# SCIENTIFIC REPORTS

Received: 16 March 2017 Accepted: 28 June 2017 Published online: 03 August 2017

## **OPEN** Rapid *de novo* assembly of the European eel genome from nanopore sequencing reads

Hans J. Jansen<sup>1</sup>, Michael Liem<sup>2</sup>, Susanne A. Jong-Raadsen<sup>1</sup>, Sylvie Dufour<sup>3</sup>, Finn-Arne Weltzien <sup>[]</sup>, William Swinkels<sup>5</sup>, Alex Koelewijn<sup>5</sup>, Arjan P. Palstra<sup>6</sup>, Bernd Pelster<sup>7</sup>, Herman P. Spaink<sup>2</sup>, Guido E. van den Thillart<sup>1</sup>, Ron P. Dirks<sup>1</sup> & Christiaan V. Henkel<sup>2,8,9</sup>

We have sequenced the genome of the endangered European eel using the MinION by Oxford Nanopore, and assembled these data using a novel algorithm specifically designed for large eukaryotic genomes. For this 860 Mbp genome, the entire computational process takes two days on a single CPU. The resulting genome assembly significantly improves on a previous draft based on short reads only, both in terms of contiguity (N50 1.2 Mbp) and structural guality. This combination of affordable nanopore sequencing and light weight assembly promises to make high-guality genomic resources accessible for many non-model plants and animals.

Just ten years ago, having one's genome sequenced was the privilege of a handful of humans and model organisms. Spectacular improvements in high-throughput technology have since made personal genome sequencing a reality and prokaryotic genome sequencing routine. In addition, sequencing the larger genomes of non-model eukaryotes has opened up a wealth of information for plant and animal breeding, conservation, and fundamental research.

As an example, we and others<sup>1-3</sup> have previously established genomic resources for the European eel (Anguilla *anguilla*), an iconic yet endangered fish species that remains resistant to efficient farming in aquaculture<sup>4, 5</sup>. A draft genome<sup>2</sup>, several transcriptomes<sup>1, 3-10</sup>, and reduced representation genome sequencing<sup>11</sup> have already shed light on its evolution and developmental biology<sup>2, 12, 13</sup>, endocrinological control of maturation<sup>7,9</sup>, metabolism<sup>14</sup>, disease mechanisms<sup>10</sup>, and population structure<sup>15, 16</sup>, thereby supporting both breeding and conservation efforts. However, compared to established model organisms, funds for eel genomics are naturally limited, and consequently the quality of current genome assemblies of Anguilla species is modest at best by today's standards (Table 1).

The recent availability of affordable long-read sequencing technology<sup>17, 18</sup> by Oxford Nanopore Technologies (ONT) presents excellent opportunities for generating high-quality genome assemblies for any organism<sup>19</sup>. Flow cells for the miniature MinION sequencing device employ a maximum of 512 nanopores concurrently for reading single-stranded DNA at up to 450 nucleotides per second, resulting in several gigabases of sequence during a two day run. As the technology does not rely on PCR or discrete strand synthesis events, DNA fragments can be of arbitrarily long length. The single-molecule reads are of increasingly good quality, with a sequence identity of ~75% for the older R7.3 chemistry<sup>17</sup>, to ~89% for the newer R9 chemistry (MinION Analysis and Reference Consortium, in preparation). Optionally, DNA can be read twice (along both strands) to yield a consensus '2D' read of higher accuracy (up to ~94% for R9).

Long-read sequencing technology is also offered by Pacific Biosciences (PacBio). This platform employs advanced optics to detect a polymerase operating on single DNA molecules, and has been commercially available

<sup>1</sup>ZF-screens B.V., Leiden, The Netherlands. <sup>2</sup>Institute of Biology, Leiden University, Leiden, The Netherlands. <sup>3</sup>Muséum National d'Histoire Naturelle, Sorbonne Universités, Research Unit BOREA, Biology of Aquatic Organisms and Ecosystems, CNRS, IRD, UCN, UA, Paris, France. <sup>4</sup>Norwegian University of Life Sciences, Faculty of Veterinary Medicine, Department of Basic Science and Aquatic Medicine, Oslo, Norway. <sup>5</sup>DUPAN Foundation, Wageningen, The Netherlands. <sup>6</sup>Animal Breeding and Genomics Centre, Wageningen Livestock Research, Wageningen University & Research, Wageningen, The Netherlands. <sup>7</sup>Institute of Zoology and Center for Molecular Biosciences, University of Innsbruck, Innsbruck, Austria. <sup>8</sup>University of Applied Sciences Leiden, Leiden, The Netherlands. <sup>9</sup>Generade Centre of Expertise in Genomics, Leiden, The Netherlands. Correspondence and requests for materials should be addressed to C.V.H. (email: c.v.henkel@biology.leidenuniv.nl)

Species	Reference	NCBI WGS reference	Assembly methods	Contigs/ scaffolds sum	Contig/scaffold N50	Scaffold gaps
A. anguilla	2	AZBK01	CLC bio+SSPACE	969/923 Mbp*	1.7/77.6 kbp	134 Mbp
A. japonica	34	AVPY01	CLC bio+SSPACE	1.13/1.15 Gbp*	3.3/52.8 kbp	127 Mbp
A. rostrata	37	LTYT01	Ray + SSPACE	1.19/1.41 Gbp	7.4/86.6 kbp	223 Mbp

 Table 1. Previous genome assemblies of Anguilla species. \*Not all contigs obtained by de novo assembly were used in scaffold construction.

Species	Haploid genome size*	Repetitive fraction*	Heterozygous fraction*
A. anguilla	854.0-866.5 Mbp	15.5-20.0%	1.48-1.59%
A. japonica**	1.022 Gbp	38.7%	2.74%
A. rostrata	799.0-813.0 Mbp	12.2-16.9%	1.50-1.60%

**Table 2.** Anguilla genome size predictions. \*Ranges are the minimum and maximum values reported for three model fits at different *k*-mer lengths. Apparent repetitive sequence decreases with *k*-mer length, and heterozygosity increases with *k*-mer length. \*\*For *A. japonica*, the model did not converge in most cases, presumably because of low coverage. These results are for k = 19.

since 2011. Both long-read technologies deliver roughly comparable quality and data volumes. PacBio sequencing has the advantages of an established, stable platform (which includes bioinformatics), as well as less bias in the error profile. Advantages of ONT include the much lower equipment cost, and currently rapidly improving quality, read length and throughput. Comprehensive comparisons of both technologies are scarce<sup>20</sup>.

