

Additional file

## **An improved genome assembly uncovers prolific tandem repeats in Atlantic cod**

**Ole K. Tørresen<sup>1†</sup>, Bastiaan Star<sup>1</sup>, Sissel Jentoft<sup>1,2</sup>, William Brynildsen Reinart<sup>1</sup>, Harald Grove<sup>3</sup>, Jason R. Miller<sup>4</sup>, Brian P. Walenz<sup>5</sup>, James Knight<sup>6</sup>, Jenny M. Ekholm<sup>7</sup>, Paul Peluso<sup>7</sup>, Rolf B. Edvardsen<sup>8</sup>, Ave Tooming-Klundrerud<sup>1</sup>, Morten Skage<sup>1</sup>, Sigbjørn Lien<sup>3</sup>, Kjetill S. Jakobsen<sup>1</sup> and Alexander J. Nederbragt<sup>1,9†</sup>**

### **Addresses:**

<sup>1</sup>Centre for Ecological and Evolutionary Synthesis (CEES), Department of Biosciences, University of Oslo, NO-0316 Oslo, Norway.

<sup>2</sup>Department of Natural Sciences, University of Agder, NO-4604 Kristiansand, Norway.

<sup>3</sup> Centre for Integrative Genetics (CIGENE), Department of Animal and Aquacultural Sciences, Norwegian University of Life Sciences, NO-1432, Ås, Norway.

<sup>4</sup> J. Craig Venter Institute, 9704 Medical Center Drive, 20850, Rockville, MD, USA.

<sup>5</sup> Genome Informatics Section, Computational and Statistical Genomics Branch, National Human Genome Research Institute, National Institutes of Health, 20892, Bethesda, MD, USA.

<sup>6</sup> Yale School of Medicine, Yale University, 06520, New Haven, CT, USA.

<sup>7</sup> Pacific Biosciences, Menlo Park, CA, USA.

<sup>8</sup> Institute of Marine Research, Nordnes, NO-5817, Bergen, Norway.

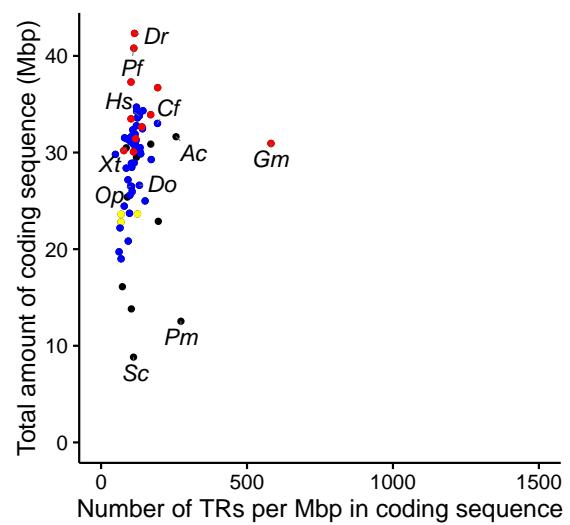
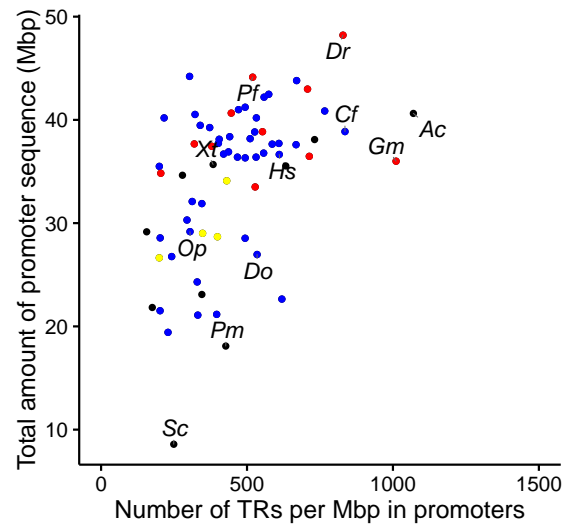
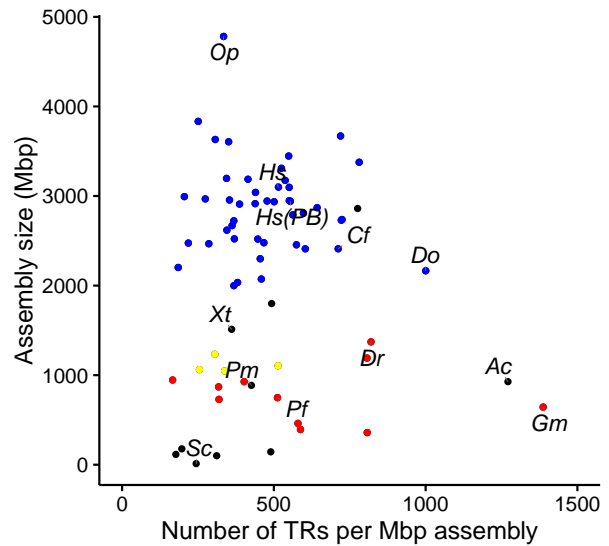
<sup>9</sup> Biomedical Informatics Research Group, Department of Informatics, University of Oslo, NO-0316 Oslo, Norway

