

# De novo assembly of the Atlantic herring (*Clupea harengus*) genome and comparison to a previously published assembly

Sunnvør í Kongsstovu<sup>1</sup>, Svein-Ole Mikalsen<sup>2</sup>, Hannes Gislason<sup>2</sup>, Hóraldur Joensen<sup>2</sup>, Eydna í Homrum<sup>3</sup>, Jan Arge Jacobsen<sup>3</sup>, Debes H. Christiansen<sup>4</sup>, Thomas Damm Als<sup>5</sup>, Masa Roller<sup>6</sup>, David Martín-Gálvez<sup>6</sup>, Paul Flicek<sup>6</sup> and Hans Atli Dahl<sup>7</sup>

<sup>1</sup> Amplexa Genetics A/S, Hoyvíksvegur 51, FO-100 Tórshavn, Faroe Islands. <sup>2</sup> University of the Faroe Islands, Dept. of Natural Sciences, Nóatún 3, FO-100 Tórshavn. <sup>3</sup> Faroe Marine Research Institute, Nóatún 1, FO-100 Tórshavn, Faroe Islands. <sup>4</sup> Faroese Food and Veterinary Authority, National Reference Laboratory for Fish Diseases, Tórshavn, Faroe Islands. <sup>5</sup> Department of Biomedicine, Bartholins Allé 6, 8000 Aarhus C, Denmark. <sup>6</sup> European Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome Genome Campus, Hinxton, Cambridgeshire, CB10 1SD, UK. <sup>7</sup> Amplexa Genetics A/S, Tolderlundsvej 3B, 2. DK-5000 Odense C, Denmark.

## Introduction

The Atlantic herring (*Clupea harengus*) is one of the most abundant fish species in the world and is an important economical and nutritional resource. The Faroe Islands exported herring worth more than 400 million DKK in 2016, which amounts to 5% of the total value of exported goods that year<sup>1</sup>. The species is a highly migratory pelagic species with a complicated population structure. In order to keep the fisheries sustainable, knowledge of subpopulations, their migrations and mixing is important and would assist in the fight against illegal, unreported and unregulated fishing<sup>2</sup>.

In 2016 Barrio *et al.* published the assembled herring genome from a herring caught in the Baltic sea<sup>3</sup>. Nevertheless as this is a species of such economical and nutritional importance we have undertaken a second assembly of the herring genome, using an Atlantic herring from the North sea. We also used different sequencing technologies and a different bioinformatical approach, in the hope of gaining further information about population structure and standing variation.

Here, we sequenced the herring genome with both Illumina and Oxford Nanopore Technologies platforms. The herring genome was assembled using the Illumina data and then scaffolded using the MinION long reads.

The basic assembly statistics were compared to the published assembly and the assembly completeness of both assemblies evaluated.

## Results

**Table 1.** Summary statistics for published assembly (Barrio *et al.*) and the assemblies generated during this study. V1 was assembled using the AllPaths-LG<sup>4</sup> software and the Illumina reads. V2 was generated by scaffolding V1 with MinION reads from two runs, using the SSPACE-LongRead<sup>5</sup> software. V3 was generated by closing gaps in V1 using GapFiller<sup>6</sup> software and scaffolded with four runs of MinION reads.

Assembly	No. Scaffolds	N50	L50	Assembly length	Gap length
Barrio <i>et al.</i>	6,915	1,860,920	113	807,711,962	82,755,438
This Study V1	15,378	175,084	993	702,694,152	177,906,258
This Study V2	12,219	219,865	820	716,584,296	191,669,007
This Study V3	10,354	262,663	719	729,006,009	186,023,247

## Discussion

Comparing assemblies is complicated and a single comparative analysis often does not include enough parameters for a fair comparison. Therefore, we decided to compare the herring assembly from this study to a published assembly using the five different approaches below.

### Summary statistics:

According to the summary statistics the assemblies generated in this study are more fragmented than the previously published herring assembly.

### FRC:

The FRCs for the V2 and the previously published assembly are very similar. The V2 FRC is slightly steeper, indicating less assembly errors are present in the V2 assembly compared to the Barrio *et al.* assembly.

