# Long-read sequence capture elucidates the evolution of the hemoglobin gene clusters in codfishes

Siv Nam Khang Hoff*[a], Helle T. Baalsrud*[a‡], Ave Tooming-Klunderud[a], Morten Skage[a], Todd Richmond[b], Gregor Obernosterer[c], Reza Shirzadi[d], Ole Kristian Tørresen[a], Kjetill S. Jakobsen[a] and Sissel Jentoft[a‡]

[a] Centre for Ecological and Evolutionary Synthesis, Department of Biosciences, University of Oslo, Oslo, Norway
[b] Roche NimbleGen Inc., Madison, WI, USA
[c] Roche Diagnostics, Mannheim, Germany
[d] Roche Diagnostics, Oslo, Norway

* these authors contributed equally to this work

‡ Correspondence to: h.t.baalsrud@ibv.uio.no, sissel.jentoft@ibv.uio.no

# Abstract

Combining high-throughput sequencing with targeted sequence capture has become an attractive tool to study specific genomic regions of interest. Most studies have so far focused on the exome using short-read technology. This approach does not capture intergenic regions needed to reconstruct genomic organization, including regulatory regions and gene synteny. In this study, we demonstrate the power of combining targeted sequence capture with long-read sequencing technology, leading to the successful sequencing and assembling of the two hemoglobin (Hb) gene clusters LA (~100kb) and MN (~200kb) across eight species of codfishes that are separated by up to 70 million years of evolution. The highly continuous assemblies – capturing both intergenic and coding sequences – revealed overall conserved genetic organization and synteny of the *Hb* genes within this lineage, yet with several, lineage-specific gene duplications. Moreover, for some of the species examined we identified amino acid substitutions at two sites in the *Hbb1* gene as well as length polymorphisms in its regulatory region, which has previously been linked to temperature adaptation in Atlantic cod populations. Taken together, our study highlights the efficiency of targeted long-read capture for comparative genomic studies by shedding light on the evolutionary history of the *Hb* gene family across the highly divergent group of codfishes.

**Key words:** Targeted sequence capture, comparative genomics, Gadiformes, PacBio sequencing, teleosts

# Article summary

Hemoglobins (Hbs) are key respiratory proteins in most vertebrates. In fishes, Hbs are shown to be of great importance for ecological adaptation, as environmental factors including temperature, directly influences the solubility of $O_2$ in surrounding waters as well as the ability of Hb to bind $O_2$ at respiratory surfaces.

We here combine targeted sequence capture and long-read sequencing to reconstruct and resolve the organization of *Hb* genes and their flanking genes in a selection of codfishes, inhabiting different environmental conditions. Our results shed light on the evolutionary history of *Hb* genes across species separated up to 70 million years of evolution, revealing genetic variations possibly linked to thermal adaptation.

# Introduction

50

51  The rapid advancement of high-throughput sequencing has over the last decade
52  revolutionized genomic research with the increasing numbers of whole genome
53  resources available for multiple vertebrate species, including the diverse group of
54  teleost fishes (Volff 2005; Ellegren 2014; Goodwin *et al.* 2016; Malmstrøm *et al.* 2016;
55  2017). However, whole genome sequencing (WGS) and generation of high quality
56  genome assemblies are still considered costly and time-consuming (Jones and Good
57  2015). For investigations concerning specific genomic regions there is no need for
58  complete genome information, which has spurred the development of reduction
59  complexity approaches such as targeted sequence capture (Teer *et al.* 2010; Grover *et*
60  *al.* 2012; Samorodnitsky *et al.* 2015). The basic idea of targeted sequence capture
61  involves design of specific probes covering the particular genomic area of interest
62  generating an enriched coverage of the targeted sequences (Turner *et al.* 2009; Grover
63  *et al.* 2012). Most studies using targeted sequence capture have to a large extent been
64  directed towards the exome, often supported by the existence of a reference genome
65  (Ng *et al.* 2009; Broeckx *et al.* 2014; Yoshihara *et al.* 2016), or transcriptome assemblies
66  (Syring *et al.* 2016). Recent reports have, however, been focusing on off-target
67  sequences in noncoding regions (Guo *et al.* 2012; Samuels *et al.* 2013; Syring *et al.*
68  2016; Yoshihara *et al.* 2016; Morin *et al.* 2016), as they may contain crucial regulatory
69  elements varying in sequence and length between populations or species, and could
70  be of functional and evolutionary importance (Woolfe *et al.* 2004; Patrushev and
71  Kovalenko 2014).
72       To our knowledge, sequence capture studies have so far been based on short-
73  read sequencing technology (George *et al.* 2011; Samorodnitsky *et al.* 2015; Li *et al.*
74  2015; Bragg *et al.* 2016), and construction of continuous sequences enabling resolution
75  of gene organization have therefore not been possible. Comparative genetic studies
76  of gene organization or synteny requires longer, more continuous stretches of DNA
77  containing more than one gene (Huddleston *et al.* 2014). By its ability to span long
78  stretches of repeats, long-read sequencing technology has been successfully applied
79  to improve genome assembly statistics and generation of highly continuous genome
80  assemblies for a growing number of species (English *et al.* 2012; Kim *et al.* 2014;
81  Bickhart *et al.* 2017; Korlach *et al.* 2017; Tørresen *et al.* 2017a; Tørresen *et al.* 2017b). For
82  example, incorporation of long PacBio reads resulted in a significantly improved
83  version of the Atlantic cod (*Gadus morhua*) genome assembly, i.e. a 50-fold increase in
84  sequence continuity and a 15-fold reduction in the proportion of gaps (Tørresen *et al.*
85  2017b). Correspondingly, utilizing long-read sequencing technology in combination

86    with targeted capture could yield longer continuous assemblies of specific genomic

87    regions of interest, allowing in-depth comparative genetic studies, even in species

88    where a reference genome is not available.

89        In fishes, the hemoglobin (*Hb*) gene family, encoding the protein subunits Hba

90    and Hbb, has shown to be of importance for ecological adaptation, as environmental

91    factors such as temperature directly influences the ability of Hb to bind $O_2$ at

92    respiratory surfaces and its subsequent release to tissues (Wells 2005). In a recent

93    report, a full characterization of the *Hb* gene repertoire – using comparative

94    genomics analysis – uncovered a remarkably high *Hb* gene copy variation within the

95    codfishes (Baalsrud *et al.* 2017). Notably, a negative correlation between the number

96    of *Hb* genes and depth of which the species occur was observed, suggesting that the

97    more variable environment in epipelagic waters have facilitated a larger and more

98    diverse *Hb* gene repertoire, which was supported by evidence of diversifying

99    selection (Baalsrud *et al.* 2017). Furthermore, in Atlantic cod, two tightly linked

100    polymorphisms at amino acid positions 55 and 62 of the Hbb1-globin – suggested to

101    be associated with thermal adaptation – exhibit a latitudinal cline in allele frequency

102    in populations inhabiting varying temperature and oxygen regimes in the North

103    Atlantic and Baltic Sea (Andersen *et al.* 2009). Populations found in the southern

104    regions display the Hbb1-1 variant (Met55Lys62), whereas more northern

105    populations largely display the Hbb1-2 variant (Val55Ala62) (Andersen *et al.* 2009).

106    The Hbb1-1 variant has been shown to be insensitive to temperature whereas Hbb1-2

107    is temperature dependent with a higher $O_2$ affinity than Hbb1-1 at colder

108    temperatures (Karpov and Novikov 1980; Andersen *et al.* 2009), however, this has

109    been questioned by (Barlow *et al.* 2017). Additionally, an indel polymorphism within

110    the promoter of the *Hbb1* gene has been reported to be in linkage disequilibrium with

111    the above-mentioned polymorphisms (Star *et al.* 2011). Examination of multiple

112    Atlantic cod populations uncovered that a longer promoter variant is associated with

113    *Hbb1-2* and found to up-regulate its gene expression at higher temperatures, i.e.

114    aiding in the maintenance of total oxygen-carrying capacity (Star *et al.* 2011).

115        In teleosts, the *Hb* genes are found to reside at two distinct genomic regions –

116    the MN and LA cluster. Earlier reports have shown that there is a high evolutionary

117    turnover of *Hb* genes across teleosts, with lineage-specific duplications and losses,

118    which is in stark contrast to genes flanking the *Hb* genes, where the synteny is highly

119    conserved (Quinn *et al.* 2010; Opazo *et al.* 2013; Feng *et al.* 2014). In this study, the

120    overall goal was to elucidate the evolutionary past of the *Hb* clusters – including *Hb*

121    genes, flanking genes and intergenic sequences – within the phylogenetically diverse

122    group of codfishes (Gadiformes) by taking the advantage of long read sequencing

4

123    technology combined with targeted sequence capture. Eight codfish species were

124    carefully selected on the basis of both phylogenetic and habitat divergence, implying

125    that they are exposed to a variety of environmental factors as well as displaying

126    distinct life-history traits. A highly continuous genome assembly of Atlantic cod

127    (Tørresen *et al.* 2017b) as well as low coverage draft genome assemblies of all eight

128    species (Malmstrøm *et al.* 2017) were used in the design of the probes covering the

129    genomic regions of interest. To enable targeted sequence capture for PacBio RSII

130    sequencing, we modified the standard protocol for sequence capture offered by

131    NimbleGen, i.e. the SeqCap EZ (Roche NimbleGen), as well as generating custom-

132    made barcodes. This combined approach resulted in successful capturing and

133    assembling of the two *Hb* gene clusters across the codfishes examined. The

134    generation of highly continuous assemblies – for most of the species – enabled

135    reconstruction of micro-synteny revealing lineage-specific gene duplications and

136    identification of a relatively large and inter-species variable indel located in the

137    promoter region between the *Hbb1* and *Hba1* genes.