<sup>†</sup>Corresponding authors: E-mail: o.k.torresen@ibv.uio.no, lex.nederbragt@ibv.uio.no

## 1 Supplementary Figures

**Additional file 1: Figure S1:** The frequency of tandem repeats in genome assembly, promoters and coding regions.

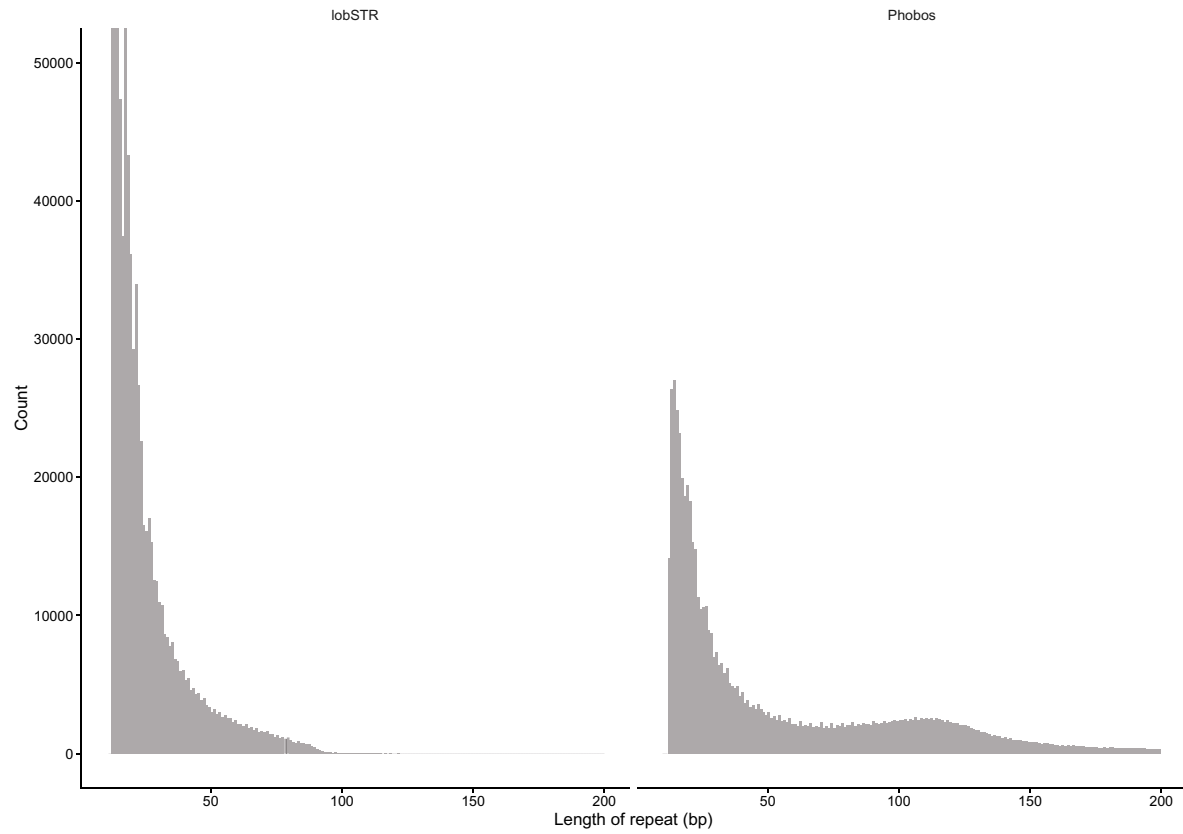
As Figure 5 in the manuscript, but plotting the number of TRs detected per Mbp assembly instead (frequency).