### Dot plot:

There are no macroscopic structural variations between V2 assembly from this study and the previously published assembly. However, there is some unexplained noise in the graph.

### BUSCO:

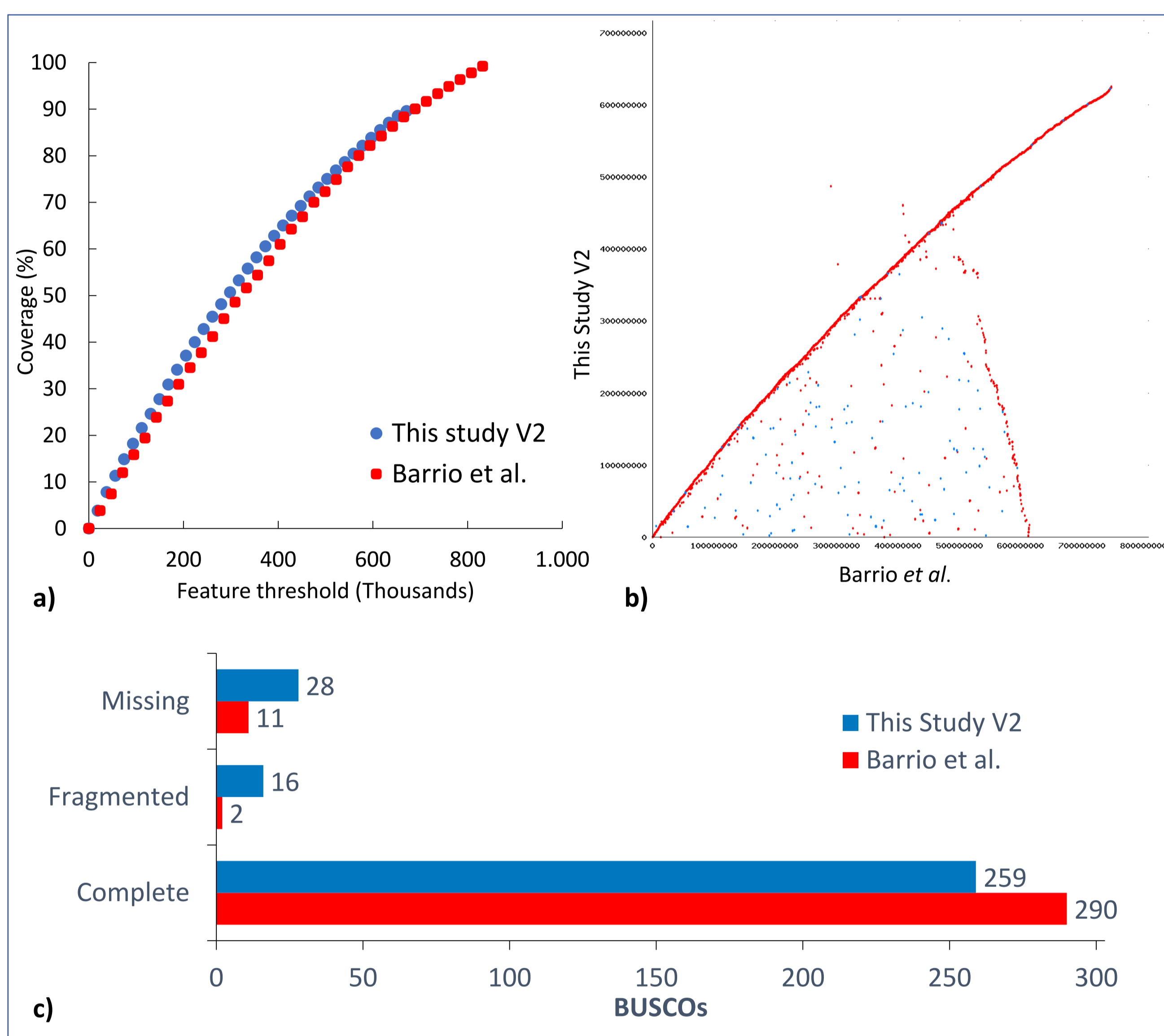
The V2 assembly has less completed BUSCOs and more fragmented and missing BUSCOs. This suggests that the Barrio *et al.* assembly is more complete than the V2 assembly. This makes sense as the Barrio *et al.* assembly has a longer total length and shorter gaps. The missing and fragmented BUSCOs in the V2 assembly could be in these gaps.