138        Our study demonstrates that this approach, combining sequence capture

139    technology with long-read sequencing is a highly efficient and versatile method to

140    investigate specific genomic regions of interest – with respect to micro-synteny,

141    regulatory regions and genetic organization – across distantly related species where

142    genome sequences are lacking.

# Results

## Capture and *de novo* assembly of the target regions

The capture probe design (workflow schematically shown in Figure 1) resulted in a total of 7057 probes based on the target region in Atlantic cod, covering 337 kbp of sequence. 26774 probes were designed for the additional codfishes, covering in total an area of 1.82 Mbp of target sequence. The target region and the *Hb* gene clusters were successfully captured and enriched for eight codfishes; Atlantic cod (*Gadus morhua*), haddock (*Melanogrammus aeglefinus)*, silvery pout (*Gadiculus argenteus*), cusk (*Brosme brosme*), burbot (*Lota lota*), European hake (*Merluccius merluccius*), marbled moray cod (*Muraenolepus marmoratus*), and roughhead grenadier (*Macrourus berglax*), with number of reads spanning from 35573 to 73005 (Table 1). The average read length was 3032 bp, varying from 2836 bp in European hake to 3265 bp in burbot, resulting in the capture of an average of 16.71 Mbp per species (Table 1). By mapping reads back to the capture target region we found that the average mapping depth was variable across the target region for all species (Figure 2 and 3). Because of the skewed distribution of mapping depth, we also calculated median depth, which was, as expected, the highest for Atlantic cod at 242x (Table S1). The median mapping depth was consistently high for most of the other species as well, with the lowest for roughhead grenadier (12x). Both median and average depths for the MN region are persistently higher than for the LA region for all species, with the exception of silvery pout (Table S1). Furthermore, positions with high degree of mapping corresponded to the location of the genes used in the design of the capture probes across all species (Figure 2 and Figure 3). The percentage of reads mapping to the target region ranged from 25-43%, however, the percentage of the target region covered by reads ranged from 53-100% with five species having more than 90% of the target region covered by reads (Figure 4c and Table S1).

To address factors influencing capture success we compared various capture statistics to overall genomic divergence between the Atlantic cod genome and independent WGS data for each species from (Malmstrøm *et al.* 2017) (Table S1). We found a strong negative correlation between genomic divergence to Atlantic cod and median mapping depth against the target region (r=-0.90, Figure 4a), percent of reads mapped to the target region (r=-0.90, Figure 4b), and percentage of reads mapped to the target region (r=-0.84, Figure 4c).

We constructed *de novo* assemblies with quite consistent assembly statistics across species. Contig N50 ranged from 8055 bp in burbot to 6523 bp in European hake and the total number of contigs varied from 205 in burbot to 455 in marbled

6

179    moray cod. However, there was some variation in the size of the largest contig,

180    which ranged from 79 kbp in Atlantic cod to 30 kbp in marbled moray cod (Table 1).

181    To evaluate whether the assemblies represent the actual target regions we mapped

182    the *de novo* assemblies for each species to the target region in Atlantic cod, for which

183    the capture design is largely based upon (Figure 2 and 3). As expected, the

184    assemblies corresponded to the regions with high coverage of reads, i.e. the areas of

185    the target region containing genes included in the probe design.

## Synteny of the *Hb* gene regions

187    Our capture design combined with long-read PacBio sequencing allowed us to

188    reconstruct micro-synteny of the MN and LA regions for Atlantic cod, haddock,

189    silvery pout, cusk, burbot, European hake, marbled moray cod and roughhead

190    grenadier (Figure 5). From the *de novo* assemblies, we were able to identify the

191    majority of the *Hb* genes and all of the flanking genes, which show that our capture

192    design was successful. However, the degree of continuity varied in the different

193    assemblies. In Atlantic cod, haddock, silvery pout, cusk, burbot and European hake

194    we could infer micro-synteny revealing that *Hb* and their flanking genes organization

195    largely followed what has previously been reported for Atlantic cod (Figure 5) (Star

196    *et al.* 2011). We found *Hbb4* only to be present in Atlantic cod (Figure 5b), which is in

197    line with (Baalsrud *et al.* 2017). Furthermore, the *de novo* assemblies confirmed a

198    linage-specific duplication of *Hbb2* in the roughhead grenadier (Baalsrud *et al.* 2017).

199    Additionally, we identified a complete *Hba4*-like gene in the assembly of the marbled

200    moray cod, not earlier identified in this species. However, the *Hba4*-like gene in

201    marbled moray cod is likely a pseudogene due to a frameshift mutation causing

202    multiple stop codons. Furthermore, we were able to identify most of the *Hb* genes

203    reported in the recent study of (Baalsrud *et al.* 2017), however, a few are missing from

204    our dataset (Figure 5a and b). Pairwise sequence alignment of these paralogous *Hb*

205    genes from (Baalsrud *et al.* 2017) revealed sequence identities up to 98 % (Table S2).

## Target region in the haddock and Atlantic cod genome assemblies

207    As a proof of concept, we reconstructed synteny of the target region in the most

208    recent genome assemblies of Atlantic cod (gadMor2 (Tørresen *et al.* 2017b)) and

209    haddock (melAeg (Tørresen *et al.* 2017a)). In Atlantic cod, the MN region is located

210    on linkage group 2 (Figure 5a) and LA on linkage group 18 (Figure 5b), in haddock

211    MN is located on scaffold MeA_20160214_scaffold_771 (Figure 5a) and LA on

212    scaffold MeA_20160214_scaffold_1676 (Figure 5b). The overall synteny in Atlantic

213    cod was congruent with (Wetten *et al.* 2010) except for the relative direction of the

214    genes *foxj1a* and *rhbdf1*. Furthermore, the organization of *Hb*s and their flanking

215    genes in the genome assembly of haddock is conserved compared to Atlantic cod

216    with the exception of *Hbb4* in the MN region, which is absent in haddock (Figure 5).

## Repetitive sequences in the in the *Hb* gene regions

218    Quantifying the amount of repetitive sequences in the target region(s) was only

219    possible for Atlantic cod (gadMor2) and haddock (melAeg), for which high-quality

220    genome assemblies exist. The amount of repetitive sequences in the target region

221    differed between the MN cluster and the LA cluster in Atlantic cod. The MN region

222    (214 kb) contained a total of 10.7% repeated sequences, including 1.0% retro-

223    elements, 1.3% transposons, 5.8% simple repeats, and 2.6% of various low complexity

224    and unclassified repeated sequences (Table S3). In comparison, in the smaller LA

225    region (123 kb) the proportion of repeated sequences was twice as high (20.3%),

226    which comprised of 2.8% retro-elements, 2.4% transposons, 13.8% simple repeats,

227    and 1.3% of various low complexity and unclassified repeated sequences.

228    Furthermore, the orthologous target regions in haddock followed the same pattern.

229    The MN region contained 16.3 % repeated sequences, in contrast to 19.8 % found in

230    the LA region (Table S3).

## Insertions and deletions in the promoter region of *Hba1 – Hbb1*

232    The previously shown 73 bp indel in the bi-directional promoter region of *Hba1* and

233    *Hbb1* – discerning the cold-adapted migratory Northeast Artic cod (NEAC) from the

234    more temperate-adapted southern Norwegian coastal cod (NCC) (Star *et al.* 2011) –

235    was confirmed by the improved version of the NEAC assembly (gadMor2). The

236    continuity of our capture assemblies (Figure 5) enabled location of the orthologous

237    captured regions in haddock, silvery pout and cusk. In each of the species an indel of

238    variable length were identified (Figure 6). Compared to the long promoter variant –

239    found to be linked with the *Hbb1-2* in Atlantic cod  – the indel is shorter in the other

240    species by 11 bp in haddock, 22 bp in silvery pout and 56 bp in cusk (Figure 6).

241    Although the indels are varying in length, the conserved flanking sequences in the

242    alignment clearly show that they represent orthologous regions. Moreover, we found

243    the amino acid positions at 55 and 62 of the *Hbb1* gene to vary between species;

244    Haddock has Val55-Lys62, silvery pout has Met55-Gln62, while cusk has Met55-

245    Lys62 similarly to NEAC (Figure 6). Additionally, we investigated amino acid

246    positions 55 and 62 in the *Hbb1* gene across a number additional codfish species for

247    which we have available gene sequences from (Baalsrud *et al.* 2017), revealing these

248    sites to be variable across this lineage (Table S4). Ancestral reconstruction of *Hbb1*

249    demonstrated that the ancestral state in position 55 was Met in codfishes, and in

250 position 62 was Lys in all codfishes except *Bregmaceros cantori* (Supplementary

251 Figures S1 and S2).

252      ## Discussion

253      ### Capture of *Hb* gene clusters with 70 million years divergence time reveal
254      ### conserved synteny and lineage-specific *Hb* duplications

255      We here demonstrate a successful in-solution targeted sequence capture and
256      assembling of coding and noncoding sequences of the *Hb* clusters from codfish
257      species separated by up to 70 million years (My) of evolution. Two features make our
258      approach unique from earlier studies. First, the target regions consisted of both
259      coding and noncoding genomic sequences. Second, we designed capture of large
260      fragments – combined with development of custom-made probes – in order to utilize
261      the long-read PacBio sequencing platform. This is in contrast to current targeted
262      capture sequencing protocols that are based on short-read sequencing technologies
263      (George *et al.* 2011; Mascher *et al.* 2013).

