



**This document is a postprint version of an article published in Marine Genomics© Elsevier after peer review. To access the final edited and published work see <https://doi.org/10.1016/j.margen.2018.01.002>**

1 **Muscle and liver transcriptome characterization and genetic marker discovery in the**  
2 **farmed meagre, *Argyrosomus regius***

3

4 **Manousaki, T.<sup>1</sup>, Tsakogiannis, A.<sup>1,2</sup>, Lagnel, J.<sup>3</sup>, Kyriakis D.<sup>1,4</sup>, Duncan, N.<sup>5</sup>, Estevez, A.<sup>5</sup>,**  
5 **Tsigenopoulos, C. S.<sup>1\*</sup>**

6

7 <sup>1</sup> *Institute of Marine Biology, Biotechnology and Aquaculture, Hellenic Centre for Marine*  
8 *Research, Heraklion, Greece*

9 <sup>2</sup> *Department of Biology, University of Crete, Heraklion, Greece*

10 <sup>3</sup> *INRA PACA, UR 1052 GAFL, 67 Allée des Chênes, CS 60 094, 84143 Montfavet Cedex,*  
11 *France*

12 <sup>4</sup> *School of Medicine, University of Crete, Heraklion, Greece*

13 <sup>5</sup> *IRTA, Sant Carles de la Rapita, Tarragona, Spain*

14

15

16 *\*Author for Correspondence: Costas S. Tsigenopoulos, Institute of Marine Biology,*  
17 *Biotechnology and Aquaculture, Hellenic Centre for Marine Research, Heraklion, Greece, tel*  
18 *+30 2810337854, fax +30 2810337870, tsigeno@hcmr.gr*

19

20 **Abstract**

21 Meagre (*Argyrosomus regius*), a teleost fish of the family Sciaenidae, is part of a group of  
22 marine fish species considered new for Mediterranean aquaculture representing the larger fish  
23 cultured in the region. Meagre aquaculture started ~25 years ago in West Mediterranean, and  
24 the supply of juveniles has been dominated by few hatcheries. This fact has raised concerns  
25 on possible inbreeding, urging the need for genetic information on the species and for an  
26 assessment of the polymorphisms found in the genome. To that end we characterized the  
27 muscle and liver transcriptome of a pool of meagre individuals, from different families and

28 phenotypic size, to obtain a backbone that can support future studies regarding physiology,  
29 immunology and genetics of the species. The assembled transcripts were assigned to a wide  
30 range of biological processes including growth, reproduction, metabolism, development,  
31 stress and behavior. Then, to infer its genetic diversity and provide a catalogue of markers for  
32 future use, we scanned the reconstructed transcripts for polymorphic genetic markers. Our  
33 search revealed a total of 42,933 high quality SNP and 20,581 STR markers. We found a  
34 relatively low rate of polymorphism in the transcriptome that may indicate that inbreeding has  
35 taken place. This study has led to a catalogue of genetic markers at the expressed part of the  
36 genome and has set the ground for understanding growth and other traits of interest in  
37 meagre.

38

39

40 Keywords: aquaculture, RNASeq, SNPs, STRs

## 41 **Introduction**

42 The meagre, *Argyrosomus regius* (Asso y del Rio 1801) is a teleost fish that belongs to the  
43 family Sciaenidae and is widely distributed along the eastern Atlantic Ocean coast and the  
44 entire Mediterranean Sea (Chao, 1986). Throughout the distribution, meagre holds an  
45 important role in fisheries and now represents one of the newly emerging and promising  
46 aquaculture species across the Mediterranean region. There appears to be few fast growing  
47 large aquaculture species in the Mediterranean region and meagre together with greater  
48 amberjack (*Seriola dumerili*) fill this niche. Meagre aquaculture started in late nineties in  
49 France and Italy and since then has expanded in other European countries (FAO, 2015).  
50 Interestingly, meagre fry production has been for years carried out through a single hatchery  
51 in France (Monfort 2010), a fact that raises concerns regarding the genetic diversity of the  
52 European aquaculture stocks and requires evaluation.

53 Coupled with the increasing interest in the aquaculture industry, meagre is being explored in  
54 various fields, such as reproduction and broodstock management (DUNCAN *et al.* 2012;  
55 MYLONAS *et al.* 2015) and spawning with (MYLONAS *et al.* 2013b; FERNÁNDEZ *et al.* 2014)  
56 and without (MYLONAS *et al.* 2013a; SOARES *et al.* 2015) hormones, larval rearing conditions  
57 (ESTEVEZ *et al.* 2007; ROO *et al.* 2010; VALLÉS AND ESTÉVEZ 2013)), larval nutritional  
58 requirements (CAMPOVERDE AND ESTEVEZ 2017; EL KERTAOUI *et al.* 2017), skeletal  
59 development (CARDEIRA *et al.* 2012) and digestion (CASTRO *et al.* 2013; PAPADAKIS *et al.*  
60 2013). Although studies are accumulating for various fields of species biology, the genetic  
61 information and stock structure are only scarcely studied with the available information being  
62 limited to only 148 nucleotide and 71 protein entries in NCBI (as of 16 May 2017).

63 The paucity of available genetic resources is currently an impediment to any future effort for  
64 genetic improvement in the species. However, through next generation sequencing (NGS)  
65 technologies, and in particular RNA-Sequencing (RNA-Seq), one can collect sequence  
66 information for thousands of genes in a single experiment (WANG *et al.* 2009). Transcriptome  
67 characterization is one of the main applications of NGS as it lays the groundwork for future  
68 studies on physiology, genetics, immunology, etc., creates inventories and gives access to

69 thousand of single nucleotide polymorphisms (SNP) and short tandem repeats (STR) markers.  
70 Up to now, it has been widely used for numerous fish species leading to a tremendous pool of  
71 genetic knowledge (e.g. see database FISHIT [<http://www.fish-it.org/hcmr/>] for 20  
72 transcriptomes). Especially for farmed species, RNA-Seq can be an invaluable source of  
73 genetic information that can facilitate research on reproduction and sex dimorphism  
74 (MANOUSAKI *et al.* 2014; PALSTRA *et al.* 2015), physiology (KAITETZIDOU *et al.* 2012;  
75 TELES *et al.* 2013; MININNI *et al.* 2014), growth (GARCIA DE LA SERRANA *et al.* 2015),  
76 metabolism (CEREZUELA *et al.* 2013; DE SANTIS *et al.* 2015; GLENCROSS *et al.* 2015),  
77 immunity and disease resistance (CALDUCH-GINER *et al.* 2012; SARROPOULOU *et al.* 2012;  
78 ALI *et al.* 2014; MARANCIK *et al.* 2015; VALENZUELA-MIRANDA *et al.* 2015) and genetic  
79 marker discovery (MANOUSAKI *et al.* 2014; YU *et al.* 2014).