In contrast to short reads, long reads offer the possibility to span repetitive or otherwise difficult regions in the genome, resulting in strongly reduced fragmentation of the assemblies. This potential advantage does require the deployment of dedicated genome assembly algorithms that are aware of long-read characteristics. In addition, as single-molecule long-read technologies (by both PacBio and ONT) do suffer from reduced sequence identity, this likewise needs to be addressed by post-sequencing bioinformatics<sup>21–23</sup>. Dealing with these challenges has reinvigorated research into genome assembly methodology, resulting in several novel strategies<sup>24–28</sup>.

However, when dealing with large eukaryotic genomes, the computational demands for long-read assembly are often higher than for short reads (using De Bruijn-graphs), even though the raw data are more informative of genome structure. Especially now that sequencing very large plant and animal genomes is finally becoming both technologically feasible and affordable, the computational costs may turn out to be prohibitive. For example, using the state-of-the-art Canu assembly software<sup>25</sup>, assembling a human genome from long reads takes tens of thousands of CPU hours, or several days on a computer cluster (https://genomeinformatics.github.io/NA12878-nanopore-assembly). As scaling behaviour is approximately quadratic with genome size, assembling a salamander<sup>29</sup> or lungfish<sup>30</sup> genome dozens of gigabases long would require several years on a cluster.

We are currently developing a computational pipeline specifically intended for future sequencing of extremely large tulip genomes<sup>31</sup> (up to 35 Gbp). Named TULIP (for *The Uncorrected Long-read Integration Process*), its primary purpose is to split up such large assembly problems into manageable subsets of long reads. Each subset can then be handled by a separate downstream *de novo* assembly process, in theory substituting quadratic scaling with nearly linear behaviour. Here, we use a prototype of this algorithm to assemble a new version of the European eel genome, based on Oxford Nanopore sequencing. The entire computational procedure takes two days on a desktop computer, and yields an assembly that is two orders of magnitude less fragmented than the previous Illumina-based draft.

#### Results

**Eel genome sizes and previous assemblies.** Before launching a genome sequencing effort, an estimate of the size of the genome of interest is needed. For the genus *Anguilla*, several studies have used flow cytometry and other methods to arrive at C-values ranging from 1.01 to 1.67 pg (http://www.genomesize.com), corresponding to haploid genome sizes in the 1–1.6 Gbp range for both *A. anguilla* and *A. rostrata*. We previously estimated a genome size of approximately 1 Gbp for *A. anguilla*, using human cells as a reference<sup>2</sup>.

Based on their assembled genomes, *Anguilla* species exhibit a similarly wide range of apparent genome sizes (see Table 1). These draft assemblies are all based on previous-generation short-read technology, and relied on Illumina mate pairs to supply long-range information used in scaffolding. The resulting assemblies remain highly fragmented, with low N50 values even considering the technology used.

We therefore examined *k*-mer profiles in the raw Illumina sequencing data, which can provide an estimate of the length of the haploid genome<sup>32, 33</sup>. Surprisingly, the predicted genome sizes are considerably – but consistently – smaller than previously estimated or assembled (Table 2 and Supplementary Fig. S1). In addition, all three examined genomes contain high levels of heterozygosity.

**Nanopore sequencing.** We isolated DNA for long-read sequencing from the blood and liver of a fresh female European eel. Using three different generations of the ONT chemistry for the MinION sequencer, we





generated 15.6 Gbp of raw shotgun genome sequencing data (see Fig. 1 and Supplementary Table S1). Assuming an 860 Mbp haploid size, this corresponds to approximately 18-fold coverage of the genome. The bulk of the

an 860 Mbp haploid size, this corresponds to approximately 18-fold coverage of the genome. The bulk of the sequence is in long or very long reads (up to hundreds of thousands of nucleotides), although a fraction is composed of very short reads or artifacts (e.g. 6 bp reads, Fig. 1). We used all raw reads for subsequent genome assembly.

**Assembly strategy.** We assembled the long nanopore sequencing reads using a prototype of an assembly strategy we are developing for very large genomes (M. Liem and C. Henkel, in preparation), named TULIP. Briefly, it takes two shortcuts compared to the established hierarchical approach<sup>21, 25</sup>. First of all, like Miniasm<sup>27</sup>, TULIP does not correct noisy single-molecule reads prior to assembly. Secondly, it does not perform an all-versus-all alignment of reads, but instead aligns reads to a sparse reference (of 'seed' sequences) that is representative for the genome. The result is a 'seed graph', which can be used to either partition the original long reads into many independent subsets for subsequent *de novo* assembly, or to immediately extract uncorrected scaffold sequences from. Here, we have chosen to use the latter functionality, and employed stand-alone post-assembly consensus applications to correct the resulting scaffolds.

Figure 2a illustrates all the steps we have taken during *de novo* assembly of the European eel genome. We employed previously generated Illumina shotgun sequencing reads as sparse seeds. Using a *k*-mer counting table, we identified merged read pairs that are suitably unique in the genome. Using strict criteria (see Methods), we could select 5019778 fragments of 270 bp, or 873058 of 285 bp, corresponding to 1.58-fold or 0.29-fold coverage of the genome, respectively. We subsequently used several random subsets of these fragments as a reference to align long nanopore reads against.

Using a custom script, we constructed a graph based on these alignments, in which the seed sequences are nodes, and edges represent long read fragments (Fig. 2b). A connection between two seeds indicates they co-align to a long read, and are therefore presumably located in close proximity in the genome. In theory, perfect alignments of very long reads to unique seeds should be sufficient to organize both sets of data into linear scaffolds.



**Figure 2.** Assembly strategy. (a) Stages in the TULIP assembly of the European eel genome. (b) Graph construction based on long read alignments to short seeds. Seeds are included in the graph as nodes if they align adjacent to each other to a long read. The apparent distance between the seeds is included as an edge property, as is the amount of evidence (i.e. number of alignments supporting the connection). (c) The initial seed graph based on alignments contains ambiguities, caused by missed alignments, repetitive seed sequences and spurious alignments. These are removed during the initial layout process, resulting in linear scaffolds. Where possible, these scaffolds are subsequently linked by further unambiguous long-distance co-alignments to long reads.