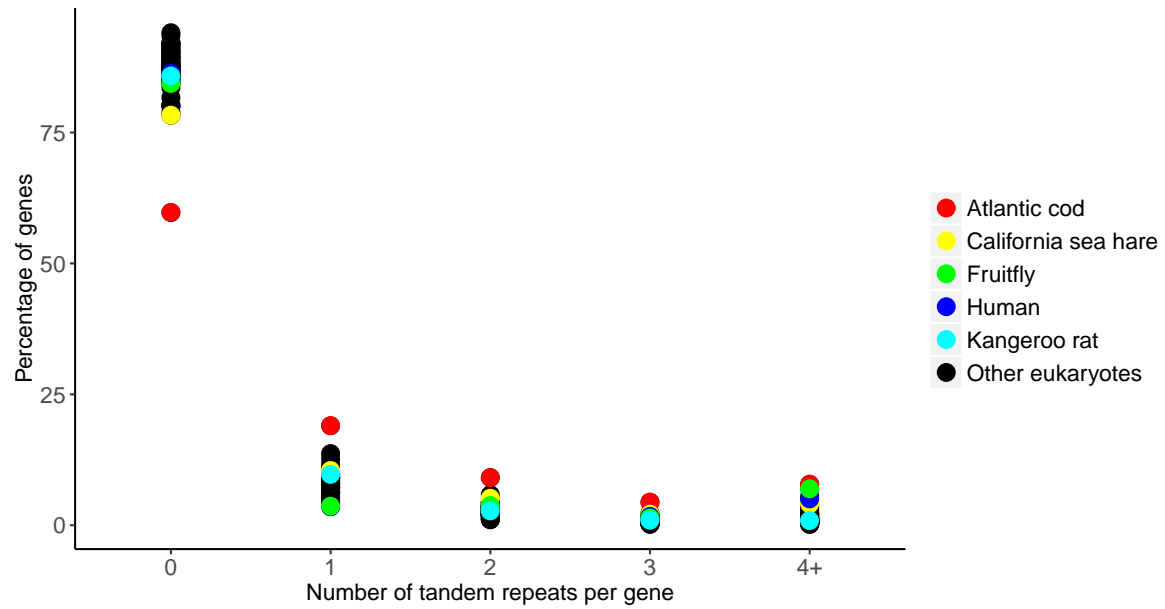
**Additional file 1: Figure S3:** The distribution of lengths of STRs in cod as found by lobSTR and Phobos.

Only repeats of unit size 1-6 bp and of total length longer than 13 bp are included.



**Additional file 1: Figure S4:** Tandem repeats in genes.

Percentage of genes (vertical axis) with a certain number of tandem repeats (i.e.; 0, 1, 2, 4 or more individual repeats within the genes, horizontal axis) in selected species.



## 2 Supplementary Tables

Additional file 1: Table S1: Read datasets, accession numbers and amount.

Technology	Insert size (bp)	Average read length (bp)	Amount bases (Gbp)	Number of pairs	Total coverage	APLILM NEWB454	CA454ILMCA454PB
Illumina	180	100	97	485,469,807	105x	52x	25x
Illumina	300	100	116	580,675,602	140x	140x	6x
Illumina	5000	100	103	513,197,070	124x	124x	124x
454	1000	170	0.4	1,185,540	0.6x	0.6x	0.6x
454	1400	170	0.3	939,375	0.5x	0.5x	0.5x
454	1800	175	0.4	1,174,856	0.6x	0.6x	0.6x
454	2300	160	0.5	1,424,229	0.7x	0.7x	0.7x
454	3000	170	1.2	3,627,219	1.9x	1.9x	1.9x
454	8000	175	1.4	3,876,715	2.1x	2.1x	2.1x
454	20000	200	0.3	818,578	0.5x	0.5x	0.5x
454	Not paired	340	23.6	NA	36.3x	36.3x	36.3x
PacBio (C2C2)	Not paired	2400	7.1	NA	11.0x		11.0x
PacBio (C2XL)	Not paired	3500	2.3	NA	3.5x		3.5x
PacBio (XLXL)	Not paired	3800	2.8	NA	4.4x		4.4x
Sanger	100000	900	0.07	39,017	0.1x		0.1x

**Additional file 1: Table S2:** Overview of assembly statistics. CEGMA annotates 458 highly conserved eukaryotic genes, REAPR analyses the discordance between the expected order, orientation and distance of mapped paired reads, with  $FRC^{bam}$  using a similar approach. Assemblies chosen for reconciliation in bold.