80 The goal of this paper was two-fold. First, we sought to characterize the transcriptome of  
81 meagre and build a solid transcriptomic reference for the species. Then, we aimed at assessing  
82 the genetic polymorphism of the species by including a thorough SNP and STR discovery  
83 from multiple individuals of farmed meagre. The discovered markers will set the groundwork  
84 for future marker-assisted selection for the species.

85

## 86 **Materials & Methods**

87

### 88 **Sample collection**

89 Animal care was carried out according to the “Guidelines for the treatment of animals in  
90 behavioural research and teaching” (Animal Behaviour 2001). Fish were selected  
91 (aquaculture facilities, IRTA, Spain, 21 August 2014, Table 1) from five different meagre  
92 crosses (families) that resulted from a mix of cultured and wild outbred parents. Muscle and  
93 liver tissues were dissected and preserved in RNAlater® (Applied Biosystems, Foster City,  
94 CA, USA). Sixteen meagre individuals were randomly selected for RNA Sequencing analysis  
95 (Supplementary Table 1).

96 **Table 1.** MIxS information for transcriptome assembly of *Argyrosomus regius*.

Item	Description
Classification	Eukaryota; Animalia; Chordata; Vertebrata; Actinopterygii; Percomorphaceae; Sciaenidae; <i>Argyrosomus regius</i>
Investigation type	Eukaryote transcriptome <sup>[1]</sup>
Project name	Meagre transcriptome
<i>Environment</i>	
Latitude, longitude	41.634502, 2.167185
Geographical location	IRTA, Spain
Collection date	21/8/2014
Biome	marine biome (ENVO_00000447)
Feature	fish farm (ENVO:00000294)
Material	sea water (ENVO:00002149)
<i>Sequencing</i>	
Sequencing method	Illumina HiSeq 2500 paired-end
Estimated size <sup>[1]</sup>	100 Mb
Organ or tissue source	Liver, muscle tissue
<i>Assembly</i>	
Method <sup>[1]</sup>	De novo assembly
Program	Trinity trinitymaseq_r2013-02-25
Finishing strategy	High quality transcriptome assembly
<i>Data accessibility</i>	
Database name	NCBI <sup>[1]</sup>
Project name	PRJNA397355, PRJNA399060
Sample name	SRR5903997, SRR5903998, SAMN07522546

97

98 **RNA extraction, library preparation and sequencing**

99 Muscle and liver tissues from the 16 individuals were collected in a sterile and RNase-free  
100 way. Following the manufacturer's recommendations, soaked tissues in RNAlater®, were  
101 stored at 4°C overnight and then were transferred to -80°C until further processing. For both  
102 tissue types the samples were grinded under liquid nitrogen using pestle and mortar. Liver is  
103 rich in RNA and thus a small amount of tissue was adequate to purify a high quality RNA  
104 using Qiagen's RNeasy Plus extraction kit (QIAGEN®). In contrary, because of the low cell  
105 density and the fibrous nature of muscle tissue, the yield of total RNA is low. In that case, a  
106 much larger proportion of tissue was grinded, focusing on pulverizing it into a fine powder  
107 while keeping it completely frozen. Complete homogenization achieved in TRIzol® reagent  
108 (Invitrogen, Carlsbad, CA, U.S.) using needle and syringe and high integrity total RNA was  
109 isolated according to the manufacturer's instructions.

110 The quantity of the isolated RNA was measured spectrophotometrically with NanoDrop®  
111 ND-1000 (Thermo Scientific), while its quality and integrity were tested on an agarose gel  
112 (electrophoresis in 1.5% w/v) and further on an Agilent Technologies 2100 Bioanalyzer  
113 (Agilent Technologies). All samples had an RNA Integrity Number (RIN) value higher than  
114 8. Following extraction, RNA from different individuals was pooled in equal quantities for  
115 each of the two tissue types. Then, an RNASeq library was constructed for each tissue  
116 following standard Illumina TruSeq protocols. The two libraries were loaded into one lane of  
117 an Illumina HiSeq2500 instrument (2x100bp). Raw reads produced are available at NCBI  
118 SRA with the project ID PRJNA397355 (Table 1).

119

120 **Raw read quality control**

121 Read quality was assessed with FastQC  
122 (<http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>) and subjected to quality control  
123 following a pipeline including multiple steps and published elsewhere (ILIAS *et al.* 2015)  
124 Briefly, we first used Scythe - a bayesian adapter trimmer (version 0.994 BETA)  
125 (<https://github.com/vsbuffalo/scythe>), to identify adapter substrings in reads. Scythe

126 recognizes adapter sequences taking into account quality information especially at the 3' end  
127 where quality falls. Thus, this step was applied prior to any quality-based trimming (prior  
128 contamination rate set in 0.1 '-p 0.1'). Then, low quality (Phred quality threshold of 20 and  
129 minimum reads length of 45 nt) reads trimming was performed with Sickle  
130 (<https://github.com/najoshi/sickle>). Sickle scans the reads in sliding windows and based on  
131 the quality it determines whether a read requires trimming in the two ends or complete  
132 removal (parameters 'pe -g -t sanger -q 20 -l 45'). The surviving reads were used as input  
133 to Trimmomatic (BOLGER *et al.* 2014) to further remove 5' and 3' adaptor sequences and  
134 apply extra filtering steps (parameters 'PE -phred33 ILLUMINACLIP:adapter\_file.fa:2:30:10  
135 LEADING:3 TRAILING:3 SLIDINGWINDOW:4:25 MINLEN:45 CROP:99'). Finally, we  
136 used PrinSeq (SCHMIEDER AND EDWARDS 2011) to filter out low complexity sequences  
137 (threshold entropy value of 30) and perform poly A/T 5' tail (minimum of 5 A/T) trimming.

138

### 139 **Transcriptome assembly and annotation**

140 Following reads pre-processing, we pooled the filtered reads from both liver and muscle  
141 samples and built a transcriptome assembly using Trinity (GRABHERR *et al.* 2011)  
142 (`trinityrnaseq_r2013-02-25`; default kmer 25; minimum contig length of 200 nucleotides).  
143 This Transcriptome Shotgun Assembly project has been deposited at DDBJ/EMBL/GenBank  
144 under the accession GFGV000000000, BioProject PRJNA399060 (Table 1).

145 To evaluate the completeness of the reconstructed assembly with used the software BUSCO  
146 v2 (SIMAO *et al.* 2015) through gVolante (<https://gvolante.riken.jp/>) selecting the Core  
147 Vertebrate Gene (CVG) set (HARA *et al.* 2015).

148 The assembled transcripts were annotated through a BLASTx similarity search against the  
149 SWISSPROT protein database (e-value threshold  $10^{-5}$ ; keeping the top twenty hits). To  
150 improve the speed of this long process, we implemented BLASTx in parallel using  
151 ParaNOblast (<https://github.com/jacqueslagnel/ParaNoBLast>) described in (LAGNEL *et al.*  
152 2009). Further, scan against protein domain signatures was done with InterProScan (JONES *et al.*  
153 *et al.* 2014), which was run in parallel splitting the query in 100 subqueries and merging the



154 output with custom scripts. Blast and InterProScan results were input in Blast2GO V.2.8.0  
155 (CONESA *et al.* 2005) where GO terms and Enzyme Commission (EC) numbers were retrieved  
156 and assigned to transcripts. Finally, sequences with EC numbers were further annotated with  
157 Kyoto Encyclopedia of Gene and Genome (KEGG) pathways using custom perl scripts.