However, because of the errors still present in long nanopore reads, the alignments are imperfect, with missed seed alignments making up the bulk of ambiguities in the seed graph (i.e. forks and joins in the seed path). Additional uncertainties are introduced by spurious alignments and residual apparently repetitive seeds. The tangles these cause in the graph can be recognized locally, and are removed during a graph simplification stage (Fig. 2c). TULIP will visit every seed that has multiple in- or outgoing connections, and attempt to simplify the local graph topology by removing connections. For example, if a single seeds fails to align to a single nanopore read, this will introduce a 'triangle' in the graph (Fig. 2c, top example), in which the neighbouring seeds now share a direct connection (based on that single read). If the intermediate seed fits between the neighbouring seeds, TULIP will then remove the connection spanning the intermediate seed. If after this stage a seed still has too many connections, it might represent repetitive content and its links are severed altogether (Fig. 2c, second example).

Finally, unambiguous linear arrangements of seeds can be extracted from the graph. Figure 3 illustrates a small fragment of the actual seed graph, with final linear paths (scaffolds) and removed connections indicated. These ordered seed scaffolds do not yet contain sequence data. These can subsequently be added from the original nanopore reads and alignments, resulting in uncorrected scaffold sequences. The scaffolds are exported bundled with their constituent nanopore reads, and can be subjected to standard nanopore sequence correction procedures.

**Assembly characteristics.** We used several combinations of short seed sequences and aligned nanopore reads to optimize the assembly process. In most cases, we did not complete the entire assembly process by adding actual nanopore sequence. Therefore, distances between seeds (and scaffold lengths) are means based on multiple nanopore reads. Adding specific sequence (and subsequently correcting scaffolds) can change these figures slightly. Supplementary Table S2 lists the assembly statistics for these experimental runs.

Both the contiguity and size of the assembly clearly improve upon adding more nanopore data (Fig. 4a,b). This suggests that at 18-fold coverage of this genome, and using the particular blend of data types available here, the assembly process is still limited by the total quantity of long read data.



**Figure 3.** Graph simplifications. Scaffolds were extracted from a graph consisting of seed sequences (nodes) linked by nanopore reads (edges). Here, a small final scaffold (number 2231, 252.2 kbp) is shown in red in the context of the initial seed graph (all seeds at a distance of up to ten links from the final scaffold). Fragments of ten other scaffolds (blues) are directly or indirectly connected to scaffold 2231 by a few incorrect links (dotted lines). Seeds and links removed during graph simplification are shown in grey. Scaffolds can be discontinuous in the initial graph, as additional long-distance links are added in a later stage. The graph was visualized using Cytoscape (version 3.4.0).

For the seeds, we investigated the effects of seed length (270 or 285 bp), as well as seed density (fractions and multiples based on the 873058 fragments available at 285 bp). There does not appear to be a clear advantage to choosing either 270 or 285 bp seeds. At identical densities, the two possibilities yield comparable assemblies in terms of size and contiguity.

For seed density, there does appear to be an optimum. As expected, low densities result in fragmentation and incompleteness (Fig. 4c,d). The assemblies with the highest seed density (1.3 or 1.7 million 270 bp sequences) do yield the highest N50 and assembly sum, but also exhibit increased fragmentation compared to lower seed densities. As Fig. 4c shows, the main difference with those assemblies is the appearance of many small scaffolds at high seed numbers. Accidentally, in this case the optimal seed density is around the 'full' set of 873058 fragments, of either 270 or 285 bp. Both also yield an assembly that is close to the estimated genome length. We selected the 285 bp version as a candidate for an updated reference genome for the European eel.

Figure 4 summarizes several characteristics of the candidate assembly (before sequence addition or correction). The length distribution of the 2366 scaffolds (Fig. 4a) shows they range in size between 431 bp and 8.7 Mbp. The lower boundary is expected, as a minimal scaffold has to consist of at least two 285 bp seeds, and the graph construction was executed with parameters allowing limited overlap between seeds. The cumulative scaffold length distributions (Fig. 4c) show that a considerable fraction of the genome is included in large scaffolds, with 232 scaffolds larger than a megabase constituting 56% of the assembly length. Seeds in the final scaffolds are connected by on average 7.4 nanopore read alignments. As can be seen in Fig. 4e, links removed during the graph simplification stage (mostly based on local graph topology only) were predominantly those supported by less evidence.

The final assembly retains 637792 seeds of 285 bp, equivalent to a maximum of 181.8 Mbp of Illumina-derived sequence. If the seed distribution is assumed to be essentially random (with local genomic architecture responsible for exceptions), the initial 873058 seeds should be spaced at a mean interval of 700 bp. As seeds are removed during simplification, larger 'gaps' filled with nanopore-derived sequence should appear. However, as Fig. 4f shows, gap lengths are heavily biased towards low and negative lengths (i.e. overlapping seeds). In this case, this could be an artifact of the very stringent seed selection procedure.



**Figure 4.** Characteristics of the final assembly. (a) Size distribution of final scaffolds, based on 285 bp seeds. Colours indicate alternative assembly runs, using subsets of the long read data. (b) Cumulative size of the final scaffolds, sorted by size. (c) and (d) Size distributions and cumulative size distributions for final scaffolds, based on both 270 and 285 bp seeds. Colours indicate alternative assembly runs, using different seeds sets. (e) Link evidence distribution in the initial graph (purple) and the final graph (orange) for the candidate assembly (285 bp seeds). (f) Distances between seeds in the initial graph (purple) and the final graph (orange) for the candidate assembly (285 bp seeds).

**Assembly quality.** In order to assess its completeness and structural correctness, we added nanopore sequence to the selected TULIP assembly and aligned it to the Illumina-based draft genome<sup>2</sup>. As a high-quality reference genome for the European eel is not yet available, such a comparison need take into account the possibility of error in either assembly. However, with appropriate caution, agreement between the assemblies – which are completely independent in both sequencing data and assembly algorithms – can confirm the integrity of both.

Figure 5a shows a full-genome alignment of the new (uncorrected) nanopore-based assembly to the 2012 draft<sup>2</sup>, based on best pairwise matches. This confirms that at this large scale, all sequence in the new assembly is also present in the older assembly. At first sight, the converse does not appear to be the case: the Illumina-based draft is 923 Mbp in size, and contains approximately 96 Mbp in scaffolds that have no reciprocal best match in the nanopore assembly (863.3 Mbp after sequence addition, see Supplementary Table S3). However, the non-matching sequences consist almost exclusively of very small scaffolds (mean/N50 664/987 bp). Since the Illumina-based draft assembly also contains 134 Mbp in gaps, these small scaffolds are plausibly sequences that could not be integrated correctly during the SSPACE scaffolding process<sup>34, 35</sup>. Both assemblies therefore roughly span the entire predicted genome of 860 Mbp.