264            The organization and orientation of the *Hb* flanking genes that we identified
265      were conserved across all species (Figure 5a and b). However, in concordance with
266      earlier studies of the *Hb* region, we found significant variation in copy numbers of
267      the *Hb* genes, with linage specific duplications and losses (Star *et al.* 2011; Opazo *et al.*
268      2013; Feng *et al.* 2014; Baalsrud *et al.* 2017). We only found *Hbb4* in Atlantic cod,
269      supporting earlier studies showing that *Hbb4* is the result of a recent duplication in
270      this species (Borza *et al.* 2009; Baalsrud *et al.* 2017). Interestingly, the presence of two
271      copies of *Hbb2* on the same contig in the roughhead grenadier *de novo* assembly
272      confirmed a lineage specific gene duplication of *Hbb2*, which was found in a recent
273      study of *Hb*s in codfishes (Baalsrud *et al.* 2017). Additionally, a copy of the *Hba4* was
274      found in the *de novo* assembly of the marbled moray cod not found in (Baalsrud *et al.*
275      2017). The presence of a frame-shifting mutation that is causing multiple stop codons
276      indicated that this *Hba4* gene is most likely a pseudogene. *Hba4* is also a pseudogene
277      in the closely related species *Mora moro*, *Trachyrincus scabrus*, *T. murrayi* and
278      *Melanonus zugmayeri* (Baalsrud *et al.* 2017). Although we identified most of the *Hb*
279      genes from (Baalsrud *et al.* 2017), a few were absent from this dataset (Figure 5a and
280      b), which we suspect may be due to collapse of paralogous *Hb* genes, as they may
281      have as high as 98% sequence identity (Table S2).

282      ### Length variation in the bi-directional *Hba1-Hbb1* promoter within the
283      ### codfishes

284      The discovery of a promoter of variable length between *Hba1* and *Hbb1* in different
285      species (Figure 6) was concordant with earlier findings of length variation in the
286      homologous region in different populations of Atlantic cod (Star *et al.* 2011). The

287  migratory NEAC population has been shown to harbor the 73 bp longer variant at a
288  higher frequency compared to coastal cod populations (see Figure 6 and (Star *et al.*
289  2011)). Interestingly, we found relatively long promoters with high sequence
290  similarity to the NEAC indel in haddock and silvery pout. In contrast, cusk
291  displayed a relatively short promoter, however, still 17 bp longer than in NCC
292  (Figure 6). Furthermore, we found the amino acid positions 55 and 62 in *Hbb1*,
293  known to be polymorphic in Atlantic cod, to be variable across all codfishes included
294  in this study (Figure 6). Investigations of the same positions in a number additional
295  codfishes for which we have available gene sequences (Baalsrud *et al.* 2017), revealed
296  that these positions are highly variable across this linage (Table S2). Notably, the
297  most likely ancestral state of codfish *Hbb1* is Met55Lys62 (Supplementary Figures S1
298  and S2). Cusk and the coastal/southern Atlantic cod thus both display the ancestral
299  state as well as a short promoter, although the cusk promoter was 17 bp longer
300  (Figure 6). Collectively, these results suggest two different scenarios for promoter
301  length evolution. Scenario 1: The short promoter represents the ancestral state of the
302  Gadidae-family (including cusk and Atlantic cod; see (Malmstrøm *et al.* 2016)) and
303  that silvery pout and some populations of Atlantic cod have evolved a longer
304  promoter. Scenario 2: The long promoter is the ancestral state with independent
305  deletions of variable lengths in cusk, silvery pout, haddock and costal/southern
306  Atlantic cod (*Hbb1-1*). To disentangle this, we would need to obtain promoter
307  sequences from additional gadiform species. Regardless, the short-long promoter
308  polymorphism has been maintained throughout speciation events based on the
309  presence of both variants in Atlantic cod. Moreover, in both scenarios, cusk and
310  Atlantic cod (*Hbb1-1*) have maintained the ancestral Met55Lys62, while silvery pout,
311  haddock and Atlantic cod (*Hbb1-2*) have acquired substitutions at these positions due
312  to similar selection pressures or genetic drift. In this regard, it could be mentioned
313  that the NEAC, haddock and silvery pout display migratory behavior (e.g. diurnally
314  feeding movements as well as seasonal spawning migrations) compared to the more
315  stationary cusk and coastal cod (Eschemeyer and Fricke 2017) which could mean that
316  they have a higher $O_2$ demand and are exposed to greater temperature variation,
317  which in turn has selected for a temperature-dependent long promoter. Furthermore,
318  given that promoter length and positions 55/62 at *Hbb1* are important genetic
319  components of temperature adaptation in Atlantic cod populations (Star *et al.* 2011),
320  they most likely play a role in temperature adaptation in the other codfishes.

321  **Assembly success affected by probe design and repeat content**

322  In some species, nearly the complete target region is assembled in large contigs
323  containing multiple genes including cusk, whereas in other species such as the more

324 distantly related roughhead grenadier, the cluster is more fragmented (Figure 5). In
325 all species, the areas of the target regions that harbor genes of which probes are
326 designed for, as well as any areas containing repeated sequences, have very high
327 depths in comparison to the areas of intergenic sequences (Figure 2 and 3). This
328 poses a challenge for the assembly software, which is based in the assumption of
329 uniform depth over the sequencing data (Miller *et al.* 2010).
330       Overall, the MN cluster seems to be more successfully assembled than the LA
331 cluster, which is more fragmented (Figure 5). Differences in assembly completeness
332 between the two regions might be a result of several factors. Firstly, the MN region
333 has more flanking genes in closer proximity to the *Hb* region, which results in a
334 higher density of probes. Secondly, the overall repeat content of the LA region is one
335 order of magnitude larger than in the MN region, largely due to the larger
336 proportion of simple repeats. Repeat content is a major interference in capture
337 experiments because unwanted repetitive DNA may be enriched for, especially if
338 there are repeated sequences included in the probe design. Furthermore, if probes
339 were not completely covered by target DNA they get single-stranded sticky ends that
340 can hybridize to repetitive or other non-target DNA (Newman and Austin 2016).
341 Lastly, unless there were some longer reads that bridged such areas, this would in
342 turn have led to gaps in the downstream *de novo* assemblies. Following that the
343 assembly success was possibly a result of read length, we reason that a future
344 increase of the average read length from 3 kbp to 5-10 kbp, would be sufficient to
345 substantially increase the completeness of the assemblies. Due to the current circular
346 consensus (CCS) PacBio sequencing technology, however, which is a trade-off
347 between accuracy and length of reads, longer reads with sufficient accuracy are not
348 feasible.

### Long-read sequencing capture across species harbors new potential for comparative genomic studies

351 The number of reads mapping to the target region was in the range of 23-43%, which
352 may seem low compared with other capture studies. For instance, a whole exome
353 capture study on humans reported 56.1% of reads mapped to the target region (Guo
354 *et al.* 2012) and a similar study in rats reported to have 78.3% of reads on target
355 (Yoshihara *et al.* 2016). In contrast to our study however, these capture experiments
356 enriched either the exome or ultra-conserved elements within a single species.
357 Furthermore, we were able to cover up to 98% of the target region with >10 reads
358 across species (Table S1) which is similar to what mentioned experiments within
359 human and rat exomes reported (Guo *et al.* 2012; Yoshihara *et al.* 2016) and the main
360 difference is the higher percentage of non-target sequences in our study.

361    We were able to capture complete genes for species with 70 My divergence
362    time from the Atlantic cod (Figure 5). As expected, we found that capture success
363    declines with increased sequence divergence between the reference genome of which
364    we chiefly based our capture probes and the genomes of the included codfishes
365    (Figure 4). It has been reported that orthologous exons were successfully captured in
366    highly divergent frog species (with 200 My of separation), nevertheless the capture
367    success greatly decreased with increased evolutionary distance (Hedtke *et al.* 2013).
368    Similarly, it has been demonstrated that it is possible to capture >97% of orthologous
369    sequences in four species of primates that diverged from humans 40 My ago, using
370    probes entirely based on the human exome (George *et al.* 2011). Further, exomes were
371    effectively captured from skink species that diverged up to 80 My from the reference,
372    yet reporting a substantial decline in capture efficiency for sequences >10 % different
373    from the reference species (Bragg *et al.* 2016). Our study stands out from previous
374    capture experiments because intergenic, noncoding sequences in addition to genes
375    were captured. Efficient capture of intergenic sequences requires less divergence
376    time, as these regions usually evolve faster than genes (Koonin and Wolf 2010). Thus,
377    the most distantly related species from Atlantic cod for which we captured both
378    coding and noncoding sequences was burbot, which diverged from Atlantic cod 46
379    My (Figure 5). We argue, in line with (Schott *et al.* 2017), that sequence divergence
380    may be a more exact predictor of capture success than evolutionary distance, as the
381    sequence capture process is mainly influenced by the difference between the probe
382    sequence and the target sequence. European hake, marbled moray cod and
383    roughhead grenadier all diverged from cod about 70 My ago, however, the European
384    hake *Hb* regions was more successfully captured and assembled (Table 1; Figure 2).
385    This could be due to European hake having a lower genome-wide divergence to
386    Atlantic cod than marbled moray cod and roughhead grenadier (809k vs 879k and
387    907k SNPs; Table S1).
388    Finally, it should be mentioned that cusk – which diverged from Atlantic cod
389    39 My ago – was added to the experimental design after the species-specific probes
390    were generated. Thus, the successful capture of cusk was therefore solely based on
391    cross-species target enrichment, and could most likely been further improved if
392    species-specific probes for this species have been included.

