



TECHNICAL NOTE

Chiron: translating nanopore raw signal directly into nucleotide sequence using deep learning

Haotian Teng^{1,*}, Minh Duc Cao^{1,†}, Michael B. Hall¹, Tania Duarte¹, Sheng Wang² and Lachlan J.M. Coin^{1,*}

¹Institute for Molecular Bioscience, University of Queensland, St Lucia, Brisbane, QLD 4072, Australia and ²Computational Bioscience Research Center (CBRC), King Abdullah University of Science and Technology (KAUST), Thuwal, 23955, Saudi Arabia

*Correspondence address. Lachlan J.M. Coin, Institute for Molecular Bioscience, University of Queensland, St Lucia, Brisbane, QLD 4072, Australia. E-mail: l.coin@imb.uq.edu.au  <http://orcid.org/0000-0003-0337-8722>; Haotian Teng, Institute for Molecular Bioscience, University of Queensland, St Lucia, Brisbane, QLD 4072, Australia. E-mail: haotian.teng@uq.net.au  <http://orcid.org/0000-0002-4300-455X>

[†]These authors contributed equally to this work.

Abstract

Sequencing by translocating DNA fragments through an array of nanopores is a rapidly maturing technology that offers faster and cheaper sequencing than other approaches. However, accurately deciphering the DNA sequence from the noisy and complex electrical signal is challenging. Here, we report Chiron, the first deep learning model to achieve end-to-end basecalling and directly translate the raw signal to DNA sequence without the error-prone segmentation step. Trained with only a small set of 4,000 reads, we show that our model provides state-of-the-art basecalling accuracy, even on previously unseen species. Chiron achieves basecalling speeds of more than 2,000 bases per second using desktop computer graphics processing units.

Keywords: ONT nanopore sequencing; deep learning; artificial neural network; comparative performance

Introduction

DNA sequencing via bioengineered nanopores, recently introduced to the market by Oxford Nanopore Technologies (ONT), has profoundly changed the landscape of genomics. A key innovation of the ONT nanopore sequencing device, MinION, is that it measures the changes in electrical current across the pore as a single-stranded molecule of DNA passes through it. The signal is then used to determine the nucleotide sequence of the DNA strand [1–3]. Importantly, this signal can be obtained and analyzed by the user while the sequencing is still in progress. A large number of pores can be packed into a MinION device that is the size of a stapler, making the technology extremely portable. The small size and real-time nature of the sequencing opens up new opportunities in time-critical genomics applications [4–7] and in remote regions [8–11].

While nanopore sequencing can be massively scaled up by designing large arrays of nanopores and allowing faster translocation of DNA fragments, one of the bottlenecks in the analysis pipeline is the translation of the raw signal into nucleotide sequence, or basecalling. Prior to the release of Chiron, basecalling of nanopore data involved two stages. Raw data series are first divided into segments corresponding to signals obtained from a k-mer (segmentation) before a model is then applied to translate segment signals into k-mers. DeepNano[12] introduced the idea of using a bidirectional recurrent neural network (RNN) that uses the basic statistics of a segment (mean signal, standard deviation, and length) to predict the corresponding k-mer. The official basecallers released by ONT, nanonet, and Albacore (prior to v2.0.1) also employ similar techniques. As k-mers from successive segments are expected to overlap by k-1 bases, these methods use a dynamic programming algorithm to find the most

Received: 9 November 2017; Revised: 7 February 2018

© The Author(s) 2018. Published by Oxford University Press. This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

probable path, which results in the basecalled sequence data. BasecRAWller [13] uses a pair of unidirectional RNNs; the first RNN predicts the probability of segment boundary for segmentation, while the second one translates the discrete event into base sequence. As such, BasecRAWller is able to process the raw signal data in a streaming fashion.

In this article we present Chiron, which is the first deep neural network model that can translate raw electrical signal directly to nucleotide sequence. Chiron has a novel architecture that couples a convolutional neural network (CNN) with an RNN and a connectionist temporal classification (CTC) decoder [13]. This enables it to model the raw signal data directly, without use of an event segmentation step. ONT has also developed a segmentation free basecaller, Albacore v2.0.1, which was released shortly after Chiron v0.1.