Figure 5b-f show detailed alignments, based on the 5 largest nanopore scaffolds (6.1–8.9 Mbp uncorrected) and their best matches only. These alignments confirm that in this sample both assemblies are mostly collinear, with the smaller Illumina draft scaffolds usually aligning end-to-end on the larger TULIP scaffolds. Therefore, both presumably reflect the actual genomic organization. However, at this level of detail several structural incongruities between both assemblies also become apparent (indicated by arrowheads). For 16 scaffolds from the 2012 draft, only part of the sequence is present in the selected TULIP scaffolds. In other words, at these loci both assembly protocols made different choices, based on the available sequencing information.

We therefore examined the evidence for the decisions made by TULIP. For each discrepancy, we examined the local neighbourhoods in the initial nanopore-based seed graphs (as in Fig. 3). If a draft scaffold is correct, at the inconsistency there should be multiple alternatives for the TULIP algorithm to choose from (Supplementary Fig. S2). As these subgraphs (Supplementary Figs S3–S7) show, there is no evidence in the nanopore data for the older draft structure for any of the 16 cases examined. On the contrary, most local graph neighbourhoods appear relatively simple and support unambiguous scaffolding paths. The links at these suspect junctions are supported by at least two (average six) independent nanopore reads, which reduces the likelihood of accidental connections (caused by e.g. chimeric reads).

Alternatively, the order of the draft scaffolds in the alignments already suggests which of the two assemblies is correct. If one of the 16 problematic scaffolds were to reflect the legitimate genome structure, this error in the new assembly would usually also affect the next aligning scaffold. However, in almost all cases, the neighbouring draft scaffold aligns end-to-end. This suggests that either the TULIP assembly intermittently features very large rearrangements that accidentally always end at draft scaffold boundaries, or that the draft scaffolds are occasion-ally misconstrued.

The distribution of draft scaffolds along the nanopore-based scaffolds reveals an interesting pattern. The distribution of draft scaffold length along the genome is clearly non-random, with some regions assembled into just a few large scaffolds, whereas other regions (often up to a Mbp in size) are highly fragmented into very small scaffolds. This indicates that using short-read technology, certain genomic features are intrinsically harder to assemble than using long reads.

Finally, we assessed the completeness of the nanopore assembly using BUSCO<sup>36</sup>. This method assumes complete assemblies to contain a high fraction of genes that are highly conserved in related species. From a set of 2586 common vertebrate genes, BUSCO was only able to recover 78 complete and 106 fragmented genes (3.0% and 4.1%, respectively). 92.9% of orthologues are missing from the nanopore assembly, indicating very poor completeness. In this case, however, this is a result of the sequence characteristics of ONT data.

**Sequence correction.** Currently, the ONT platform does not yield reads of perfect sequence identity. Like with PacBio data, therefore, at some point in the assembly process the single-molecule-derived sequence needs to be corrected by extracting a consensus from multiple reads covering every genomic position. Here, we opted for a standalone post-assembly correction step with Racon, which extracts a consensus from nanopore reads<sup>23</sup>. As some positions in the assembly are based on a single nanopore read (Fig. 4e), in this case this correction may not be sufficient. Therefore, we subsequently corrected with Pilon, which extracts a consensus based on alignment of Illumina reads to the noisy sequence<sup>37, 38</sup>.

To assess the changes made by these correction algorithms, we counted and compared the occurrence of 6-mers in the draft Illumina-based assembly, the uncorrected TULIP assembly, and after correction (Fig. 6). These frequencies reveal several expected patterns<sup>17</sup>, specifically a slight underrepresentation of high CG content in Illumina-based sequence (draft and Pilon), and an underrepresentation of homopolymer sequence in nanopore-based sequence (TULIP and Racon). Overall, the correction steps bring the sequence similarity of the nanopore-based assembly closer to the Illumina-based draft, with the final corrected assembly having a high correlation to the draft (Fig. 6 lower left panel).

Sequence correction also has a strong positive impact on the BUSCO completeness assessment. As BUSCO relies on the prediction of gene structures, small artefactual deletions and insertions might cause it to miss genes. After correction with Racon, the BUSCO scores increased to 10.8% complete, 21.6% fragmented and 67.6% missing; correction with Pilon resulted in a further increase to 77.5% complete, 14.1% fragmented and 8.4% missing. An additional round of Pilon polishing resulted in a BUSCO assessment of 79.8% complete, 12.9% fragmented and 7.3% missing.

Sequence correction remains the most time-consuming stage of the assembly process, requiring 22 and 24 hours (on a single CPU) for Racon and Pilon, respectively (Supplementary Table S3). As TULIP bundles uncorrected scaffolds with its constituent nanopore reads, this process could still be sped up by parallelization, with individual scaffolds distributed over concurrent correction threads.

#### Discussion

In this study, we have evaluated whether it is possible to sequence a vertebrate genome using Oxford Nanopore long-read technology, and quickly assemble it by means of a relatively simple and lightweight procedure. Using



**Figure 5.** Full-genome alignment of the final assembly. (**a**) The final uncorrected scaffolds (N50 = 1.19 Mbp, y-axis) were aligned to the 2012 *A. anguilla* assembly (N50 = 77.6 kbp, x-axis) using nucmer<sup>51</sup> with minimum match length 100, filtered for best pairwise matches between scaffolds (*delta-filter -1*), and plotted using the mummerplot *--layout* option. The grey area corresponds to small scaffolds in the 2012 assembly that are not part of a best reciprocal match. (**b**-**f**) More detailed alignments between the five largest nanopore scaffolds (*y*-axes) and their best matches in the 2012 draft assembly (x-axes). Grey vertical lines indicate scaffold boundaries. These figures were generated in R (version 3.3.1) based on mummerplot output. 2012 draft scaffolds with minimal contributions to the overall alignment were removed manually. Arrowheads indicate discrepancies between both assemblies.

.....

our original TULIP methodology, we were able to assemble the 860 Mbp genome of the European eel using 18-fold nanopore coverage and sparse pre-selected Illumina reads in three and a half hours on a modest desktop computer. Including subsequent sequence correction, the entire process takes two days. This yields an assembly that is essentially complete and of high structural quality (Fig. 5).



**Figure 6.** Sequence identity in nanopore-based assemblies. The sequence similarity to the older draft of different stages of the nanopore assembly process (uncorrected TULIP, corrected by Racon<sup>23</sup>, and additionally corrected by Pilon<sup>37, 38</sup>) is illustrated by 6-mer frequency counts (generated using Jellyfish<sup>46</sup>). With every point a discrete 6-mer, colours indicate CG-content, and open circles indicate the two homo-6-mers. Scales are logarithmic. Also shown are Pearson correlation coefficients between the frequency distributions.