393    **Concluding remarks and future perspectives**
394    Here, we have successfully demonstrated that combining targeted sequence capture
395    with long-read sequencing technology is as an efficient approach to obtain high
396    quality sequence data of a specific genomic region, including both coding and

397    noncoding sequences, across evolutionary distant species. We show that genome-
398    wide divergence is of importance for capture success across species. Furthermore, the
399    use of long-read sequencing augmented the *de novo* assembly of regions containing
400    repeated sequences that would otherwise fragment assemblies based on short-read
401    sequencing. This is crucial for capturing complete intergenic sequences that may be
402    highly divergent compared to genic regions even among fairly closely related
403    species. Given the rapid development in sequencing technologies future methods
404    will enable read-through of repeated regions and thus further increase the
405    completeness of assemblies. Moreover, a less stringent hybridization protocol should
406    make it possible to capture sequences across even deeper evolutionary time. In sum,
407    our approach has the potential of enhancing comparative genomic studies of
408    continuous genic and intergenic regions between any eukaryotic species-group
409    where genomic resources are scarce.

# Material and methods

## Defining target region and probe design

The probe design was chiefly based on the high-quality genome of Atlantic cod, known as gadMor2 (Tørresen *et al.* 2017b). In addition, species-specific probes were designed based on low-coverage assembled genomes (Malmstrøm *et al.* 2016) for ten selected species representing six families in the Gadiformes order. These species were Atlantic cod (*Gadus morhua*), Alaskan Pollock (*Gadus chalcogrammus*), polar cod (*Boreogadus saida*), haddock (*Melanogrammus aeglefinus*), Silvery pout (*Gadiculus argenteus*), burbot (*Lota lota*), European hake (*Merluccius merluccius*), roughhead grenadier (*Macrourus berglax*), roughsnout grenadier (*Trachyrincus scabrus*) and marbled moray cod (*Muraenolepus marmoratus*).

To retrieve relevant sequence data for the probe design, the MN and LA *Hb* regions were extracted from gadMor2 (Figure 1). These sequences, hereby known as the target region, were then used queries in BLAST (Altschul *et al.* 1990) searches with an E-value threshold of <0.1 against the genome assembly data of all ten species.

In total, 5604 sequences from the chosen species were supplied to NimbleGen probe design. Protein coding genes from the ENSEMBL database were used to define the regions to be tiled in the probe design (Table S5) within the target region of the Atlantic cod, and the unitigs for each of the additional codfishes.

NimbleGen SeqCap EZ capture probes were designed by NimbleGen (Roche, Madison, USA) using a proprietary design algorithm. NimbleGen offers an in-solution sequence capture protocol, which includes custom made probes. Uniquely, the capture probes from NimbleGen are tiled to overlap the target area. 50 – 100 bp (average 75 bp) probes where designed tiled over the target region (subset of gadMor2) resulting in each base, on average, being covered by two probes. Additionally, raw reads from Illumina sequencing from (Malmstrøm *et al.* 2017) were used for each species to estimate repetitive sequences in each of the species' genomes, aiming to discard probes containing any repeats.

## Sample collection and DNA extraction

Our goal working with animals is always to limit any harmful effects of our research on populations and individuals. Whenever possible we try to avoid animals being euthanized to serve our scientific purpose alone by collaborating with commercial fisheries or museums. The tissue samples used in this study are either from commercially fished individuals intended for human consumption or museum specimen. The commercially caught fish were immediately stunned, by bleeding

15

445   following standard procedures by a local fisherman. There is no specific legislation

446   applicable to this manner of sampling in Norway, however it is in accordance with

447   the guidelines set by the 'Norwegian consensus platform for replacement, reduction

448   and refinement of animal experiments' (www.norecopa.no).

449       DNA was extracted from tissue samples using High Salt DNA Extraction

450   method by Phill Watts (https://www.liverpool.ac.uk/~kempsj/IsolationofDNA.pdf ,

451   last day accessed: December 2017). The concentration and purity of the DNA samples

452   were quantified using NanoDrop (Thermo Scientific, Thermo Fisher Scientific,

453   Waltham, MA, USA) and a Qubit fluorometer (Invitrogen, Thermo Fisher Scientific,

454   Waltham, MA, USA). Due to poor DNA quality, three species included in the probe

455   design; Alaskan Pollock, polar cod and roughsnout grenadier were excluded from

456   further analysis. In total, eight species were sequenced; seven of these species were

457   included in the probe design and one closely related species (cusk, *Brosme brosme*),

458   which serves as a cross species capture experiment without species-specific probes.

### Capture, library preparation and sequencing

460   The sequencing libraries were prepared following a modified Pacific Biosciences

461   SeqCap EZ protocol. As multiplexing of the samples before capture was required,

462   barcodes were designed at the Norwegian Sequencing Centre

463   (http://www.sequencing.uio.no) using guidelines from Pacific Biosciences

464   (Supplementary Materials and methods). Genomic DNA was sheared to 5 kb

465   fragments using MegaRuptor (Diagenode, Seraing (Ougrée), Belgium). Due to poorer

466   DNA quality, fragmenting was not done for European hake. For this sample together

467   with fragmented DNA from roughhead grenadier, short fragments were removed

468   using BluePippin (Sage Science, Beverly, MA, USA) before library preparation.

469   Illumina libraries were prepared using KAPA Hyper Prep kit (Kapa Biosystems,

470   Wilmington, MA, USA) and barcoded using different Illumina barcodes. PacBio

471   barcodes were implemented during pre-capture amplification of libraries. After

472   amplification, fragment length distribution was evaluated using Bioanalyzer (Agilent

473   Technologies, Santa Clara, CA, USA) and samples were pooled in equimolar ratio.

474   During hybridization, SeqCap EZ Developer Reagent (universal repeat blocker for

475   use on vertebrate genomes) and oligos corresponding to Illumina and PacBio

476   barcodes were used for blocking. Captured gDNA was amplified to ensure that

477   sufficient amount of DNA was available for PacBio library preparation. Size selection

478   of the libraries was performed using BluePippin. Final libraries were quality checked

479   using Bioanalyzer and Qubit fluorometer (Invitrogen, Thermo Fisher Scientific,

480   Waltham, MA, USA) and sequenced on RS II instrument (PacBio, Menlo Park, CA,

481    USA) using P6-C4 chemistry with 360 minutes movie time. In total, 9 SMRT cells

482    were used for sequencing.

### *De novo* assemblies

484    Reads were filtered and de-multiplexed using the 'RS_reads of insert.1' pipeline on

485    SMRT Portal (SMRT Analysis version smrtanalysis_2.3.0.140936.p2.144836). Each set

486    of reads corresponding to a given species was crossed-checked with their respective

487    six-nucleotide Illumina adapter. Reads containing an incorrect Illumina adapter were

488    removed. Adapter sequences were then trimmed using the application Prinseq-lite

489    v0.20.4 (Schmieder and Edwards 2011). The trimmed reads were assembled *de novo*

490    using Canu v1.4 +155 changes (r8150 c0a988b6a106c27c6f993dfe586d2336282336a6)

491    (Berlin *et al.* 2015). The Canu software is optimized for assembling single molecule

492    high noise sequence data. We specified genome size as the size of the target region

493    (300 kbp). Additionally, we ran PBJelly (English *et al.* 2012) on the Canu *de novo*

494    assemblies, using the raw reads to possible bridge gaps between scaffolds, settings

495    given in Supplementary Materials and Methods.

496          We assessed the assemblies by running Assemblathon 2 (Bradnam *et al.* 2013),

497    which reports assembly metrics such as the longest contig, the number of contigs,

498    and the N50 value. *De novo* assemblies of the MN and LA regions of Atlantic cod and

499    haddock were aligned and compared to their reference genomes, gadMor2 and

500    melAeg respectively, using BLAST and BWA v0.7.10 (Li and Durbin 2009) to

501    determine syntenic similarities and assembly completeness.

### Estimating capture success

503    PacBio reads for all the species were mapped back to the Atlantic cod genome

504    assembly (gadMor2) in order to determine sequence capture success and target

505    mapping depths. Mapping was done using BWA-MEM v0.7.10 (Li and Durbin 2009).

506    Target-area read depth for all the species based on mapping against gadMor2, were

507    calculated using Samtools v1.3.1 (Li *et al.* 2009). We calculated both average and

508    median mapping depth against the target region as a whole and for the MN and LA

509    region separately. We also calculated percentage of reads that mapped to the target

510    region, and the percentage of the target regions covered by reads to a minimum

511    depth of 10x. To compare assemblies to the target region we additionally mapped the

512    assemblies to the target region. In order to verify the sequence capture process,

513    sequence data for Atlantic cod and haddock were mapped back to their reference

514    genomes using BWA-mem v0.7.10 (Li and Durbin 2009). The results were visualized

515    using Integrative Genome Viewer (Robinson *et al.* 2011).

516    To obtain an independent measure of divergence between species in the
517    capture experiment we calculated genome wide level of divergence of each species to
518    the reference genome of Atlantic cod using low-coverage whole-genome sequence
519    data from (Malmstrøm *et al.* 2017). We mapped raw reads to Atlantic cod using
520    BWA-MEM (Li and Durbin 2009) and called SNPs using the Freebayes variant caller
521    (Garrison and Marth 2017). Some species are more closely related to Atlantic cod
522    than others, which could introduce a bias in mapping. To avoid this, we only looked
523    at genomic regions where all species mapped. The number of SNPs was then used as
524    an estimate of genome-wide divergence of each species to Atlantic cod. We also
525    mapped a low-coverage genome of Atlantic cod to the Atlantic cod reference genome
526    as a control.

527    In pursuance of factors explaining capture success we tested for correlations
528    and plotted the relationship between the genome wide level of divergence and the
529    following variables; median mapping depth against the target region (for total, LA
530    and MN, respectively); percentage of reads that mapped to the target region; and the
531    percentage of the target region covered by reads. All tests and plots were done using
532    R version 3.2.5(Team 2013).