158

### 159 **Genetic marker discovery**

160 To detect single nucleotide polymorphisms (SNPs) we used GATK pipeline (MCKENNA *et al.*  
161 2010) according to the GATK best practices (DANECEK *et al.* 2011; VAN DER AUWERA *et al.*  
162 2013). The implemented steps included mapping of the filtered reads to the assembled  
163 transcriptome using the highly accurate and fast aligner STAR (DOBIN *et al.* 2013), duplicate  
164 marking and sorting with Picard (<https://github.com/broadinstitute/picard>) and finally variant  
165 calling and filtering (options -window 35 -cluster 3 -filterName FS -filter "FS > 30.0" -  
166 filterName QD -filter "QD < 2.0"). The filtering options chosen filtered out SNPs that form  
167 clusters (more than 3 SNPs in a window of 35 bases), and variants with QualByDepth (QD) <  
168 2.0 and FisherStrand (FS) > 30 accounting for variant quality and strand bias. Finally only  
169 SNPs with at least 15 reads coverage were kept.

170 Following the filtering steps conducted within GATK, variants without a "PASS" filter tag  
171 were excluded, ii. the variants that included insertions and deletions (indels), and SNPs with  
172 more than two alleles. Further, to avoid sampling the same SNP locus twice due to alternative  
173 splicing, we kept only those identified in the longest transcript of each locus. Finally, to check  
174 whether they belong to non-coding (3'UTR and 5'UTR) or coding regions (first, second or  
175 third codon positions), we excluded those that were identified in transcripts without ORF and  
176 analyzed the rest with a custom python script taking into account the SNP position in the  
177 longest predicted ORF and the ORF coordinates in each transcript.

178 To detect short tandem repeats (STRs) we scanned the longest transcript of each assembled  
179 locus using the software Phobos ([http://www.ruhr-uni-  
180 bochum.de/ecoevo/cm/cm\\_phobos.htm](http://www.ruhr-uni-bochum.de/ecoevo/cm/cm_phobos.htm)). In particular, we detected non-exact STRs with 2–10  
181 repeat unit length and a minimum length of 20 nucleotides. A custom Perl script was used to

182 parse the output. Once again, for markers included on transcripts with ORFs, STRs were  
 183 categorized in coding, 3'UTR or 5'UTR according to the position in relation to the longest  
 184 ORF within the longest transcripts per locus using python scripts.

185

186

187 **Results & Discussion**

188

189 **Meagre transcriptome reconstruction and annotation**

190 Illumina sequencing of the multi-individual liver and muscle libraries yielded in total  
 191 523,137,020 raw reads that were subjected to a series of quality control filters (Table 2).  
 192 Following filtering, 341,439,304 paired reads were kept and used for assembly and  
 193 downstream analyses.

194

195 **Table 2.** The raw read quality control process and the read survival following each filtering  
 196 step.

Filtering steps	Surviving reads Muscle	Surviving reads Liver
Raw	280,804,390	242,332,630
Scythe*	280,804,390	242,332,630
Sickle	250,526,202	216,487,756
Trimmomatic	209,252,073	181,232,317
PrinSeq**	182,802,502	158,636,802

197 \*NOTE: Scythe does not eliminate sequences

198 \*\*NOTE: Only paired reads surviving PrinSeq filtering step were used for assembly and downstream  
 199 analyses

200

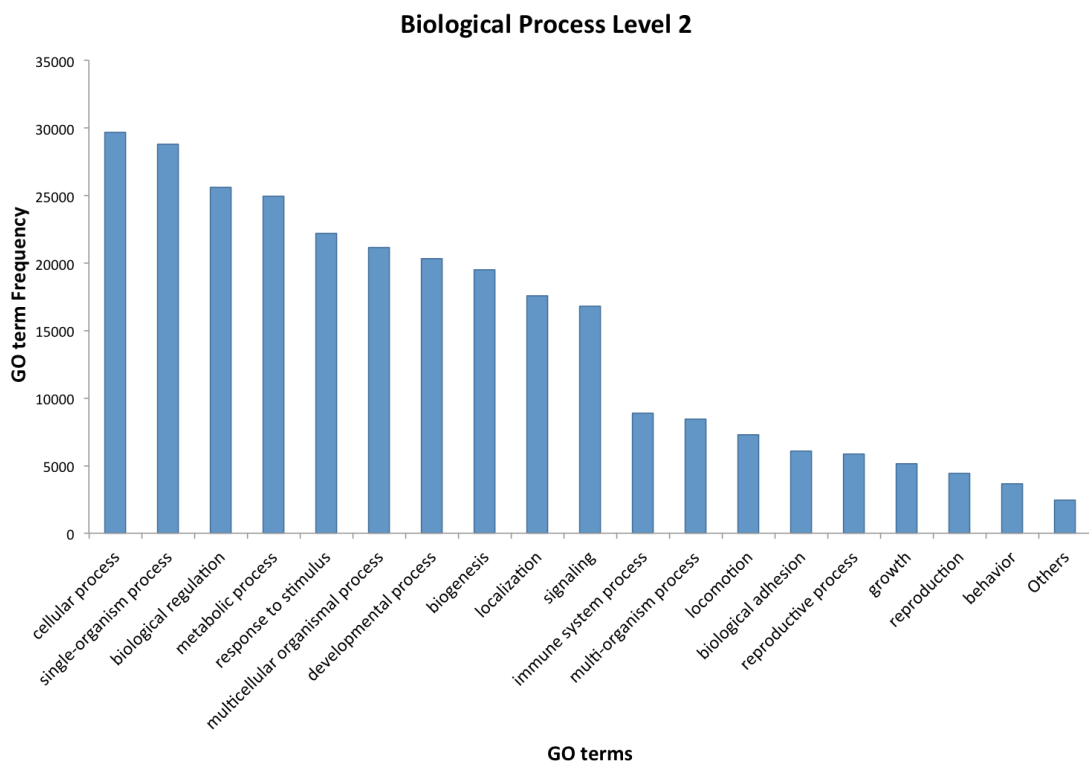
201 The high quality reads from both muscle and liver libraries were pooled and used for  
202 reconstructing the reference transcriptome of meagre. The strategy followed involved pooling  
203 all liver information from multiple individuals in a single library and the same for muscle.  
204 This design enables the identification of genetic markers across the individuals' transcriptome  
205 in a cost-effective manner, as only two libraries are constructed and sequenced deeply.  
206 However, it holds limitations, first in terms of identifying markers at the individual level, and  
207 second pooling individuals bears drawbacks especially regarding the unequal representation  
208 of each individual's alleles in the final read count leading to erroneous allele frequency  
209 estimations in population studies (SCHLÖTTERER *et al.* 2014), which were not the scope of  
210 this study.

