. *Gadus macrocephalus* was used as the outgroup. Number of individuals $n = 31$. Theoretical expectation under Kingman coalescent (dotted line), Beta($2 - \alpha, \alpha$) coalescent (dashed line), and point-mass coalescent (solid line). The bars represent the observed spectrum. The spectrum represents the genetic variability from an alignment of 31 Atlantic cod sequences (25 *Pan I^A* and 6 *Pan I^B*) 4.2 kb long.

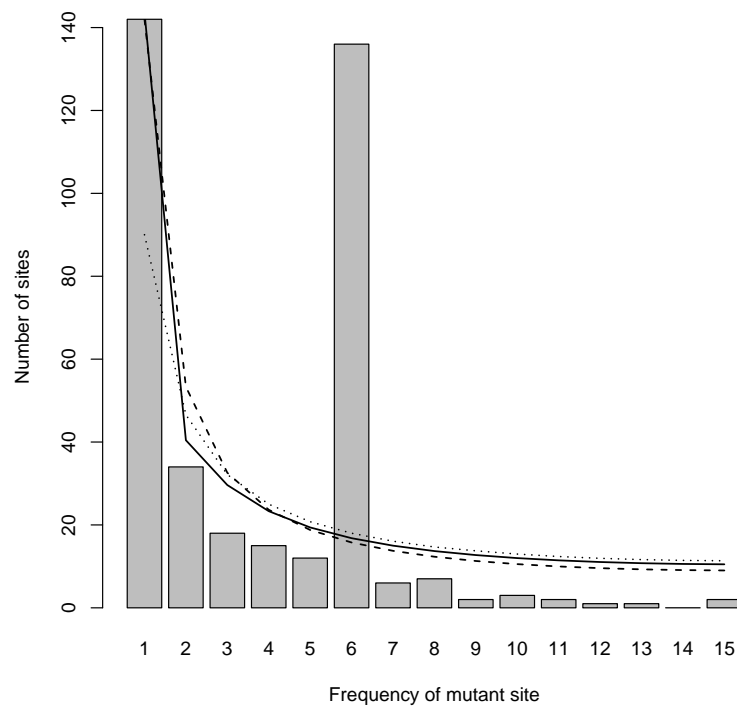


Figure 5. Folded site frequency spectrum of Atlantic cod *Pan I*, *Sort1* and *Atxn7l2* genes. Number of individuals $n = 31$. Theoretical expectation under Kingman coalescent (dotted line), Beta($2 - \alpha$, α) coalescent (dashed line), and point-mass coalescent (solid line). The bars represent the observed spectrum. The spectrum represents the genetic variability from an alignment of 31 Atlantic cod sequences (25 *Pan I*^A and 6 *Pan I*^B) 12.5 kb long.

Table 1. Summary statistics of nucleotide polymorphism at *Pan I* and its peripheric region, the *Sort1* and *Atxn7l2* loci (partial segments). The region analyzed is 12.56 kb. n is the number of sequences used, \hat{S} is the number of segregating sites, \hat{k} is the average number of nucleotide differences, $\hat{\pi}$ is nucleotide diversity, $\hat{\theta}$ (per site) is based on S , \hat{h} is number of haplotypes, and \hat{Hd} is haplotype diversity.

Allelic Type	n	\hat{S}	$\hat{\pi}$	$\hat{\theta}$	\hat{k}	\hat{h}	\hat{Hd}
<i>Pan I^A</i> alleles only	25	209	0.00284	0.00455	34.580	25	1.000
<i>Pan I^B</i> alleles only	6	31	0.00103	0.00109	12.800	6	1.000
<i>Pan I^A</i> and <i>Pan I^B</i> combined	31	349	0.00593	0.00723	71.626	31	1.000