533    Assembly continuity is very often hampered by the presence of repeats, which
534    create gaps. We therefore quantified repeat-content in the target region extracted
535    from gadMor2 and orthologous regions in haddock using Repeatmasker Open 3.0
536    (Smit *et al.* 2010) for the MN region and the LA region separately.

## Identifying gene location and synteny

538    In order to identify the genes of interest and their location in the assembly we used
539    local sequence alignment algorithm BLAST v2.4.0 (Altschul *et al.* 1990) with protein
540    sequences of the genes of interest (Table S5) as queries. tblastn was used with an e-
541    value of 0.1. Investigation of *Hbb1-Hba1* promoter region was done for four species,
542    Atlantic cod, haddock, silvery cod and cusk. Sequences were aligned with ClustalW
543    default settings using MEGA7 (Kumar *et al.* 2016). Ancestral sequence reconstruction
544    was carried out for *Hbb-1* gene sequences from 24 species of codfishes from (Baalsrud
545    *et al.* 2017) using a maximum likelihood method implemented in MEGA7 (Kumar *et*
546    *al.* 2016).

547    Additionally, we estimated sequence identity using EMBOSS Needle (Rice *et*
548    *al.* 2000) with default settings, between *Hbb* gene sequences from (Baalsrud *et al.*
549    2017) that where missing and present in the *de novo* assemblies to evaluate similarity
550    (Table S2).

## Acknowledgements

## Author contributions

H.T.B. and S.J. initially conceived and designed the study, with input from S.N.K.H, A.T.-K., M.S., G.O., R.S., and K.S.J. Tissue samples were provided by S.J. and H.T.B. Probe design was carried out by T.R. with assistance from S.N.K.H and H.T.B. DNA extraction and sequence library preparation was performed by S.N.K.H and A.T.-K, respectively. Sequence capture was carried out by S.N.K.H, A.T.-K., M.S. and G.O. Filtering, mapping of sequences and *de novo* assemblies was done by S.N.K.H., assisted by O.K.T and H.T.B. Annotation of genes, synteny analyses, statistical analyses and construction of all figures and tables was done by S.N.K.H and H.T.B. The manuscript was written by S.N.K.H and H.T.B. with input from S.J. and K.S.J.

## Competing interests

The authors declare that they have no competing interests.

## Data and materials availability

All reads and assemblies (unitigs) reported on here, and the target region, subset of gadMor2 have been deposited at figshare under doi/xxx.

## References

Altschul, S. F., W. Gish, W. Miller, E. W. Myers, and D. J. Lipman, 1990 Basic local alignment search tool. J. Mol. Biol. 215: 403–410.

Andersen, Ø., O. F. Wetten, M. C. De Rosa, C. Andre, C. Carelli Alinovi *et al.*, 2009 Haemoglobin polymorphisms affect the oxygen-binding properties in Atlantic cod populations. Proc. Biol. Sci. 276: 833–841.

Baalsrud, H. T., K. L. Voje, O. K. Tørresen, M. H. Solbakken, M. Matschiner *et al.*, 2017 Evolution of Hemoglobin Genes in Codfishes Influenced by Ocean Depth. Sci Rep 7: 168–10.

Barlow, S. L., J. Metcalfe, D. A. Righton, and M. Berenbrink, 2017 Life on the edge: O 2binding in Atlantic cod red blood cells near their southern distribution limit is not sensitive to temperature or haemoglobin genotype. Journal of Experimental Biology 220: 414–424.

Berlin, K., S. Koren, C.-S. Chin, J. P. Drake, J. M. Landolin *et al.*, 2015 Assembling large genomes with single-molecule sequencing and locality-sensitive hashing. Nature Biotechnology 33: 623–630.

Bickhart, D. M., B. D. Rosen, S. Koren, B. L. Sayre, A. R. Hastie *et al.*, 2017 Single-molecule sequencing and chromatin conformation capture enable de novo reference assembly of the domestic goat genome. Nat Genet 431: 931–11.

Borza, T., C. Stone, A. K. Gamperl, and S. Bowman, 2009 Atlantic cod (*Gadus morhua*) hemoglobin genes: multiplicity and polymorphism. BMC Genetics 10: 51.

Bradnam, K. R., J. N. Fass, A. Alexandrov, P. Baranay, M. Bechner *et al.*, 2013 Assemblathon 2: evaluating de novo methods of genome assembly in three vertebrate species. Giga Sci 2: 10.

Bragg, J. G., S. Potter, K. Bi, and C. Moritz, 2016 Exon capture phylogenomics: efficacy across scales of divergence. Molecular Ecology Resources 16: 1059–1068.

Broeckx, B. J. G., F. Coopman, G. E. C. Verhoeven, V. Bavegems, S. De Keulenaer *et al.*, 2014 Development and performance of a targeted whole exome sequencing enrichment kit for the dog (Canis Familiaris Build 3.1). Sci Rep 4: 1522–4.

Ellegren, H., 2014 Genome sequencing and population genomics in non-model organisms. Trends in Ecology & Evolution 29: 51–63.

English, A. C., S. Richards, Y. Han, M. Wang, V. Vee *et al.*, 2012 Mind the Gap: Upgrading Genomes with Pacific Biosciences RS Long-Read Sequencing

611 Technology (Z. Liu, Ed.). PLoS ONE 7: e47768.

612 Eschemeyer, W. N., and R. Fricke, 2017 Catalog of fishes.
613  http://research.calacademy.org/researchichthyology/catalog/fishcatmain.asp.

614 Feng, J., S. Liu, X. Wang, R. Wang, J. Zhang *et al.*, 2014 Channel catfish hemoglobin
615  genes: Identification, phylogenetic and syntenic analysis, and specific induction
616  in response to heat stress. Comparative Biochemistry and Physiology Part D:
617  Genomics and Proteomics 9: 11–22.

618 Garrison, E., and G. Marth, 2017 Haplotype-based variant detection from short-read
619  sequencing. 1–9.

620 George, R. D., G. McVicker, R. Diederich, S. B. Ng, A. P. MacKenzie *et al.*, 2011 Trans
621  genomic capture and sequencing of primate exomes reveals new targets of
622  positive selection. Genome Research 21: 1686–1694.

623 Goodwin, S., J. D. McPherson, and W. R. McCombie, 2016 Coming of age: ten years
624  of next-generation sequencing technologies. Nat Rev Genet 17: 333–351.

625 Grover, C. E., A. Salmon, and J. F. Wendel, 2012 Targeted sequence capture as a
626  powerful tool for evolutionary analysis. American Journal of Botany 99: 312–319.

627 Guo, Y., J. Long, J. He, C.-I. Li, Q. Cai *et al.*, 2012 Exome sequencing generates high
628  quality data in non-target regions. BMC Genomics 13:.

629 Hedtke, S. M., M. J. Morgan, D. C. Cannatella, and D. M. Hillis, 2013 Targeted
630  Enrichment: Maximizing Orthologous Gene Comparisons across Deep
631  Evolutionary Time (U. Joger, Ed.). PLoS ONE 8: e67908.

632 Huddleston, J., S. Ranade, M. Malig, F. Antonacci, M. Chaisson *et al.*, 2014
633  Reconstructing complex regions of genomes using long-read sequencing
634  technology. Genome Research 24: 688–696.

635 Jones, M. R., and J. M. Good, 2015 Targeted capture in evolutionary and ecological
636  genomics. Molecular Ecology 25: 185–202.

637 Karpov, A. K., and G. G. Novikov, 1980 Hemoglobin alloforms in cod, *Gadhus morhua*
638  (Gadiformes, Gadidae), their functional characteristics and occurrence in
639  populations. Journal of Ichthyology 20: 45–50.

640 Kim, K. E., P. Peluso, P. Babayan, P. J. Yeadon, C. Yu *et al.*, 2014 Long-read, whole-
641  genome shotgun sequence data for five model organisms. Sci. Data 1: 140045–10.

642 Koonin, E. V., and Y. I. Wolf, 2010 Constraints and plasticity in genome and
643  molecular-phenome evolution. Nat Rev Genet 11: 487–498.

644   Korlach, J., G. Gedman, S. B. Kingan, C.-S. Chin, J. T. Howard *et al.*, 2017 Giga Sci 6:
645        1–16.

646   Kumar, S., G. Stecher, and K. Tamura, 2016 MEGA7: Molecular Evolutionary
647        Genetics Analysis Version 7.0 for Bigger Datasets. Molecular Biology and
648        Evolution 33: 1870–1874.

649   Li, C., S. Corrigan, L. Yang, N. Straube, M. Harris *et al.*, 2015 DNA capture reveals
650        transoceanic gene flow in endangered river sharks. Proc. Natl. Acad. Sci. U.S.A.
651        112: 13302–13307.

652   Li, H., and R. Durbin, 2009 Fast and accurate short read alignment with Burrows-
653        Wheeler transform. Bioinformatics 25: 1754–1760.

654   Li, H., B. Handsaker, A. Wysoker, T. Fennell, J. Ruan *et al.*, 2009 The Sequence
655        Alignment/Map format and SAMtools. Bioinformatics 25: 2078–2079.

656   Malmstrøm, M., M. Matschiner, O. K. Tørresen, K. S. Jakobsen, and S. Jentoft, 2017
657        Whole genome sequencing data and de novo draft assemblies for 66 teleost
658        species. Sci. Data 4: 1–13.

659   Malmstrøm, M., M. Matschiner, O. K. Tørresen, B. Star, L. G. Snipen *et al.*, 2016
660        Evolution of the immune system influences speciation rates in teleost fishes. Nat
661        Genet 48: 1204–1210.

662   Mascher, M., G. J. Muehlbauer, D. S. Rokhsar, J. Chapman, J. Schmutz *et al.*, 2013
663        Anchoring and ordering NGS contig assemblies by population sequencing
664        (POPSEQ). Plant J 76: 718–727.