Chiron was trained on a small dataset sequenced from a viral and bacterial genome and yet it is able to generalize to a range of genomes such as other bacteria and Human. Chiron is as accurate as the ONT-designed and -trained Albacore v2.0.1 on bacterial and viral basecalling and outperforms all other existing methods. Moreover, unlike Albacore, Chiron allows users to train their own neural network, and it is also fully open-source, enabling development of specialized basecalling applications, such as detection of base modifications.

Results

Deep neural network architecture

We have developed a deep neural network (NN) for end-to-end, segmentation-free basecalling that consists of two sets of layers: a set of convolutional layers and a set of recurrent layers (Fig. 1). The convolutional layers discriminate local patterns in the raw input signal, whereas the recurrent layers integrate these patterns into basecall probabilities. At the top of the neural network is a CTC decoder [14] to provide the final DNA sequence according to the base probabilities (Fig. 2). More details pertaining to the NN are provided in the Methods section.

Chiron presents an end-to-end basecaller in that it predicts a complete DNA sequence from raw signal. It translates sliding windows of 300 raw signals to sequences of roughly 10-20 base pairs (which we call slices). These overlapping slices are stacked together to get a consensus sequence in real time. The window is shifted by 30 raw signals; by processing the slices in parallel, the basecalling accuracy can be improved with little speed loss.

Performance comparison

For training and evaluating the performance of Chiron, a phage Lambda virus sample (*Escherichia virus Lambda*) provided by ONT and an *Escherichia coli* (K12 MG1655) sample using 1D protocol on R9.4 flowcells were sequenced for calibrating the MinION device (see the Methods section). A total of 34,383 reads were obtained for the Lambda sample and 15,012 reads were obtained for *E. coli*, but only 2,000 reads were randomly picked from each sample to train Chiron. It took the model 10 hours to train 3 epoch with 4,000 reads (~4 Mbp) on a Nvidia K80 Graphics Processing Unit (GPU). Then Chiron was cross-validated on the remainder of the reads from two runs, and the model was further evaluated by testing its basecalling accuracy on other species. A *Mycobacterium tuberculosis* sample was sequenced and a set of Human data was downloaded from chromosome 21 part 3 from the Nanopore WGS Consortium [15], to be used in testing the generality of Chiron.

In order to establish the ground-truth of the data, the *E. coli* and *M. tuberculosis* samples were sequenced using Illumina technology (see the Methods section) and assembled, which provided a high per-base accuracy reference. The reference sequence for the Phage Lambda virus was National Center for Biotechnology Information (NCBI) reference sequence NC_001416.1; for the Human data, the GRCh38 reference was used. The raw signals were labeled by identifying the raw signal segment corresponding to the nucleotide assumed to be in the pore at a given time point (see the Methods section).

Table 1 presents the accuracy of the four basecalling methods, including the Metrichor basecaller (ONT cloud service), Albacore v1.1 (ONT official local basecaller), BasecRAWller [13], and Chiron, with a greedy decoder (Chiron) and beam-search decoder (Chiron-BS), on the data. Chiron had the highest identity rate on the Lambda, *E. coli*, and *M. tuberculosis* samples. Additionally, it had the lowest deletion rate; mismatch rate on Lambda, *M. tuberculosis*, and *E. coli*; and the lowest insertion rate on Lambda and *E. coli*. In the Human dataset where Chiron did not have the highest identity rate, it was no more than 0.01 from the best.

In addition, we compared the segmentation-free ONT basecaller Albacore v2.0.1 with Chiron-BS in Table 1. Chiron-BS had a consistently lower insertion rate across all species tested, as well as a lower deletion rate on Lambda and *E. coli*; however, it suffered a slightly higher mismatch rate on all species except *E. coli*. The performance is comparable to Albacore v2.0.1 on all species except for Human; however, this is likely at least partially due to the fact that it had not been trained on any Human DNA.