211 Following assembly, the resulted transcriptome comprised of 95,945 transcripts belonging to  
212 80,807 loci with N50 value of 2,183, average length of 1,059 nucleotides and 46.19% GC  
213 content. To evaluate to completeness of the assembly, we ran BUSCOv2 and gVolante to find  
214 that out of 233 queried genes, 208 (89.27) were complete, 15 were partial (summing up to  
215 95.71% complete and partial genes) and only 10 genes (4.29%) were missing. The results  
216 revealed a satisfying assembled transcriptome covering the great majority of meagre genes.  
217 However, future sequencing of more tissue types would lead to a more complete  
218 transcriptome in the species.

219 To annotate the assembly, we conducted a blastx similarity search against the highly curated  
220 SWISSPROT database. The results revealed that 33,638 out of 95,945 transcripts were  
221 significantly homologous to a known SWISSPROT sequence. Finally, targeted BLASTN  
222 search of meagre transcripts against tilapia cDNA retrieved 15,589 unique tilapia genes as top  
223 hits, once again confirming the thorough representation of the expected geneset in the  
224 transcriptome.

225 Following similarity search through blast, GO mapping resulted in 31,986 annotated  
226 sequences. The most important GO terms in the 'biological process' ontology at the level 2  
227 are shown in Figure 1. Search for InterPro domains resulted in 46,647 sequences annotated  
228 with protein domains and raised the number of GO annotated transcripts to 34,252. Then, EC

229 number mapping through GO terms resulted in 1,016 potential enzymes in 8,682 total  
 230 transcripts with EC:6.3.2.19 ‘ubiquitin-protein ligase’ as the top enzyme group in meagre  
 231 transcriptome, followed by EC:3.6.1.3 (adenosine triphosphatase) and EC:2.3.1.48 (histone  
 232 acetyltransferase). Based on EC mapping, we identified the corresponding KEGG pathways  
 233 to find that the EC-annotated 8,682 sequences are involved in 382 total KEGG pathways. The  
 234 most highly represented pathway was MAPK signaling pathway, followed by Purine  
 235 metabolism and PI3K-Akt signaling pathway. A summary of the annotation results is  
 236 presented in Table 3 and detailed annotation is given in Supplementary Table 2.  
 237



238

239 **Figure 1.** Gene Ontology functional characterization of meagre assembled transcriptome.

240 Terms are shown for biological process level 2.

241

242

243

**Table 3.** Annotation Summary

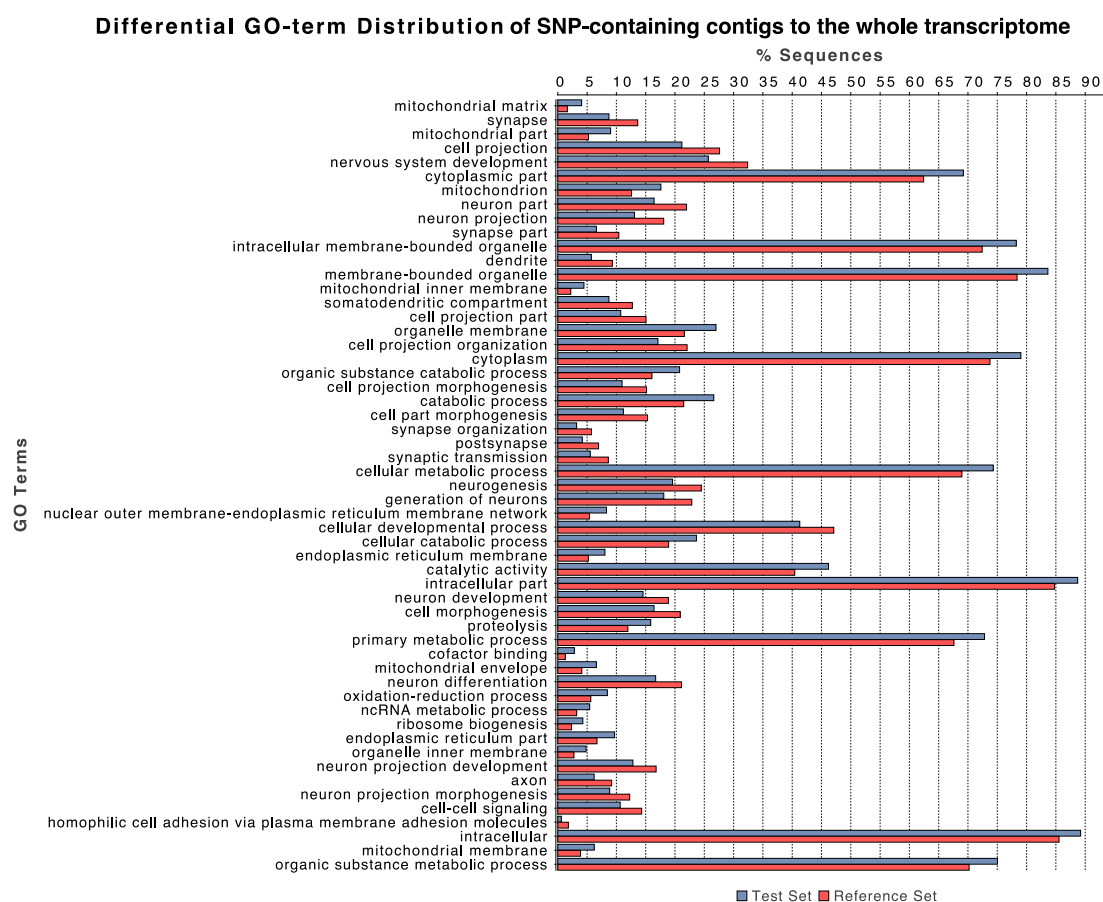
Annotation steps	Contigs	% Contigs
	95,945	100
BLAST	33,638	35.05
InterPro	46,647	48.61
With IPR number	27,907	29.09
With $\geq 1$ GO	34,252	35.70
Blast2Go annotated	33,220	34.62
EC	8,440	8.80
KEGG 380 pathways	2,475	2.58

244

245 **Meagre genetic markers**

246 Following the transcriptome characterization of meagre, we aimed at scanning for both SNPs  
 247 and STRs across meagre transcriptome.

248 Our SNP search revealed a total of 42,933 high quality markers located in 14,544 transcripts  
 249 (Supplementary Table 3). A GO enrichment analysis (FDR 0.05) of the contigs containing  
 250 SNPs compared to the assembly revealed a significant underrepresentation of genes related to  
 251 the nervous system (Figure 2; Supplementary Table 4), which might reflect the evolutionary  
 252 pressure for conservation in this group of genes.