665   Miller, J. R., S. Koren, and G. Sutton, 2010 Assembly algorithms for next-generation
666        sequencing data. Genomics 95: 315–327.

667   Morin, A., T. Kwan, B. Ge, L. Letourneau, M. Ban *et al.*, 2016 Immunoseq: the
668        identification of functionally relevant variants through targeted capture and
669        sequencing of active regulatory regions in human immune cells. BMC Med
670        Genomics 9: 7–12.

671   Newman, C. E., and C. C. Austin, 2016 Sequence capture and next-generation
672        sequencing of ultraconserved elements in a large-genome salamander. Molecular
673        Ecology 25: 6162–6174.

674   Ng, S. B., E. H. Turner, P. D. Robertson, S. D. Flygare, A. W. Bigham *et al.*, 2009
675        Targeted capture and massively parallel sequencing of 12 human exomes. Nature
676        461: 272–276.

677   Opazo, J. C., G. T. Butts, M. F. Nery, J. F. Storz, and F. G. Hoffmann, 2013 Whole-

678    genome duplication and the functional diversification of teleost fish
679    hemoglobins. Molecular Biology and Evolution 30: 140–153.

680    Patrushev, L. I., and T. F. Kovalenko, 2014 Functions of noncoding sequences in
681    mammalian genomes. Biochemistry (Mosc) 79: 1442–1469.

682    Quinn, N. L., K. A. Boroevich, K. P. Lubieniecki, W. Chow, E. A. Davidson *et al.*, 2010
683    Genomic organization and evolution of the Atlantic salmon hemoglobin
684    repertoire. BMC Genomics 11: 539.

685    Rice, P., I. Longden, and A. Bleasby, 2000 EMBOSS: the European Molecular Biology
686    Open Software Suite. Trends in Genetics 16: 276–277.

687    Robinson, J. T., H. Thorvaldsdottir, W. Winckler, M. Guttman, E. S. Lander *et al.*, 2011
688    Integrative genomics viewer. Nature Biotechnology 29: 24–26.

689    Samorodnitsky, E., J. Datta, B. M. Jewell, R. Hagopian, J. Miya *et al.*, 2015 Comparison
690    of custom capture for targeted next-generation DNA sequencing. J Mol Diagn 17:
691    64–75.

692    Samuels, D. C., L. Han, J. Li, S. Quanghu, T. A. Clark *et al.*, 2013 Finding the lost
693    treasures in exome sequencing data. Trends in Genetics 29: 593–599.

694    Schmieder, R., and R. Edwards, 2011 Quality control and preprocessing of
695    metagenomic datasets. Bioinformatics 27: 863–864.

696    Schott, R. K., B. Panesar, D. C. Card, M. Preston, T. A. Castoe *et al.*, 2017 Targeted
697    Capture of Complete Coding Regions across Divergent Species. Genome Biology
698    and Evolution 9: 398–414.

699    Smit, A., R. Hubley, and P. Green, 2010 RepeatMasker Open 3.0.
700    http://www.repeatmasker.org.

701    Star, B., A. J. Nederbragt, S. Jentoft, U. Grimholt, M. Malmstrøm *et al.*, 2011 The
702    genome sequence of Atlantic cod reveals a unique immune system. Nature 477:
703    207–210.

704    Syring, J. V., J. A. Tennessen, T. N. Jennings, J. Wegrzyn, C. Scelfo-Dalbey *et al.*, 2016
705    Targeted Capture Sequencing in Whitebark Pine Reveals Range-Wide
706    Demographic and Adaptive Patterns Despite Challenges of a Large, Repetitive
707    Genome. Front Plant Sci 7: 484.

708    Team, R. C., 2013 R Core Team. R: A language and environment for statistical
709    computing. R Foundation for Statistical Computing, Vienna, Austria: ISBN 3-
710    900051-07-0, URL http://www. R-project. org/.(3.3. 1) Software Vienna, Austria: R
711    Foundation for Statistical Computing.

712    Teer, J. K., L. L. Bonnycastle, P. S. Chines, N. F. Hansen, N. Aoyama *et al.*, 2010
713       Systematic comparison of three genomic enrichment methods for massively
714       parallel DNA sequencing. Genome Research 20: 1420–1431.

715    Turner, E. H., S. B. Ng, D. A. Nickerson, and J. Shendure, 2009 Methods for Genomic
716       Partitioning. Annu. Rev. Genom. Human Genet. 10: 263–284.

717    Tørresen, O. K., M. S. O. Brieuc, M. H. Solbakken, E. Sørhus, A. J. Nederbragt *et al.*,
718       2017a Genomic architecture of codfishes featured by expansions of innate
719       immune genes and short tandem repeats. bioRxiv 1–42.

720    Tørresen, O. K., B. Star, S. Jentoft, W. B. Reinar, H. Grove *et al.*, 2017b An improved
721       genome assembly uncovers prolific tandem repeats in Atlantic cod. BMC
722       Genomics 18: 311–23.

723    Volff, J.-N., 2005 Genome evolution and biodiversity in teleost fish. Heredity 94: 280–
724       294.

725    Wells, R. M. G., 2005 Blood-Gas Transport and Hemoglobin Function in Polar Fishes:
726       Does Low Temperature Explain Physiological Characters?, pp. 281–316 in *Fish*
727       *Physiology*, Fish Physiology, Elsevier.

728    Wetten, O. F., A. J. Nederbragt, R. C. Wilson, K. S. Jakobsen, R. B. Edvardsen *et al.*,
729       2010 Genomic organization and gene expression of the multiple globins in
730       Atlantic cod: conservation of globin-flanking genes in chordates infers the origin
731       of the vertebrate globin clusters. BMC Evolutionary Biology 10: 315.

732    Woolfe, A., M. Goodson, D. K. Goode, P. Snell, G. K. McEwen *et al.*, 2004 Highly
733       Conserved Non-Coding Sequences Are Associated with Vertebrate Development
734       (Sean Eddy, Ed.). PLoS Biol 3: e7.

735    Yoshihara, M., D. Saito, T. Sato, O. Ohara, T. Kuramoto *et al.*, 2016 Design and
736       application of a target capture sequencing of exons and conserved non-coding
737       sequences for the rat. BMC Genomics 17: 1522–11.

# Figure legends

738

739  **Figure 1: Flowchart of sequence capture approach**. a) Sequence data from the
740  Atlantic cod genome (gadMor2 (Tørresen *et al.* 2017b)) combined with gene
741  sequences of target genes and sequences from low coverage genomes of the
742  additional codfishes are combined to generate probes. b) Isolated DNA is
743  multiplexed with Illumina and PacBio barcodes. c) Raw reads for each species are
744  used to score all probes, ensuring that no repeated sequences are present. DNA
745  Probes are used in solution on isolated DNA for all of the included species,
746  hybridizing to the target sequences. Target sequences are then captures and
747  sequences on the PacBio RSII sequencing platform. d) Downstream bioinformatics
748  includes de-multiplexing of reads and trimming, making the reads ready for
749  downstream analysis such as mapping and *de novo* assembly.
750

751  **Figure 2: Mapping of reads and assemblies against the MN target region.** Each
752  panel shows the reads and *de novo* assembly mapped against the MN target region in
753  grey and orange, respectively, for species a.) Atlantic cod, b) haddock, c) silvery
754  pout, d) cusk, e) burbot, f) European hake, g) marbled moray cod and h) roughhead
755  grenadier. The positions of genes in the target region are indicated at the top.
756

757  **Figure 3: Mapping of reads and assemblies against the LA target region.** Each
758  panel shows the reads and *de novo* assembly mapped against the LA target region in
759  grey and orange, respectively, for species a) Atlantic cod, b) haddock, c) silvery pout,
760  d) cusk, e) burbot, f) European hake, g) marbled moray cod and h) roughhead
761  grenadier. The positions of genes in the target region are indicated at the top.
762

763  **Figure 4: The relationship between capture success and genomic divergence to
764  Atlantic cod**. Linear regression of the relationship between the genomic divergence
765  to Atlantic cod (SNPs x $10^5$) and a) median mapping depth for the MN region (blue),
766  LA region (red) and the combined target region (black); b) the percentage of reads
767  mapping to the target region; c) the percentage of the target region covered by reads
768  to a minimum depth of 10x. For each regression the correlation coefficient, r, is
769  shown along with a p-value. Each data point is labeled by species according to this
770  code: Ac=Atlantic cod, H=haddock, Sp=silvery pout, C=cusk, B=burbot, Eh=European
771  hake, Mm=marbled moray cod and Rg=roughhead grenadier.
772

773  **Figure 5: Synteny of the Hb gene clusters**. Genomic synteny of the hemoglobin gene
774  clusters shown at the top for the genomes of Atlantic cod (gadMor2 (Tørresen *et al.*
775  2017b)) and haddock (MelAeg (Tørresen *et al.* 2017a)). Below, the genomic synteny
776  inferred from the *de novo* assemblies for all of the species included in the capture
777  experiment. Stippled lines indicate assembly gaps – here we assume that the
778  orientation of genes corresponds to the genomes of Atlantic cod and haddock. Gray
779  boxes indicate genes that have been identified in (Baalsrud *et al.* 2017), but are absent

25

780    in the *de novo* assemblies. a) Synteny across the MN region b) Synteny across the LA

781    region.

782

783    **Figure 6: Polymorphisms in the bi-directional promoter between *Hba1* and *Hbb1***

784    **for five species in the Gadidae family.**

785    A schematic representation of *Hba1* and *Hbb1* with the promoter region between

786    them. The region contains an indel polymorphism of variable length across the five

787    species, as indicated by gaps. For each species/variant the alignment is shown along

788    with amino acid substitutions at positions 55 and 62 in the translated part of the *Hbb1*

789    gene.