In order to assess the quality of genomes assembled from reads generated by each basecaller, we used Miniasm together with Racon to generate a *de novo* genome assembly for each bacterial and viral genome (see the Methods section). The results presented in Table 2 demonstrate that Chiron assemblies for Phage lambda and *E. coli* had approximately half as many errors as those generated from Albacore (v1 or v2) reads. For *M. tuberculosis*, Chiron had fewer errors than Albacore v1 but slightly more than Albacore v2. The identity rate and relative length for each round of polishing with Racon are shown in Fig. 3.

In terms of speed on a central processing unit (CPU) processor (Table 3), Chiron is slower (21 bp/sec, 17 bp/sec using a beam-search decoder with a 50 beam width) than Albacore (2,975 bp/sec) and, to a lesser extent, slower than BasecRAWller (81bp/sec). However, when run on a Nvidia K80 GPU, a basecalling rate of 1,652 bp/sec and 1,204 bp/sec using a beam-search decoder is achieved. (Chiron was also tested on a Nvidia GTX 1080 Ti GPU, and the rate was 2,657 bp/sec). The GPU rate for the other two local basecallers are not included, as Albacore and BasecRAWller do not currently offer GPU support. Metrichor was not included in the speed benchmarking as it is not possible to gather information about CPU/GPU speed as it is a cloud basecaller.

Discussion

Segmenting the raw nanopore electrical signal into piece-wise constant regions that correspond to the presence of different k-mers in the pore is an appealing but error-prone approach. Segmentation algorithms determine a boundary between two segments based on a sharp change of signal values within a window. The window size is determined by the expected speed of the translocation of the DNA fragment in the pore. We noticed that the speed of DNA translocation is variable during a