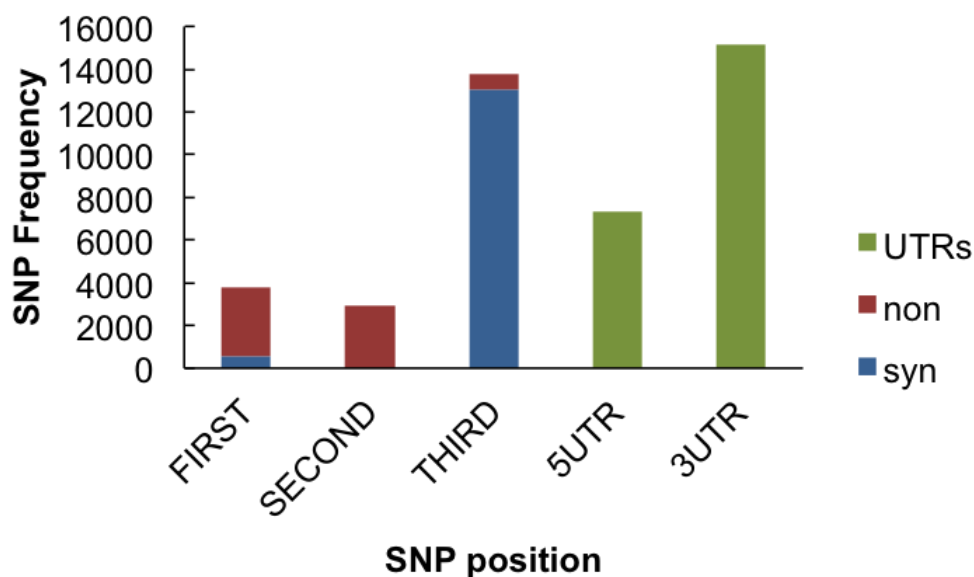
254

255 **Figure 2.** Top significantly over and under represented GO categories of SNP-containing  
 256 contigs compared to the whole meagre transcriptome (only GO terms with FDR cut-off < 10<sup>-15</sup>  
 257 are shown).

258

259 Downstream analyses showed that most SNPs were located in the UTRs (15,149 SNPs) and  
 260 at the third codon position (13,768 SNPs) in accordance to the sequence conservation pattern  
 261 observed in coding sequences. SNPs that fall within the predicted open reading frame, result  
 262 mostly to synonymous changes for the third codon position (731 non-synonymous and 13,037  
 263 synonymous SNPs), only in non-synonymous changes for the second codon position (2,937  
 264 SNPs) and mainly in non-synonymous changes for the first codon position (3,174 non-  
 265 synonymous and 580 synonymous SNPs), as expected from the genetic code degeneracy  
 266 (Figure 3).

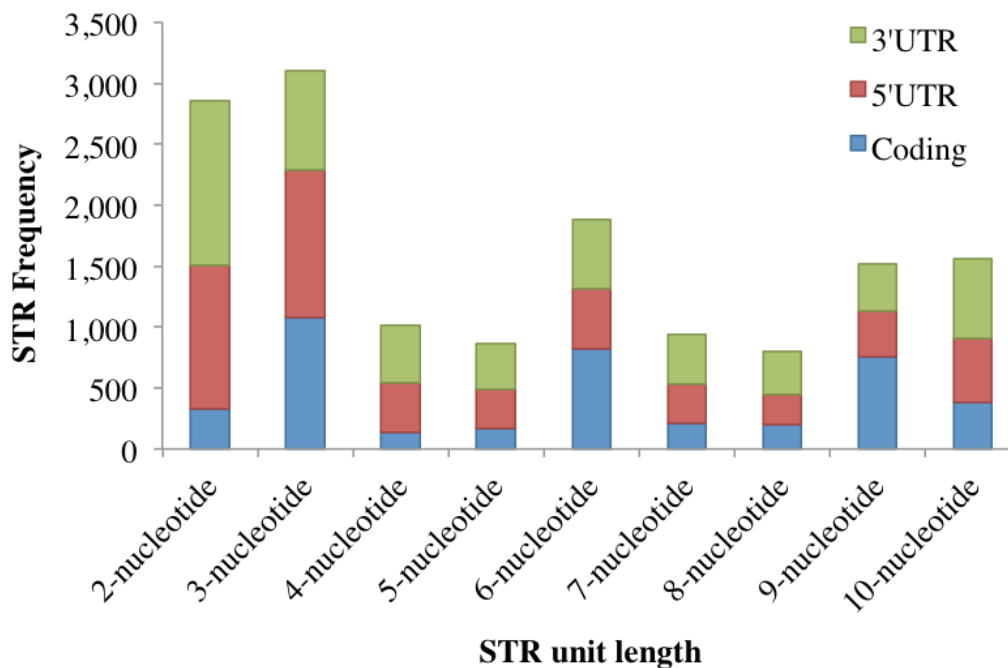
267 Based on the frequency of SNPs found (1.22 SNPs per 1,000 bp), we describe a relatively low  
 268 rate of polymorphism in the species probably reflecting the inbreeding that has taken place.  
 269 Examples of other teleosts with reported higher SNP rate are turbot (one SNP per 302 bp;  
 270 VERA et al. 2013), channel catfish (one SNP per 93 bp; LIU et al. 2016) and salmon with only  
 271 slightly higher SNP rate (2.52 SNPs per 1,000 bp; TINE et al. 2014).  
 272 Our STR search revealed 20,581 markers ranging from 2-10 unit length transcriptome-wide  
 273 (Supplementary Table 5). Following breaking down of the discovered STRs to those that fall  
 274 within the coding regions (14,546 STRs) and those that do not, we show that 3-mer STRs  
 275 (with unit length of 3, and then 6 and 9) fell into the first category (coding). The distribution  
 276 of the 3-mer STRs are significantly higher than expected in the coding regions and  
 277 significantly lower in the UTRs (chi-square p-value < 0.00) as expected due to the non-  
 278 disturbance of the open reading frame by those repeats (Figure 4). The rest were mostly in the  
 279 UTRs or in transcripts without ORFs.  
 280 The distribution of both marker types in the different regions of the transcriptome is  
 281 consistent with the expected distribution of genetic variants as discovered in other studies as  
 282 well (e.g. MANOUSAKI *et al.* 2014) and provide a high quality dataset for future genetic  
 283 analysis in this new but important to aquaculture species.  
 284



285

286 **Figure 3.** Distribution of SNPs along the coding and non-coding part of the transcripts. SNPs  
 287 found within coding regions are separated to first, second and third position and are  
 288 characterized as synonymous or nonsynonymous based on causing or not an amino acid shift  
 289 in the protein sequence.

290  
 291



292

293 **Figure 4.** Distribution of STRs along the coding and non-coding part of the transcripts.