## Supporting Information Legends

**Table S1:** For each species, the average and median depth of reads mapped against the target region (for MN, LA and total), the genomic divergence to Atlantic cod (number of SNPs), percentage of nucleotides mapped to the target and the percentage of the target region with more than 10x coverage.

**Table S2:** Estimated sequence identity using EMBOSS Needle (Rice *et al.* 2000) with default settings, between paralogous Hbb gene sequences from (Baalsrud *et al.* 2017). Genes highlighted in bold are missing from the assemblies in figure 5.

**Table S3:** Amino acids at positions 55 and 62 in Hbb1 in various codfishes taken from (Baalsrud *et al.* 2017).

**Table S4:** Amount of repeated sequences in the target region of the Atlantic cod (gadMor2 (Tørresen *et al.* 2017b)) and haddock (melAeg (Tørresen *et al.* 2017a)) given in percentage.

**Table S5:** Genes provided Nimblegen for the probe design, and used to identify genes in *de novo* assemblies. For each gene, the gene name is given with its ENSEMBL name and ENSEMBL identifier.
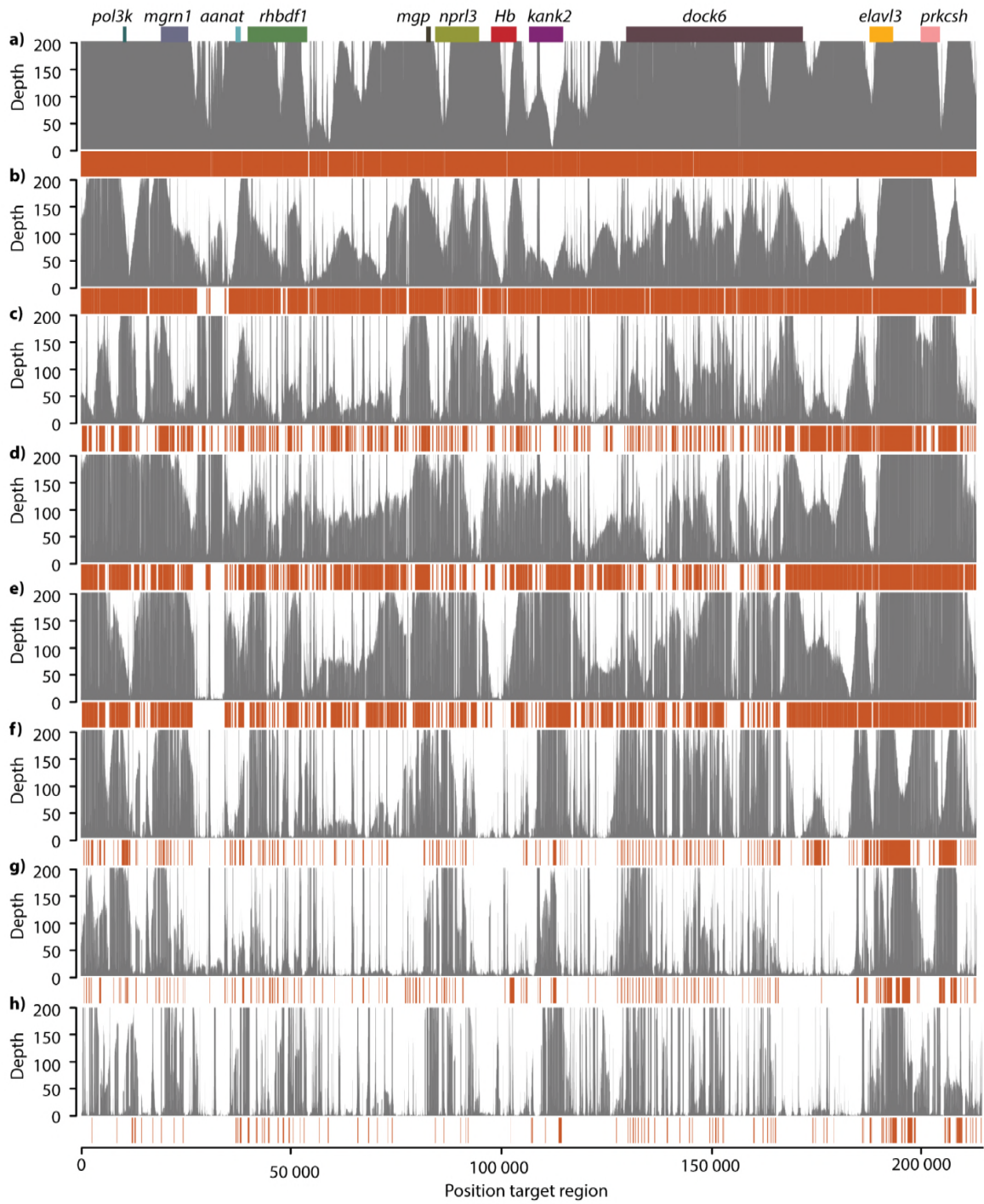
**Figure S1:** Ancestral reconstruction of amino acids at position 55 in the *Hbb-1* gene in Gadiformes. Phylogenetic trees and ancestral reconstruction was carried out in MEGA 7.0.

**Figure S2:** Ancestral reconstruction of amino acids at position 62 in the *Hbb-1* gene in Gadiformes. Phylogenetic trees and ancestral reconstruction was carried out in MEGA 7.0.
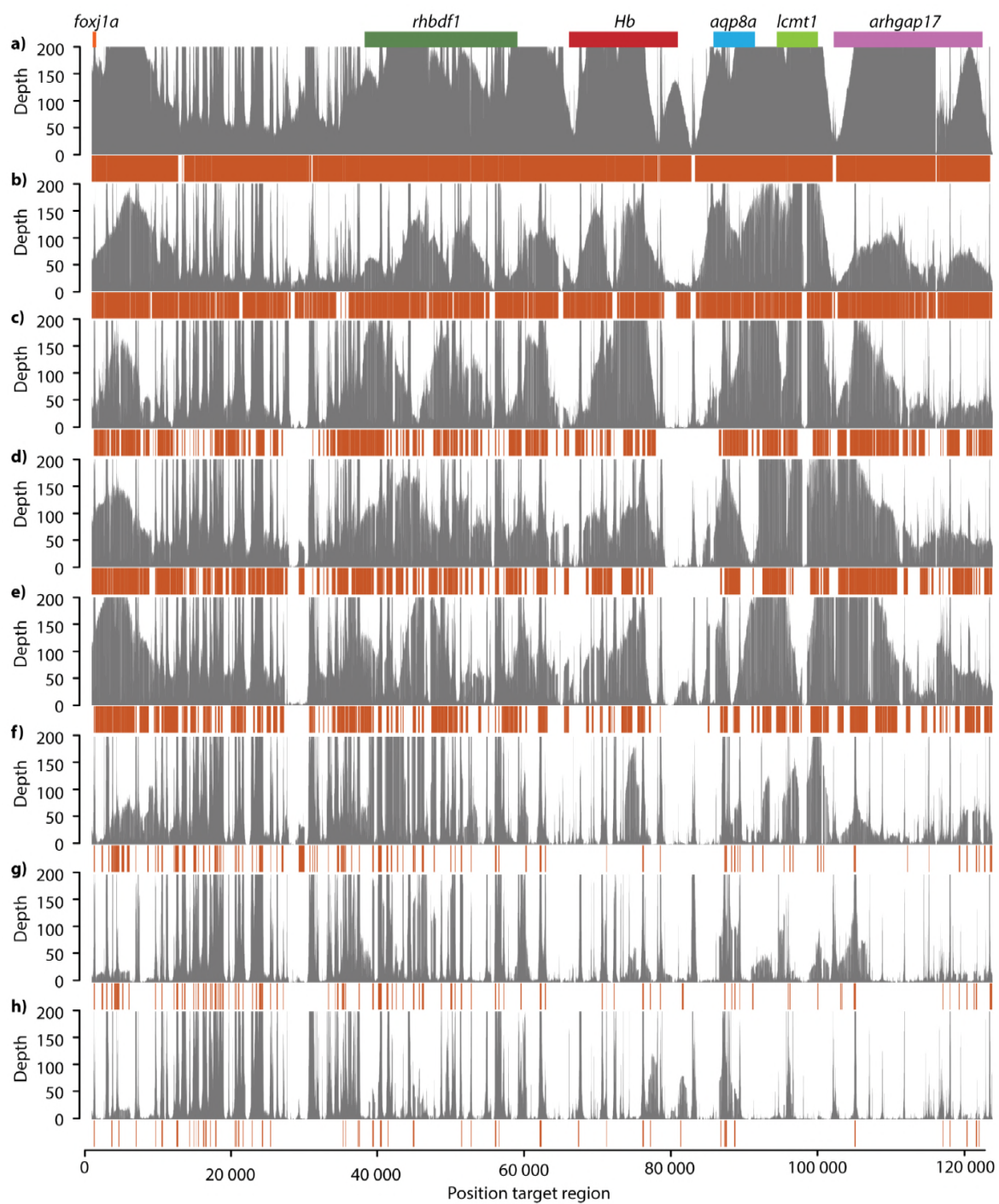
818 **Figures**