One of the most striking outcomes of this eel genome sequencing effort is the close match between the genome size predicted from k-mer analysis (~860 Mbp) and the TULIP assembly (891.7 Mbp after corrections), and their distance from short-read-based assemblies. This can be explained either by the absence of a substantial fraction of the genome from the nanopore data or assembly, or by an artificially inflated genome size for the short-read assemblies. Full-genome alignment between both assemblies (Fig. 5a) suggests the latter phenomenon is at least partially responsible, as only tiny short-read scaffolds are absent from the long-read assembly. Furthermore, BUSCO analyses indicate the new assembly is approximately complete.

An analysis of the short-read *A. anguilla*<sup>2</sup> and *A. japonica*<sup>35</sup> assembly procedures implies that the scaffolding process, based on mate pair data, is responsible for the introduction of numerous gaps (Table 1). In addition, at the time we discarded a considerable fraction of the initial contigs, which was composed primarily of very small contigs that appeared to be artefactual (based on low read coverage or very high similarity to other contigs). Plausibly, such contigs – and the high residual fragmentation of these assemblies – are the result of the high levels of heterozygosity in these genomes (Supplementary Fig. S1).

Similar processes could also explain the even larger discrepancy between the predicted and assembled size of the recently published genome<sup>39</sup> of the American eel *A. rostrata* (Table 1). As European and American eels

interbreed in the wild<sup>40</sup>, a large difference in genome size is unlikely – although it could also provide an explanation for the observed limited levels of gene flow between the species<sup>15</sup>.

The whole-genome alignments between the Illumina draft and the new nanopore-based assembly (Fig. 5) also serve to confirm the structural accuracy of both. In a representative sample (corresponding to of 4.2% of the genome), we observed 16 apparent assembly errors (Fig. 5b–f). In the absence of a high-quality reference, it is not straightforward to establish which assembly is correct. Our analyses, however, strongly suggest that in these cases the nanopore-based assembly is accurate. This is not unexpected: TULIP has access to far richer and more precise sequencing information than SSPACE, which had to rely on  $2 \times 36$  bp mate pair data. Under such circumstances, a low number of incorrect joins between contigs is inevitable<sup>41</sup>. In fact, considering the fact that the SSPACE scaffolds analyzed in Fig. 5b–f consist of on the order of ten thousand very small contigs, a result with only 16 errors signifies better scaffolding performance than expected<sup>41</sup>.

In other aspects, the TULIP assembly is likely to be suboptimal. By design, scaffolds that could be merged based on long reads remain separate if these reads do not share a fortuitous seed alignment in the correct position. Similarly, large repetitive regions in the genome, as well as (sub) telomeric repeats will not always contain frequent 285 bp islands of unique sequence, and hence could be absent from the assembly. Although counterintuitive, this should not pose a major problem for some extremely large genomes. Survey sequencing indicates that the 32 Gbp axolotl genome contains mostly unique sequence<sup>29</sup>, as do many tulip genomes (C. Henkel, unpublished data).

The selection of sparse seeds by the user adds an unusual level of flexibility to the assembly process. In an early phase of this study, we opted for essentially randomly placed Illumina-based seed sequences. This choice was motivated by their very high sequencing identity, which aids alignment quality when working with noisy long reads. This strategy should work equally well with PacBio data or early, error-prone nanopore chemistries (i.e. R7.3).

The genome assembly generated here is a hybrid, incorporating two different sequencing technologies, three generations of nanopore sequencing, and two different animals. At the time, it was unavoidable to use a combination of multiple nanopore sequencing chemistries, as these rapidly replaced each other. Although the later R9 and R9.4 chemistries have better sequencing error profiles, they still retain structural biases that cannot be resolved by taking a consensus of nanopore data only (e.g. using Racon). In the final Pilon polishing stage, the nanopore data are therefore corrected using Illumina data obtained from a different eel specimen than used for nanopore sequencing. As the European eel is highly heterozygous (Table 2), in theory this generates a consensus between up to four different haplotypes. In practice, we expect this to have little influence on the quality of the final assembly, as the variation resulting from heterozygosity is much lower than the raw nanopore error rate. In other words, Pilon will treat SNPs and small indels not occurring in the Illumina data as sequencing errors to be corrected.

With the speed at which the quality of reads produced by the ONT platform is improving<sup>18</sup>, it should soon be possible to avoid a hybrid assembly incorporating short reads altogether. A natural choice for seed sequences would then be the ends of long reads. Alternatively, seeds could be chosen to facilitate further sequence integration. If a high density genetic map is available for a species, map markers could serve as pre-ordered seeds. For example, with minor modifications, TULIP might be used to selectively add long read sequencing data only to single map marker bins (containing thousands of actual, unordered markers) resulting from a population sequencing strategy<sup>42</sup>.

The bottleneck for such strategies lies in the interplay between marker density and nanopore read length, where the latter currently appears to be limited chiefly by DNA isolation protocols<sup>43, 44</sup>. Conceivably, in the near future, the problem of genome assembly from sequencing reads will all but disappear: abundant megabase-sized reads of high sequence identity are becoming possible, which should span the vast majority of recalcitrant regions in medium-sized genomes that remain a challenge to short- and medium-read technologies.

The fulfillment of such prophesies may still lie several years in the future. Therefore, we plan to further integrate and validate the candidate assembly generated here with long-range information obtained from optical mapping<sup>45</sup>, in order to develop a high-quality reference genome for the troubled European eel.

#### Methods

**Eel samples.** Two different European eels were used to generate the genome assembly. For all Illumina sequencing, a female specimen caught in Lake Veere, The Netherlands, was used. These data were previously used for the Illumina-based draft assembly<sup>2</sup>. For nanopore sequencing, a farmed female eel was obtained from Passie voor Vis, Sevenum, The Netherlands. As the European eel is a panmictic species<sup>16</sup>, these sequenced eels belong to the same population. The experiments were approved by the animal ethical commission of Leiden University (DEC #13060), and carried out in accordance with the relevant guidelines and regulations.