294

#### 295 **Meagre transcriptome gene content**

296 Following the assembly annotation and genetic marker discovery, we sought to identify  
 297 transcripts that might be involved in important biological functions. For example, growth is  
 298 one of meagre's most important phenotypic traits for aquaculture. To that end, we extracted  
 299 the sequences associated to growth by selecting transcripts with the search term 'growth'  
 300 within GO annotation descriptions through Blast2GO. Our search revealed 7,121 sequences  
 301 (Supplementary Table 6). The SNPs and STRs of those particular genes might serve as  
 302 valuable resource for identifying variants linked to this critical trait of the species.



303 To further explore the gene content of meagre, we searched the assembled transcriptome for  
304 sequences that include representative terms in the GO annotations. More specifically, we  
305 found 3,144 sequences related to ‘immune’ functions, 15,274 genes involved in  
306 ‘development’, 4,300 genes involved in “stress”, 1,012 genes involved in reproduction, 5,196  
307 genes involved in metabolism and 2,168 genes involved in behavior (Supplementary Table 6).

308

### 309 **Conclusions**

310 Our study has built the first transcriptome assembly and at the same time the first next-  
311 generation based genomics resource for meagre. Following the annotated transcriptome, we  
312 launched a genetic marker discovery pipeline that led to the construction of a valuable dataset  
313 of SNPs and STRs transcriptome-wide coming out from a pool of 16 individuals. The low  
314 rate of polymorphism discovered imply that the transcriptome of meagre has been possibly  
315 shaped by inbreeding, a factor that raises even more the risk for further inbreeding through  
316 aquaculture. The provided assembly and genetic markers dataset will lay the groundwork for  
317 further studies of meagre biology and genetics and will set the basis for future applications of  
318 genetic breeding and marker-assisted selection for the species.

319

320

### 321 **Acknowledgements**

322 This study was funded by the European Union’s Seventh Framework Programme for  
323 research, technological development and demonstration (KBBE-2013-07 single stage, GA  
324 603121, DIVERSIFY). The sequencing service was provided by the Norwegian Sequencing  
325 Centre ([www.sequencing.uio.no](http://www.sequencing.uio.no)), a national technology platform hosted by the University of  
326 Oslo (UiO) and supported by the ‘Functional Genomics’ and ‘Infrastructure’ programs of the  
327 Research Council of Norway and the South-Eastern Regional Health Authorities.

328

### 329 **Supplementary Tables Legends**

330 **Supplementary Table 1.** Summary including weight/length of sampled individuals.

331 **Supplementary Table 2.** Detailed annotation of meagre transcriptome including sequence  
332 description as defined by Blast2GO based on the annotation of the blast hits, the number of  
333 blast hits, the minimum e-value, the mean percentage of similarity, the number of GO terms,  
334 the assigned GO terms, EC numbers and InterProScan results.

335 **Supplementary Table 3.** The high quality SNP dataset discovered in meagre transcriptome.  
336 SNPs that remain after filtering and are located in the longest open reading frame of each  
337 gene are reported. For each SNP, provided information include: the respective contig, the  
338 open reading frame selected (ORF\_Region), the starting point of the ORF (Start\_ORF) and  
339 the position of the SNP (SNP\_pos). Each SNP is characterized as coding or noncoding  
340 according to whether it falls inside or outside the coding regions and each noncoding SNP is  
341 annotated as upstream or downstream depending on whether it is found in the 5' or the 3'  
342 UTR of the gene. Further, SNPs that fall within the coding regions are broken down to those  
343 that fall within the first, second or third codon position and are also characterized as  
344 synonymous or nonsynonymous depending on whether the two alleles code for the same  
345 amino acid or not.

346 **Supplementary Table 4.** GO terms that are over- or under-represented in the SNP-containing  
347 genes compared to the whole assembly through a Fisher's exact test (FDR threshold 0.05).  
348 The test is run through Blast2GO for the three GO categories (P: Biological Process, C:  
349 Cellular Component, F: Molecular Function).

350 **Supplementary Table 5.** The high quality STR dataset discovered in meagre transcriptome.  
351 For each STR information regarding the respective contig (Seq ID), the unit length, the  
352 number of units in the reference (# of units), the start and stop position in the contig (start,  
353 stop), the total length (length), the length ignoring insertions/deletions (norm\_length), number  
354 of mismatches (mis), number of insertions (ins), number of deletions (del), the unit motif  
355 (motif) and the total STR sequence (seq) are given.

356 **Supplementary Table 6.** The gene content of meagre transcriptome. The list and annotation  
357 of contigs that include "growth", "behaviour", "development", "reproduction", "metabolism",

358 “stress”, “immune” within the GO terms description and may have a possible role in the  
359 respective functions.

360

## 361 **References**

- 362 Ali, A., C. Rexroad, G. Thorgaard, J. Yao and M. Salem, 2014 Characterization of the  
363 rainbow trout spleen transcriptome and identification of immune-related genes.  
364 *Frontiers in Genetics* 5.
- 365 Andrews S. (2010). FastQC: a quality control tool for high throughput sequence data.  
366 Available online at: <http://www.bioinformatics.babraham.ac.uk/projects/fastqc>
- 367 Bolger, A. M., M. Lohse and B. Usadel, 2014 Trimmomatic: a flexible trimmer for Illumina  
368 sequence data. *Bioinformatics* 30: 2114-2120.
- 369 Calduch-Giner, J. A., A. Sitja-Bobadilla, G. C. Davey, M. T. Cairns, S. Kaushik *et al.*, 2012  
370 Dietary vegetable oils do not alter the intestine transcriptome of gilthead sea bream  
371 (*Sparus aurata*), but modulate the transcriptomic response to infection with  
372 *Enteromyxum leei*. *BMC Genomics* 13: 470.
- 373 Campoverde, C., and A. Estevez, 2017 The effect of live food enrichment with  
374 docosahexaenoic acid (22:6n-3) rich emulsions on growth, survival and fatty acid  
375 composition of meagre (*Argyrosomus regius*) larvae. *Aquaculture* 478: 16-24.
- 376 Cardeira, J., R. Vallés, G. Dionísio, A. Estévez, E. Gisbert *et al.*, 2012 Osteology of the axial  
377 and appendicular skeletons of the meagre *Argyrosomus regius* (Sciaenidae) and early  
378 skeletal development at two rearing facilities. *Journal of Applied Ichthyology* 28:  
379 464-470.
- 380 Castro, C., A. Pérez-Jiménez, F. Coutinho, P. Pousão-Ferreira, T. M. Brandão *et al.*, 2013  
381 Digestive enzymes of meagre (*Argyrosomus regius*) and white seabream (*Diplodus*  
382 *sargus*). Effects of dietary brewer's spent yeast supplementation. *Aquaculture* 416-  
383 417: 322-327.
- 384 Cerezuela, R., J. Meseguer and M. A. Esteban, 2013 Effects of dietary inulin, *Bacillus subtilis*  
385 and microalgae on intestinal gene expression in gilthead seabream (*Sparus aurata* L.).  
386 *Fish Shellfish Immunol* 34: 843-848.
- 387 Chao, L. N. (1986). Sciaenidae. In *Fishes of the Eastern Atlantic and Mediterranean Poissons*  
388 *de l'Atlantique du nord-est et de la Méditerranée* (Whitehead, P. J. P., Bauchot, M.-  
389 L., Hureau, J. C. & Tortonese, E., eds), pp. 865–874. Paris: Unesco.
- 390 Conesa, A., S. Gotz, J. M. Garcia-Gomez, J. Terol, M. Talon *et al.*, 2005 Blast2GO: a  
391 universal tool for annotation, visualization and analysis in functional genomics  
392 research. *Bioinformatics* 21.
- 393 Danecek, P., A. Auton, G. Abecasis, C. A. Albers, E. Banks *et al.*, 2011 The variant call  
394 format and VCFtools. *Bioinformatics* 27.
- 395 De Santis, C., J. F. Taylor, L. Martinez-Rubio, S. Boltana and D. R. Tocher, 2015 Influence  
396 of development and dietary phospholipid content and composition on intestinal  
397 transcriptome of Atlantic salmon (*Salmo salar*). *PLoS One* 10: e0140964.
- 398 Dobin, A., C. A. Davis, F. Schlesinger, J. Drenkow, C. Zaleski *et al.*, 2013 STAR: ultrafast  
399 universal RNA-seq aligner. *Bioinformatics* 29: 15-21.
- 400 Duncan, N., A. Estévez, J. Porta, I. Carazo, F. Norambuena *et al.*, 2012 Reproductive  
401 development, GnRH-induced spawning and egg quality of wild meagre  
402 (*Argyrosomus regius*) acclimatised to captivity. *Fish Physiology and Biochemistry*  
403 38: 1273-1286.
- 404 El Kertaoui, N., C. M. Hernández-Cruz, D. Montero, M. J. Caballero, R. Saleh *et al.*, 2017  
405 The importance of dietary HUFA for meagre larvae (*Argyrosomus regius*; Asso,  
406 1801) and its relation with antioxidant vitamins E and C. *Aquaculture Research* 48:  
407 419-433.