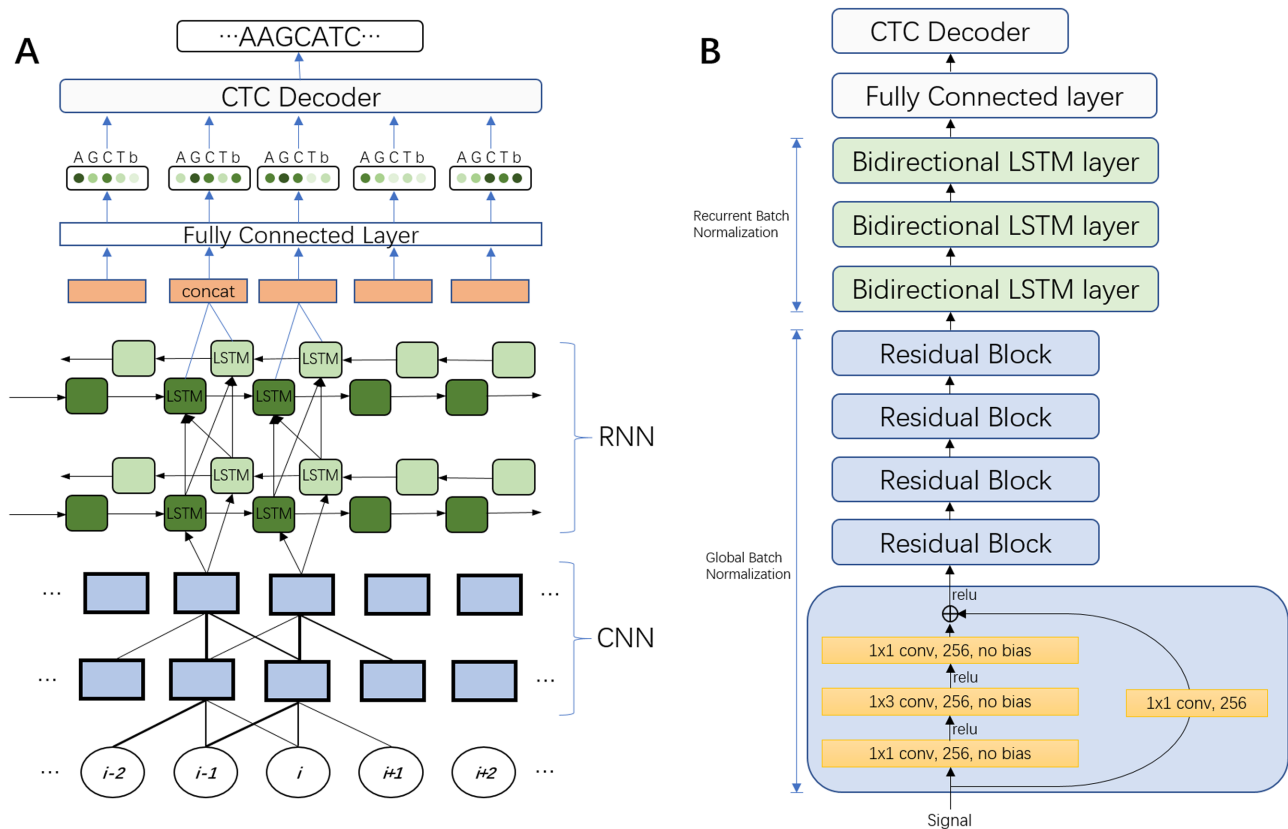


Figure 1: (A) An unrolled sketch of the NN architecture. The circles at the bottom represent the time series of raw signal input data. Local pattern information is then discriminated from this input by a CNN. The output of the CNN is then fed into an RNN to discern the long-range interaction information. A FC layer is used to get the base probability from the output of the RNN. These probabilities are then used by a CTC decoder to create the nucleotide sequence. The repeated component is omitted. (B) Final architecture of the Chiron model. Variants of this architecture were explored by varying the number of convolutional layers from 3 to 10 and recurrent layers from 3 to 5. We also explored networks with only convolutional layers or recurrent layers, 1×3 conv, 256, no bias means a convolution operation with a 1×3 filter and a 256-channel output with no bias added. LSTM = long-term short memory.

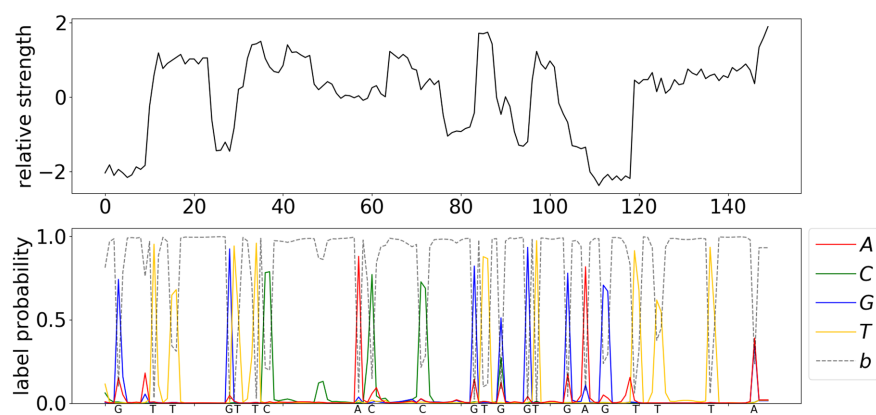


Figure 2: Visualization of the predicted probability of bases and the readout sequence. The upper panel is a normalized raw signal from the MinION nanopore sequencer, normalized by subtracting the mean of the whole signal and then dividing by the standard deviation. The bottom panel shows the predicted probability of each base at each position from Chiron. The final output DNA sequence is annotated on the x-axis of the bottom plane.

sequencing run; coupled with the high level of signal-to-noise in the raw data, this can result in low segmentation accuracy. As a result, the segmentation algorithm often makes conservative estimates of the window size, resulting in segments that are smaller than the actual signal group for k-mers. While dynamic programming can correct this by joining several segments together for a k-mer, this effects the prediction model.

All existing nanopore basecallers prior to Chiron use a segmentation step. The first nanopore basecalling algorithms [16, 17] used a hidden Markov model, which maintains a table of event models for all possible k-mers. These event models were learned from a large set of training data. More recent methods (DeepNano [12], nanonet) train a deep neural network for inferring k-mers from segmented raw signal data.