**Genome size estimation and** *k***-mer analyses.** We used Jellyfish<sup>46</sup> version 2.2.6 to count *k*-mers in sequencing reads and assemblies. In order to estimate genome size, we obtained frequency histograms for 19-to 25-mers in raw Illumina sequencing data. Reads were truncated to a uniform length of 76 nt, except for *A. japonica*, for which we used 100 nt (the model did not converge for short lengths). For the American eel, which has been sequenced to much higher coverage than the European and Japanese species, we used a subset of the available data (NCBI Sequence Read Archive SRR2046741 and SRR2046672). Histograms were analyzed using the GenomeScope<sup>33</sup> website in order to obtain estimates for genome sizes, heterozygosity and duplication levels.

**Illumina seed selection.** We selected unique seed sequences from 11.9 Gbp in sequence previously generated at  $2 \times 151$  nt on an Illumina Hiseq 2000 (NCBI Sequence Read Archive SRR5235521). Pairs were merged using FLASh<sup>47</sup>, requiring a minimum of 15 nt terminal overlaps, resulting in 29.16% merged fragments. In these, 25-mers were counted using Jellyfish. We used a custom script to filter out all fragments that contained 25-mers occurring over 25 times in the remaining data. This corresponds to a maximum occurrence of approximately  $6.25 \times$  in the 860 Mbp genome. Finally, fragments were selected based on size (either 270 nt or 285 nt). **MinION library preparation and sequencing.** High MW chromosomal DNA was isolated from European eel blood and liver samples using a genomic tip 100 column according to the manufacturer's instructions (Qiagen). For each nanopore sequencing library, we used  $2-3 \mu g$  genomic DNA, approximately twice the recommended quantity. In this way, we compensated for the decreased molar quantities of DNA ends at increased fragment lengths (see below).

First the DNA was sequenced on R7.3 flow cells. Subsequently multiple R9 and R9.4 flow cells were used to sequence the DNA. For R7.3 sequencing runs we prepared the library using the SQK-MAP006 kit from Oxford Nanopore Technologies. Briefly, high molecular weight DNA was sheared with a g-TUBE (Covaris) to an average fragment length of 20 kbp. The sheared DNA was repaired using the FFPE repair mix according to the manufacturer's instructions (New England Biolabs, Ipswich, USA). After cleaning up the DNA with an extraction using a ratio of 0.4:1 Ampure XP beads to DNA the DNA ends were polished and an A overhang was added with the the NEBNext End Prep Module and again cleaned up with an extraction using a ratio of 1:1 Ampure XP beads to DNA the DNA prior to ligation. The adaptor and hairpin adapter were ligated using Blunt/TA Ligase Master Mix (New England Biolabs). The final library was prepared by cleaning up the ligation mix using MyOne C1 beads (Invitrogen).

To prepare 2D libraries for R9 sequencing runs we used the SQK-NSK007 kit from Oxford Nanopore Technologies. The procedure to prepare a library with this kit is largely the same as with the SQK-MAP006 kit. 1D library preparation was done with the SQK-RAD001 kit from Oxford Nanopore Technologies. In short, high molecular weight DNA was tagmented with a transposase. The final library was prepared by ligation of the sequencing adapters to the tagmented fragments using the Blunt/TA Ligase Master Mix (New England Biolabs).

Library preparation for R9.4 sequencing runs was done with the SQK-LSK108 and the SQK-RAD002 kits from Oxford Nanopore Technologies. The procedure to prepare libraries using the SQK-RAD002 kit was the same as for the SQK-RAD001 kit. For SQK-LSK108 the procedure was essentially the same as for SQK-NSK007 except that only adapters and no hairpins were ligated to the DNA fragments. As a consequence the final purification step was done using Ampure XP beads instead of MyOne C1 beads. Libraries for R7.3 and R9 flow cells were directly loaded on the flow cells. To load the library on the R9.4 flow cell the DNA fragments were first bound to beads which were then loaded on the flow cell.

The MinKNOW software was used to control the sequencing process and the read files were uploaded to the cloud based Metrichor EPI2ME platform for base calling. Base called reads were downloaded for further processing and assembly.

**Nanopore read alignment.** From the base called read files produced by the Metrichor EPI2ME platform sequence files in FASTA format were extracted using the R-package poRe version 0.17 (ref. 48). We used BWA-MEM<sup>49</sup> (version 0.7.15-r1140) to align nanopore reads to selected seeds, using specific settings for each nanopore chemistry. The built-in *-x ont2d* setting (*-k 14 - W 20 -r 10 -A 1 -B 1 -O 1 -E 1 -L 0*) is too tolerant for newer chemistries. We therefore optimized alignment settings (*-k* and *-W* only) on small subsets to yield the highest recall (number of aligning reads) at the highest precision (number of seeds detected/number of alignments). With all other settings as before, this yielded the following parameters: *-k 14 - W 45* (R7.3 2D); *-k 16 - W 50* (R9 1D); *-k 19 - W 60* (R9 2D); *-k 16 - W 60* (R9.4 1D).

**Genome assembly using TULIP.** Currently, TULIP consists of two prototype scripts in Perl: *tulipseed.perl* and *tulipbulb.perl* (version 0.4 'European eel'). The *tulipseed* script constructs the seed graph based on input SAM files and a set seed length, and outputs a simplified graph and seed arrangements (scaffold models). *tulipbulb* adds seed and long read sequence to the scaffolds, and exports either a complete set of uncorrected scaffolds, or for each scaffold two separate files: the uncorrected sequence, and a FASTA 'bundle' consisting of all long reads associated with that scaffold.

For each scaffold, we used the long read bundle and Illumina data to polish it according to ONT guidelines (https://github.com/nanoporetech/ont-assembly-polish). We first corrected nanopore-derived scaffolds with nanopore data using Racon<sup>22</sup>, based on alignments produced by Graphmap<sup>50</sup> version 0.3.0. Ultimately Racon sequence correction is performed by SPOA<sup>51</sup>, which is a partial order alignment algorithm that generates consensus sequences.

Subsequently, we used previously generated<sup>2</sup> Illumina data (NCBI Sequence Read Archive SRR5235521–SRR5235523), trimmed to Phred 30 quality values (using Sickle version 1.33, https://github.com/najoshi/sickle) in a second correction step using Pilon (version 1.21), an integrated software tool for assembly improvement<sup>37, 38</sup>. Pilon uses evidence from the alignment between short-read data and Racon-corrected scaffolds to identify events that are different in the draft genome compared to the support of short-read data.

All genome assembly steps and analyses were performed on a desktop computer equipped with an Intel Xeon E3-1241 3.5 GHz processor, in a virtual machine (Oracle VirtualBox version 4.3.26) running Ubuntu 16.04 LTS with 28 GB RAM and 4 processor threads available. For the final candidate assembly, the TULIP scripts required a maximum of 4.4 GB RAM.