Assembly	Total size assembly (Mbp)	N50 contig (kbp)	N50 scaffold (Mbp)	Percentage gap bases	CEGMA	REAPR <sup>1</sup>	$FRC^{bam2}$	Potential conflict (sequences) <sup>3</sup>
ALPILM	660	4.4	0.16	28.7	424 (92.6 %)	19,787	2,182,096	122
+ Pilon	660	4.5	0.16	28.5	427 (93.2 %)	18,668	2,171,880	123
+ PBJelly	620	8.3	0.16	9.7	431 (94.1 %)	23,994	1,878,873	134
<b>+ PBJelly + Pilon</b>	<b>620</b>	<b>8.5</b>	<b>0.16</b>	<b>9.6</b>	<b>431 (94.1 %)</b>	<b>24,066</b>	<b>1,828,800</b>	<b>134</b>
NEWB454	656	6.2	1.30	24.4	435 (95.0 %)	18,117	2,044,008	26
+ Pilon	656	6.6	1.30	24.0	430 (93.9 %)	15,917	2,018,862	19
+ PBJelly	646	10.2	1.30	15.4	437 (95.4 %)	16,930	1,875,518	28
<b>+ PBJelly + Pilon</b>	<b>645</b>	<b>10.4</b>	<b>1.30</b>	<b>15.1</b>	<b>437 (95.4 %)</b>	<b>17,534</b>	<b>1,822,739</b>	<b>28</b>
CA454ILM	647	9.9	0.50	3.5	447 (97.5 %)	7,406	1,351,500	96
<b>+ Pilon</b>	<b>648</b>	<b>10.2</b>	<b>0.50</b>	<b>3.4</b>	<b>444 (97.0 %)</b>	<b>7,025</b>	<b>1,339,572</b>	<b>83</b>
+ PBJelly	672	15.3	0.52	2.5	447 (97.5 %)	14,755	1,449,619	98
+ PBJelly + Pilon	673	15.6	0.52	2.5	444 (97.0 %)	14,750	1,438,035	92
CA454PB	682	95	0.27	1.62	431 (97.6 %)	8,617	1,508,054	188
+ Pilon	683	95	0.27	1.6	441 (96.3 %)	7,754	1,426,588	163
+ PBJelly	687	96	0.27	1.1	436 (95.2 %)	8,565	1,502,582	163
<b>+ PBJelly + Pilon</b>	<b>684</b>	<b>97</b>	<b>0.27</b>	<b>1.1</b>	<b>439 (95.6 %)</b>	<b>9,043</b>	<b>1,418,020</b>	<b>165</b>

<sup>1</sup> detected potential errors, fewer is better

<sup>2</sup> total number of features (i.e., potential assembly problems), fewer is better

<sup>3</sup> number of sequences mapping to more than one linkage group or to multiple linkage groups, fewer is better

**Additional file 1: Table S3:** Linkage groups and their sizes.

Linkage group	Size (bp)
1	28,303,952
2	24,054,406
3	29,451,055
4	34,805,322
5	24,074,055
6	25,464,620
7	31,232,877
8	26,796,886
9	25,382,314
10	25,304,306
11	28,942,968
12	27,297,974
13	25,676,735
14	29,296,932
15	26,597,959
16	31,093,243
17	19,149,207
18	22,554,255
19	21,176,260
20	24,149,133
21	22,510,304
22	21,735,703
23	23,264,654
Unplaced	46,128,564



**Additional file 1: Table S4:** Calculating of genome size using odd-sized kmers from 17 through 31 with SGA PreQC.

k	Estimated genome size (bp)
17	633,173,903
19	617,492,869
21	615,747,892
23	621,292,036
25	612,150,017
27	606,607,539
29	601,318,671
31	597,207,477