Table 1: Results from the experimental validation and benchmarking of Chiron against three segmentation-based nanopore basecallers and Albacore v2, which is also a segmentation-free basecaller

Dataset	Basecaller	Deletion rate (%)	Insertion rate (%)	Mismatch rate (%)	Identity rate (%)	Error rate (%)
Lambda	Metrichor	8.93	2.38	4.57	86.50	15.88
	Albacore v1.1	6.35	3.82	4.69	88.96	14.86
	Albacore v2	6.19	3.38	3.98	89.82	13.55
	BasecRAWller	7.89	10.01	10.56	81.54	28.46
	Chiron	8.20	2.13	4.03	87.76	14.36
	Chiron-BS	6.20	2.13	4.20	89.60	12.53
E. coli	Metrichor	7.52	1.93	3.84	88.64	13.29
	Albacore v1.1	5.76	3.27	4.14	90.10	13.17
	Albacore v2	5.21	2.99	3.57	91.22	11.77
	BasecRAWller	7.16	10.40	10.30	82.54	27.86
	Chiron	6.36	1.81	3.07	90.57	11.24
	Chiron-BS	4.94	2.36	3.16	91.90	10.46
M. tuberculosis	Metrichor	7.63	2.40	4.35	88.02	14.38
	Albacore v1.1	6.12	3.57	4.68	89.19	14.37
	Albacore v2	5.05	3.58	4.05	90.90	12.68
	BasecRAWller	7.17	10.85	10.42	82.41	28.44
	Chiron	7.16	2.50	4.33	88.51	13.99
	Chiron-BS	5.84	3.05	4.50	89.66	13.39
Human	Metrichor	12.95	4.15	7.65	79.4	24.75
	Albacore v1.1	8.62	6.51	7.52	83.86	22.65
	Albacore v2	8.71	6.03	6.05	85.24	20.79
	BasecRAWller	8.41	10.28	10.10	81.49	28.79
	Chiron	9.13	5.14	9.33	81.54	23.60
	Chiron-BS	9.30	5.62	7.87	82.83	22.79

Deletion, insertion, and mismatch rates (%) are defined as the number of deleted, inserted, and mismatched bases divided by the number of bases in the reference genome (the lower the better). Identity rate (%) is defined as the number of matched bases divided by the number of bases in the reference genome for that sample (the higher the better; identity rate = 1 - deletion rate - mismatch rate). Error rate (%) is defined as the sum of deletion, insertion, and mismatch rates (the lower the better; error rate = deletion rate + insertion rate + mismatch rate). This statistic effectively summarizes the basecalling accuracy of the associated model. The best result in each category is indicated in bold.

Table 2: Assembly identity rate and relative length benchmark,

Sample (coverage)	Albacore	Albacore.2	Chiron-BS	Metrichor	Albacore	Albacore.2	Chiron-BS	Metrichor
E. coli-S18 (27X)	99.004	99.162	99.533	87.678	100.055	99.715	99.720	94.253
E. coli-S10 (40X)	99.106	99.316	99.646	88.745	100.144	99.739	99.811	94.829
M. tuberculosis (130X)	99.541	99.628	99.554	84.736	100.126	100.029	99.900	90.875
Lambda Phage (790X)	97.926	99.207	99.507	99.164	101.104	100.123	99.800	99.335

Draft genome generated by Miniasm and polished 10 rounds by Racon. Assembly identity rates are presented in the left 4 columns; relative lengths are presented in the right 4 columns. Identity rate (%) is calculated by first shredding the assembly contigs into 10K read fragments and then obtaining the mean of the identity rate of the aligned reads. Relative length (%) is defined as the sum of the length of all the aligned pieces divided by the length of the reference genome. E. coli-S10 and E. coli-S18 are reads from two independent sequencing events.

A recent basecaller named BasecRAWller [13] was used an initial neural network (referred to as a *raw* network) to output probabilities of boundaries between segments. A segmentation algorithm was then applied to segment these probabilities into discrete events. BasecRAWller then used a second neural network (referred to as the *fine-tune* network) to translate the segmented data into the base sequence.

Our proposed model is a departure from the above approaches in that it performs base prediction directly from raw data without segmentation. Moreover, the core model is an end-to-end basecaller in the sense that it predicts the complete base sequence from raw signal. This is made possible by combining a multilayer convolutional neural network to extract the local

features of the signal, with a recurrent neural network to predict the probability of nucleotides in the current position. Finally, the complete sequence is called by a simple greedy algorithm, based on a typical CTC-style decoder [14], reading out the nucleotide in each position with the highest probability. Thus, the model need not make any assumption of the speed of DNA fragment translocation and can avoid the errors introduced during segmentation.