819
820     Figure 1
821
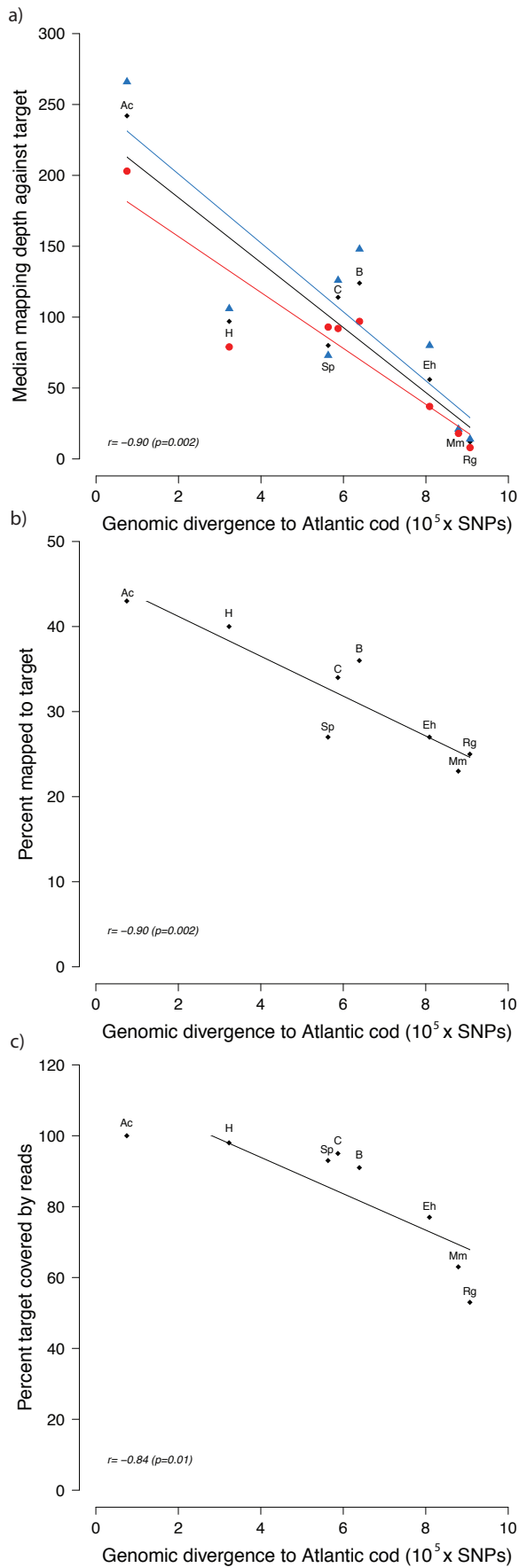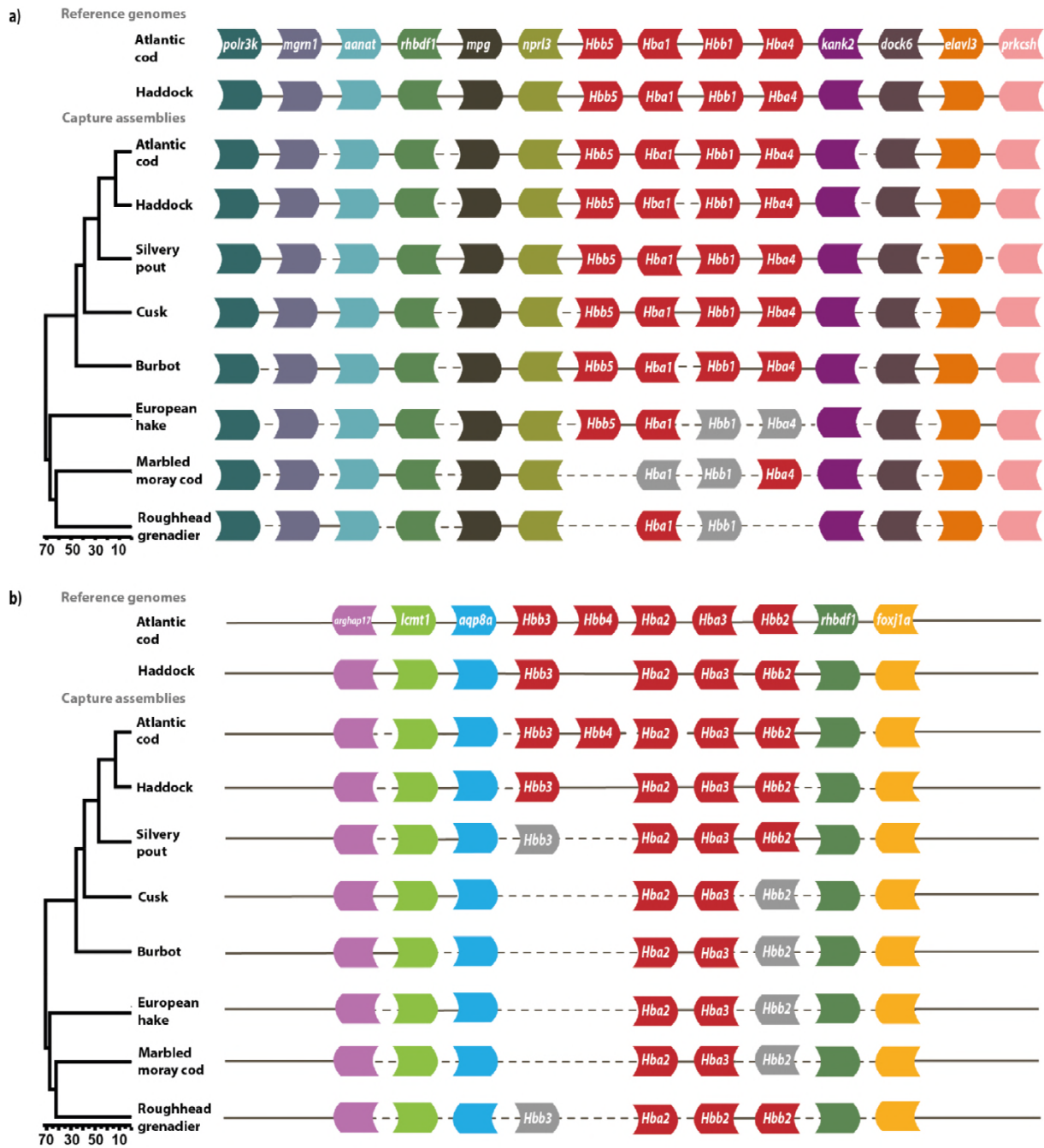
822
823    Figure 2
824

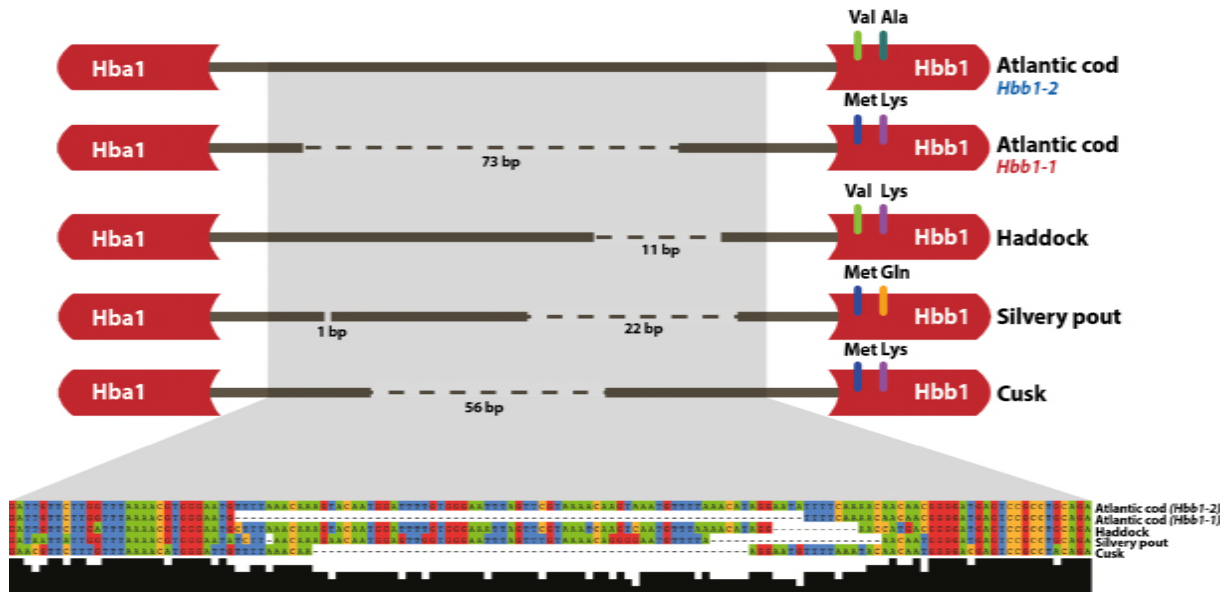Figure 3

828
829　　Figure 4

830
831    Figure 5
832

32

833
834    Figure 6
835

## Tables

**Table 1.** Number of reads and bases captured and sequenced for each species, and number of utgs, largest utg and N50 in the assemblies.

| Species | Latin name | Number of reads | Number of bases | Number of utgs | Largest utg (bp) | N50 (bp) |
|---|---|---|---|---|---|---|
| Atlantic cod | *Gadus morhua* | 73005 | 217252583 | 278 | 79 020 | 7 728 |
| Haddock | *Melanogrammus aeglefinus* | 35573 | 107839552 | 227 | 52 433 | 7 227 |
| Silvery pout | *Gadiculus argenteus* | 69775 | 212519845 | 410 | 35 801 | 7 098 |
| Cusk | *Brosme brosme* | 55348 | 175883008 | 394 | 64 145 | 7 322 |
| Burbot | *Lota lota* | 56155 | 165360828 | 205 | 70 602 | 8 055 |
| European hake | *Merluccius merluccius* | 65661 | 180558336 | 311 | 31 558 | 6 523 |
| Marbled moray cod | *Muraenolepus marmoratus* | 52076 | 148100933 | 455 | 30 019 | 6 632 |
| Roughhead grenadier | *Macrourus berglax* | 46195 | 129085001 | 325 | 35 216 | 7122 |