408 Estévez, A., Treviño, L., Gisbert, E. 2007. La densidad larvaria inicial afecta al crecimiento  
409 pero no a la supervivencia de las larvas de corvina (*Argyrosomus regius*) en cultivo.  
410 Paper presented at the XI Spanish Aquaculture Congress, Vigo, Spain, 24-28  
411 September 2007.

412 FAO, 2015 <http://www.fao.org/home/en/>

413 Fernández, J., M. Á. Toro, A. K. Sonesson and B. Villanueva, 2014 Optimizing the creation  
414 of base populations for aquaculture breeding programs using phenotypic and genomic  
415 data and its consequences on genetic progress. *Frontiers in Genetics* 5.

416 Fernandez-Palacios, H., Schuchardt, D., Roo, J., Izquierdo, M., Hernandez-Cruz, C.M.,  
417 Duncan, N. 2014. Dose-dependent effect of a single GnRH $\alpha$  injection on the  
418 spawning of meagre ("*Argyrosomus regius*") broodstock reared in captivity. *Spanish*  
419 *journal of agricultural research, Instituto Nacional de Investigación y Tecnología*  
420 *Agraria y Alimentaria (INIA)*, 12(4): 1038-1048.

421 Garcia de la Serrana, D., R. H. Devlin and I. A. Johnston, 2015 RNAseq analysis of fast  
422 skeletal muscle in restriction-fed transgenic coho salmon (*Oncorhynchus kisutch*): an  
423 experimental model uncoupling the growth hormone and nutritional signals  
424 regulating growth. *BMC Genomics* 16: 564.

425 Glencross, B. D., C. De Santis, B. Bicskei, J. B. Taggart, J. E. Bron *et al.*, 2015 A  
426 comparative analysis of the response of the hepatic transcriptome to dietary  
427 docosahexaenoic acid in Atlantic salmon (*Salmo salar*) post-smolts. *BMC Genomics*  
428 16: 684.

429 Grabherr, M. G., B. J. Haas, M. Yassour, J. Z. Levin, D. A. Thompson *et al.*, 2011 Full-length  
430 transcriptome assembly from RNA-Seq data without a reference genome. *Nat*  
431 *Biotech* 29: 644-652.

432 Guidelines for the treatment of animals in behavioural research and teaching. *Anim Behav.*  
433 2001, 61 (1): 271-275.

434 Hara, Y., K. Tatsumi, M. Yoshida, E. Kajikawa, H. Kiyonari *et al.*, 2015 Optimizing and  
435 benchmarking de novo transcriptome sequencing: from library preparation to  
436 assembly evaluation. *BMC Genomics* 16: 977.

437 Ilias, A., J. Lagnel, D. E. Kapantaidaki, E. Roiditakis, C. S. Tsigenopoulos *et al.*, 2015  
438 Transcription analysis of neonicotinoid resistance in Mediterranean (MED)  
439 populations of *B. tabaci* reveal novel cytochrome P450s, but no nAChR mutations  
440 associated with the phenotype. *BMC Genomics* 16: 939.

441 Joshi NA, Fass JN. 2011. <https://github.com/najoshi/sickle>.

442 Jones, P., D. Binns, H.-Y. Chang, M. Fraser, W. Li *et al.*, 2014 InterProScan 5: genome-scale  
443 protein function classification. *Bioinformatics* 30: 1236-1240.

444 Kaitetzidou, E., D. Crespo, Y. Vraskou, E. Antonopoulou and J. V. Planas, 2012  
445 Transcriptomic response of skeletal muscle to lipopolysaccharide in the gilthead  
446 seabream (*Sparus aurata*). *Mar Biotechnol (NY)* 14: 605-619.

447 Lagnel, J., C. S. Tsigenopoulos and I. Iliopoulos, 2009 NOBLAST and JAMBLAST: new  
448 options for BLAST and a java application manager for BLAST results.  
449 *Bioinformatics* 25.

450 Liu, Z., S. Liu, J. Yao, L. Bao, J. Zhang *et al.*, 2016 The channel catfish genome sequence  
451 provides insights into the evolution of scale formation in teleosts. *Nat Commun* 7:  
452 11757.

453 Manousaki, T., A. Tsakogiannis, J. Lagnel, E. Sarropoulou, J. Z. Xiang *et al.*, 2014 The sex-  
454 specific transcriptome of the hermaphrodite sparid sharpnose seabream (*Diplodus*  
455 *puntazzo*). *BMC Genomics* 15: 1-16.

456 Marancik, D., G. Gao, B. Paneru, H. Ma, A. G. Hernandez *et al.*, 2015 Whole-body  
457 transcriptome of selectively bred, resistant-, control-, and susceptible-line rainbow  
458 trout following experimental challenge with *Flavobacterium psychrophilum*.  
459 *Frontiers in Genetics* 5.