**Genome alignment.** Uncorrected scaffolds were aligned against the 2012 scaffolds using nucmer<sup>52</sup> version 3.23, with settings --maxmatch and --minmatch 100, filtered for optimal correspondence (*delta-filter -1*), and visualized using mummerplot (with the --layout option). The five largest scaffolds were likewise aligned against the 2012 scaffolds, but with settings encouraging longer alignments (*--breaklen 1000* and *--minmatch 25*) and not filtered. The 285 nt seeds were aligned against the 2012 draft scaffolds using BWA-MEM with default settings.

**BUSCO assembly assessment.** The completeness of the genome assemblies was tested with BUSCO<sup>36</sup> (version 3.0.0), which tries to find orthologues of a curated dataset of near-universal genes in new assemblies. A more complete assembly will result in a higher percentage of genes retrieved. As the European eel is a primitive teleost, we used the vertebrate-specific orthologue catalogue (*vertebrata\_odb9*, creation date 13-2-2016, 2586 genes) instead of *actinopterygii\_odb9*, which is based predominantly on the genome sequences of advanced teleosts.

**Data availability.** The nanopore sequencing data are available in the European Nucleotide Archive (accession number PRJEB20018). The Racon- and Pilon-corrected candidate assembly is available at http://www.eelge-nome.com. The TULIP-scripts are available at https://github.com/Generade-nl

#### References

- Coppe, A. et al. Sequencing, de novo annotation and analysis of the first Anguilla anguilla transcriptome: EeelBase opens new
  perspectives for the study of the critically endangered European eel. BMC Genomics 11, 635 (2010).
- 2. Henkel, C. V. et al. Primitive duplicate Hox clusters in the European eel's genome. PLoS One 7, e32231 (2012).
- 3. Pujolar, J. M. *et al.* Surviving in a toxic world: transcriptomics and gene expression profiling in response to environmental pollution in the critically endangered European eel. *BMC Genomics* **13**, 507 (2012).
- Minegishi, Y., Henkel, C. V., Dirks, R. P. & van den Thillart, G. E. Genomics in eels towards aquaculture and biology. *Mar Biotechnol (NY)* 14, 583–590 (2012).
- 5. IUCN Red List. doi:10.2305/IUCN.UK.2014-1.RLTS.T60344A45833138.en (2014).
- 6. Ager-Wick, E. et al. The pituitary gland of the European eel reveals massive expression of genes involved in the melanocortin system. PLoS One 8, e77396 (2013).
- 7. Dirks, R. P. et al. Identification of molecular markers in pectoral fin to predict artificial maturation of female European eels (Anguilla anguilla). Gen Comp Endocrinol 204, 267–276 (2014).
- Churcher, A. M. *et al.* Deep sequencing of the olfactory epithelium reveals specific chemosensory receptors are expressed at sexual maturity in the European eel *Anguilla anguilla. Mol Ecol* 24, 822–834 (2015).
- 9. Burgerhout, E. *et al.* Changes in ovarian gene expression profiles and plasma hormone levels in maturing European eel (*Anguilla anguilla*); biomarkers for broodstock selection. *Gen Comp Endocrinol* 225, 185–196 (2016).
- 10. Pelster, B., Schneebauer, G. & Dirks, R. P. Anguillicola crassus infection significantly affects the silvering related modifications in steady state mRNA levels in gas gland tissue of the European eel. Front Physiol 7, 175 (2016).
- 11. Pujolar, J. M. *et al.* A resource of genome-wide single-nucleotide polymorphisms generated by RAD tag sequencing in the critically endangered European eel. *Mol Ecol Resour* 13, 706–714 (2013).
- 12. Pasquier, J. *et al.* Multiple kisspeptin receptors in early osteichthyans provide new insights into the evolution of this receptor family. *PLoS One* 7, e48931 (2012).
- 13. Maugars, G. & Dufour, S. Demonstration of the coexistence of duplicated LH receptors in teleosts, and their origin in ancestral actinopterygians. *PLoS One* **10**, e0135184 (2015).
- 14. Morini, M. et al. Duplicated leptin receptors in two species of eel bring new insights into the evolution of the leptin system in vertebrates. PLoS One 10, e0126008 (2015).
- 15. Jacobsen, M. W. et al. Genomic footprints of speciation in Atlantic eels (Anguilla anguilla and A. rostrata). Mol Ecol 23, 4785–4798 (2014).
- 16. Pujolar, J. M. *et al.* Genome-wide single-generation signatures of local selection in the panmictic European eel. *Mol Ecol* 23, 2514–2528 (2014).
- 17. Ip, C. L. et al. MinION Analysis and Reference Consortium: Phase 1 data release and analysis. F1000Res 4, 1075 (2015).
- Jain, M., Olsen, H. E., Paten, B. & Akeson, M. The Oxford Nanopore MinION: delivery of nanopore sequencing to the genomics community. *Genome Biol* 17, 239 (2016).
- 19. Tyson, J. R. *et al.* Whole genome sequencing and assembly of a *Caenorhabditis elegans* genome with complex genomics rearrangements using the MinION sequencing device. *BioRxiv*, doi:10.1101/099143 (2017).
- Weirather, J. L. *et al.* Comprehensive comparison of Pacific Biosciences and Oxford Nanopore Technologies and their applications to transcriptome analysis [version 1; referees: 2 approved with reservations]. *F1000Res* 6, 100 (2017).
- Koren, S. *et al.* Reducing assembly complexity of microbial genomes with single-molecule sequencing. *Genome Biol* 14, R101 (2013).
   Loman, N. J., Quick, J. & Simpson, J. T. A complete bacterial genome assembled de novo using only nanopore sequencing data. *Nat*
- Methods 12, 733–735 (2015).
  23. Vaser, R., Sović, I., Nagarajan, N. & Šikić, M. Fast and accurate *de novo* genome assembly from long uncorrected reads. *Genome Res* 27, 737–746 (2017).
- 24. Chin, C. S. et al. Phased diploid genome assembly with single-molecule real-time sequencing. Nat Methods 13, 1050–1054 (2016).
- Koren, S., Walenz, B. P., Berlin, K., Miller, J. R. & Phillippy, A. M. Canu: scalable and accurate long-read assembly via adaptive k-mer weighting and repeat separation. *Genome Res* 27, 722–736 (2017).
- Kamath, G. M., Shomorony, I., Xia, F., Courtade, T. A. & Tse, D. N. HINGE: long-read assembly achieves optimal repeat resolution. Genome Res 27, 747–756 (2017).
- 27. Li, H. Minimap and miniasm: fast mapping and de novo assembly for noisy long sequences. Bioinformatics 32, 2103–2110 (2016).
- Ye, C., Hill, C. M., Wu, S., Ruan, J. & Ma, Z. DBG2OLC: efficient assembly of large genomes using long erroneous reads of the third generation sequencing technologies. Sci Rep 6, 31900 (2016).
- Keinath, M. C. et al. Initial characterization of the large genome of the salamander Ambystoma mexicanum using shotgun and laser capture chromosome sequencing. Sci Rep 5, 16413 (2015).
- 30. Biscotti, M. A. *et al.* The lungfish transcriptome: a glimpse into molecular evolution events at the transition from water to land. *Sci Rep* **6**, 21571 (2016).
- Zonneveld, B. J. The systematic value of nuclear genome size for all species of *Tulipa* L. (Liliacaeae). *Plant Syst Evol* 281, 217–245 (2009).
- Li, X. & Waterman, M. S. Estimating the repeat structure and length of DNA sequences using *l*-tuples. *Genome Res* 13, 1916–1922 (2003).
- Vuture, G. W. et al. GenomeScope: fast reference-free genome profiling from short reads. Bioinformatics, doi:10.1093/bioinformatics/ btx153 (2017).
- Boetzer, M., Henkel, C. V., Jansen, H. J., Butler, D. & Pirovano, W. Scaffolding pre-assembled contigs using SSPACE. *Bioinformatics* 27, 578–579 (2011).
- 35. Henkel, C. V. et al. First draft genome sequence of the Japanese eel. Anguilla japonica. Gene 511, 195–201 (2012).
- Simão, F. A., Waterhouse, R. M., Ioannidis, P., Kriventseva, E. V. & Zdobnov, E. M. BUSCO: assessing genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics* 31, 3210–3212 (2015).
- Walker, B. J. et al. Pilon: an integrated tool for comprehensive microbial variant detection and genome assembly improvement. PLoS One 9, e112963 (2014).