### REAPR:

REAPR reported less errors and warnings and more error free bases for the V2 assembly, compared to Barrio *et al.* assembly. REAPR outputs broken assemblies that are broken based on the assembly errors called. The broken V2 assembly was less fragmented than the broken Barrio *et al.* assembly. This suggests that even though the original V2 assembly is more fragmented than the published assembly, when assembly errors are taken into account, it has better summary statistics.

### Conclusion:

Our V2 assembly of the herring genome has considerably lower N50 value than that of Barrio *et al.*, but several assembly parameters suggest that our assembly is of higher quality. Thus, the results show that our new assembly is a good representation of the Atlantic herring genome. The V3 assembly has even better summary statistics than V2 but the other parameters also need to be evaluated.



**Figure 1.** Comparison of Barrio *et al.* herring assembly and assembly V2 from this study **a)** Comparison of features in assembly with a Feature-Response curve (FRC). The FRC was generated using FRC\_align<sup>7</sup> software. **b)** Structural variation between assemblies shown as a dotplot generated with MUMmer3<sup>8</sup>. **c)** Assembly completeness assessed using the BUSCOv3<sup>9</sup> software. Bar chart shows the number of complete, fragmented and missing Benchmarking Universal Single-Copy Orthologs (BUSCOs) in each assembly.

**Table 2.** Assembly errors and warnings called by REAPR<sup>10</sup>. Errors are calculated by coverage and fragment coverage distribution (FCD). Subgroups for warnings are reported by REAPR but are not included here for simplicity.

Errors and warnings	This Study V2	Barrio <i>et al.</i>
Error free bases	59.45%	11.42%
Errors	6,085	6,719
FCD errors within a contig	558	1,637
FCD errors over a gap	5,104	5,068
Low fragment coverage within a contig	89	6
Low fragment coverage over a gap	334	8
Warnings	509,349	2,160,429

**Table 3.** Summary statistics for broken assemblies generated by REAPR based on errors in table 2.

Assembly	No. Scaffolds	N50	L50	Assembly length
Barrio <i>et al.</i>	50,731	45,390	3,993	792,634,161
This study V2	23,172	105,152	1,784	706,219,224

## References

- Hagstovan [http://statbank.hagstova.fo/pxweb/fo/H2/H2\\_UH/](http://statbank.hagstova.fo/pxweb/fo/H2/H2_UH/)
- Nielsen, E.E., *et al.* *Nature communications* 3 (2012): 851.
- Barrio, A.M., *et al.* *Elife* 5 (2016): e12081.
- Butler, J., *et al.* *Genome Research* 18 (2008):810–20.
- Boetzer, M. and Pirovano, W. *BMC Bioinformatics* (2014)15:211.
- Boetzer, M. and Pirovano, W. *Genome biology* 13.6(2012): R56.
- Vezi, F. [https://github.com/vezi/FRC\\_align](https://github.com/vezi/FRC_align).
- Kurtz, S., *et al.* *Genome biology* 5.2 (2004): R12.
- Simão, F.A., *et al.* *Bioinformatics* 31.19 (2015): 3210–3212.
- Hunt, M., *et al.* *Genome biology* 14.5 (2013): R47.

## Method outline

### Data generation

NextSeq500      MinION

### QC

Adapter trimming (Trimmomatic)  
Remove polyG reads (AfterQC)  
Sort Mate pair reads (NextClip)

Align to reference (bwa)  
Filter mapped reads (samtools)

### Assembly

Allpaths-LG  
SPAdes  
SGA

SPAdes

### Scaffolding

SSPACE-LongRead

### Comparison

Summary statistics (N50, scaffold no. etc.)  
Feature-Response curve (FRC) and Dot plot  
BUSCO and REAPR