460 Mayer, Christoph, Phobos 3.3.11, 2006-2010 Available online at: [http://www.ruhr-uni-](http://www.ruhr-uni-bochum.de/ecoevo/cm/cm_phobos.htm)  
461 [bochum.de/ecoevo/cm/cm\\_phobos.htm](http://www.ruhr-uni-bochum.de/ecoevo/cm/cm_phobos.htm)

462 McKenna, A., M. Hanna, E. Banks, A. Sivachenko, K. Cibulskis *et al.*, 2010 The Genome  
463 Analysis Toolkit: A MapReduce framework for analyzing next-generation DNA  
464 sequencing data. *Genome Research* 20: 1297-1303.

465 Mininni, A. N., M. Milan, S. Ferrarresso, T. Petochi, P. Di Marco *et al.*, 2014 Liver  
466 transcriptome analysis in gilthead sea bream upon exposure to low temperature. *BMC*  
467 *Genomics* 15: 765.

468 Monfort M.C. 2010 Present market situation and prospects of meagre (*Argyrosomus regius*),  
469 as an emerging species in Mediterranean aquaculture. General Fisheries Commission  
470 for the Mediterranean: Studies and Reviews, No. 89, Food and Agriculture  
471 Organization of the United Nations, Rome.

472 Mylonas, C. C., E. Fatira, P. Karkut, M. Papadaki, I. Sigelaki *et al.*, 2015 Reproduction of  
473 hatchery-produced meagre *Argyrosomus regius* in captivity III. Comparison between  
474 GnRH $\alpha$  implants and injections on spawning kinetics and egg/larval performance  
475 parameters. *Aquaculture* 448: 44-53.

476 Mylonas, C. C., N. Mitrizakis, C. A. Castaldo, C. P. Cerviño, M. Papadaki *et al.*, 2013a  
477 Reproduction of hatchery-produced meagre *Argyrosomus regius* in captivity II.  
478 Hormonal induction of spawning and monitoring of spawning kinetics, egg  
479 production and egg quality. *Aquaculture* 414: 318-327.

480 Mylonas, C. C., N. Mitrizakis, M. Papadaki and I. Sigelaki, 2013b Reproduction of hatchery-  
481 produced meagre *Argyrosomus regius* in captivity I. Description of the annual  
482 reproductive cycle. *Aquaculture* 414: 309-317.

483 Nishimura & Kuraku, 2016 gVolante. Available online at: <https://gvolante.riken.jp/>

484 Palstra, A. P., K. Fukaya, H. Chiba, R. P. Dirks, J. V. Planas *et al.*, 2015 The olfactory  
485 transcriptome and progression of sexual maturation in homing Chum salmon  
486 *Oncorhynchus keta*. *PLoS One* 10: e0137404.

487 Papadakis, I. E., M. Kentouri, P. Divanach and C. C. Mylonas, 2013 Ontogeny of the  
488 digestive system of meagre *Argyrosomus regius* reared in a mesocosm, and  
489 quantitative changes of lipids in the liver from hatching to juvenile. *Aquaculture* 388-  
490 391: 76-88.

491 Roo, J., C. M. Hernández-Cruz, C. Borrero, D. Schuchardt and H. Fernández-Palacios, 2010  
492 Effect of larval density and feeding sequence on meagre (*Argyrosomus regius*; Asso,  
493 1801) larval rearing. *Aquaculture* 302: 82-88.

494 Sarropoulou, E., J. Galindo-Villegas, A. Garcia-Alcazar, P. Kasapidis and V. Mulero, 2012  
495 Characterization of European sea bass transcripts by RNA SEQ after oral vaccine  
496 against *V. anguillarum*. *Mar Biotechnol (NY)* 14: 634-642.

497 Schlotterer, C., R. Tobler, R. Kofler and V. Nolte, 2014 Sequencing pools of individuals -  
498 mining genome-wide polymorphism data without big funding. *Nat Rev Genet* 15:  
499 749-763.

500 Schmieder, R., and R. Edwards, 2011 Quality control and preprocessing of metagenomic  
501 datasets. *Bioinformatics* 27: 863-864.

502 Simao, F. A., R. M. Waterhouse, P. Ioannidis, E. V. Kriventseva and E. M. Zdobnov, 2015  
503 BUSCO: assessing genome assembly and annotation completeness with single-copy  
504 orthologs. *Bioinformatics* 31: 3210-3212.

505 Soares, F., L. Ribeiro, M. Gamboa, S. Duarte, A. C. Mendes *et al.*, 2015 Comparative  
506 analysis on natural spawning of F1 meagre, *Argyrosomus regius*, with wild  
507 broodstock spawns in Portugal. *Fish Physiology and Biochemistry* 41: 1509-1514.

508 Teles, M., S. Boltana, F. Reyes-Lopez, M. A. Santos, S. Mackenzie *et al.*, 2013 Effects of  
509 chronic cortisol administration on global expression of GR and the liver  
510 transcriptome in *Sparus aurata*. *Mar Biotechnol (NY)* 15: 104-114.

511 Tine, M., H. Kuhl, P. A. Gagnaire, B. Louro, E. Desmarais *et al.*, 2014 European sea bass  
512 genome and its variation provide insights into adaptation to euryhalinity and  
513 speciation. *Nat Commun* 5: 5770.

514 Valenzuela-Miranda, D., S. Boltana, M. E. Cabrejos, J. M. Yanez and C. Gallardo-Escarate,  
515 2015 High-throughput transcriptome analysis of ISAV-infected Atlantic salmon  
516 *Salmo salar* unravels divergent immune responses associated to head-kidney, liver  
517 and gills tissues. *Fish Shellfish Immunol* 45: 367-377.  
518 Vallés, R., and A. Estévez, 2013 Light conditions for larval rearing of meagre (*Argyrosomus*  
519 *regius*). *Aquaculture* 376: 15-19.  
520 Van der Auwera, G. A., M. O. Carneiro, C. Hartl, R. Poplin, G. Del Angel *et al.*, 2013 From  
521 FastQ data to high confidence variant calls: the Genome Analysis Toolkit best  
522 practices pipeline. *Curr Protoc Bioinformatics* 43: 11 10 11-33.  
523 Vera, M., J. A. Alvarez-Dios, C. Fernandez, C. Bouza, R. Vilas *et al.*, 2013 Development and  
524 validation of single nucleotide polymorphisms (SNPs) markers from two  
525 transcriptome 454-runs of turbot (*Scophthalmus maximus*) using high-throughput  
526 genotyping. *Int J Mol Sci* 14: 5694-5711.  
527 Wang, Z., M. Gerstein and M. Snyder, 2009 RNA-Seq: a revolutionary tool for  
528 transcriptomics. *Nat Rev Genet* 10.  
529 Yu, Y., J. Wei, X. Zhang, J. Liu, C. Liu *et al.*, 2014 SNP Discovery in the Transcriptome of  
530 White Pacific Shrimp *Litopenaeus vannamei* by Next Generation Sequencing. *PLOS*  
531 *ONE* 9: e87218.  
532