- Goodwin, S. et al. Oxford Nanopore sequencing, hybrid error correction, and de novo assembly of a eukaryotic genome. Genome Res 25, 1–7 (2015).
- 39. Pavey, S. A. et al. Draft genome of the American eel (Anguilla rostrata). Mol Ecol Resour, doi:10.1111/1755-0998.12608 (2016).
- 40. Albert, V., Jónsson, B. & Bernatchez, L. Natural hybrids in Atlantic eels (*Anguilla anguilla, A. rostrata*): evidence for successful reproduction and fluctuating abundance in space and time. *Mol Ecol* **15**, 1903–1916 (2006).
- 41. Hunt, M., Newbold, C., Berriman, M. & Otto, T. D. A comprehensive evaluation of assembly scaffolding tools. *Genome Biol* 15, R42 (2014).
- Chapman, J. A. et al. A whole-genome shotgun approach for assembling and anchoring the hexaploidy bread wheat genome. Genome Biol 16, 26 (2015).
- Urban, J. M., Bliss, J., Lawrence, C. E. & Gerbi, S. A. Sequencing ultra-long DNA molecules with the Oxford Nanopore MinION. *BioRxiv*, doi:10.1101/019281 (2015).
- 44. Datema, E. *et al.* The megabase-sized fungal genome of *Rhizoctonia solani* assembled from nanopore reads only. *BioRxiv*, doi:10.1101/084772 (2016).
- 45. Mostovoy, Y. et al. A hybrid approach for de novo human genome sequence assembly and phasing. Nat Methods 13, 587-590 (2016).
- 46. Marçais, G. & Kingsford, C. A fast, lock-free approach for efficient parallel counting of occurrences of *k*-mers. *Bioinformatics* **27**, 764–770 (2011).
- Magoc, T. & Salzberg, S. FLASH: Fast length adjustment of short reads to improve genome assemblies. *Bioinformatics* 27, 2957–2963 (2011).
- Watson, M. et al. poRe: an R package for the visualization and analysis of nanopore sequencing data. Bioinformatics 31, 114–115 (2015).
- 49. Li, H. & Durbin, R. Fast and accurate long-read alignment with Burrows-Wheeler transform. Bioinformatics 26, 589–95 (2010).
- 50. Sović, I. et al. Fast and sensitive mapping of nanopore sequencing reads with GraphMap. Nat Commun 7, 11307 (2016).
- 51. Lee, C. Generating consensus sequences from partial order multiple sequence alignment graphs. *Bioinformatics* **19**, 999-1008 (2003).
- 52. Kurtz, S. et al. Versatile and open software for comparing large genomes. Genome Biol 5, R12 (2004).

#### Acknowledgements

We are grateful to the DUPAN Foundation for making this project possible, as well as to Rosemary Dokos and Oliver Hartwell at Oxford Nanopore for technical support and encouragement. This project was funded by grants from the DUPAN Foundation for sustainable eel farming and fishing, the Dutch Ministry of Economic Affairs (to A.P.P., KB-21-001-001), the Austrian Science Foundation (to B.P., FWF P26363-B25), the European Union's Horizon 2020 research and innovation programme under the Marie Sklodowska-Curie Actions: Innovative Training Network IMPRESS, grant agreement No 642893 (to F.-A.W.), ZonMW/DTL (to C.V.H.), and local funds from CNRS (to S.D.) and Generade, the Leiden Centre of Expertise in Genomics (to C.V. H.).

#### **Author Contributions**

H.J.J., S.D., F.-A.W., W.S., A.K., A.P.P., B.P., H.P.S., G.E.V.D.T., R.P.D. and C.V.H. conceived the research. R.P.D. coordinated the project. H.J.J. and S.A.J.-R. performed sequencing, M.L. and C.V.H. assembled the genome, H.J.J., R.P.D. and C.V.H. analyzed the data. H.J.J., M.L., R.P.D. and C.V.H. wrote the paper with input from all other authors.

### **Additional Information**

Supplementary information accompanies this paper at doi:10.1038/s41598-017-07650-6

**Competing Interests:** H.J.J. and C.V.H. are members of the Nanopore Community, and have previously received flow cells free of charge (used for some of the R7.3 data of this project), as well as travel expense reimbursements from Oxford Nanopore Technologies.

**Publisher's note:** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit http://creativecommons.org/licenses/by/4.0/.

© The Author(s) 2017