To improve the basecalling speed and minimize its memory requirements, the neural network is run on a 300-signal sliding window (equivalent to approximately 20bp), overlapping the sequences on these windows and generating a consensus sequence. Chiron has the potential to stream these input raw

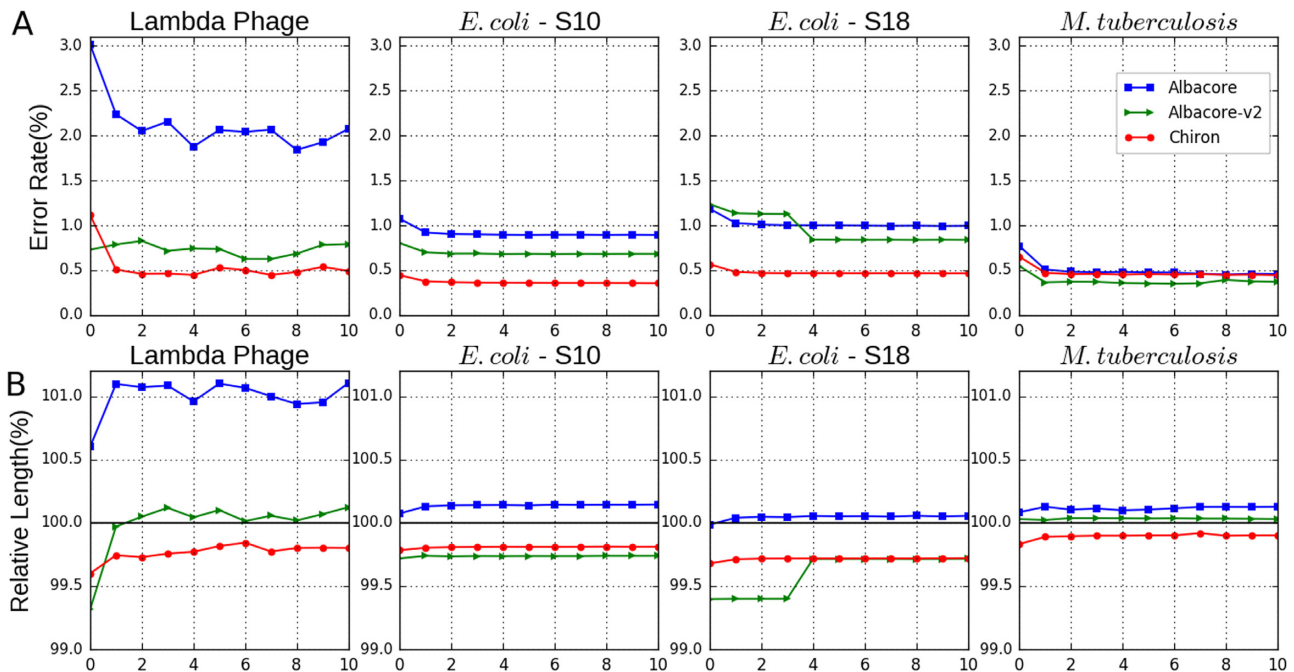


Figure 3: (A) Assembly error rate (%) for each polishing round using Racon. Two individually sequenced *E. coli* samples are included (S10, S18). All basecallers have a similar performance on the *M. tuberculosis* dataset due to its high sequencing depth (130X). (B) Relative assembly length (%) after each round of polishing. Relative length is defined as the length of the assembly divided by the length of reference genome.

signal "slices" into output sequence data, which will become an increasingly important aspect of basecalling very long reads (100kb+), particularly if used in conjunction with the read-until capabilities of the MinION.

Our model was either the best or second-best in terms of accuracy on all of the datasets we tested in terms of read-level accuracy. This includes the Human dataset, despite the fact that the model had not seen Human DNA during training. Our model had only been trained on a mixture of 2,000 bacterial and 2,000 viral reads. The most accurate basecaller is the proprietary ONT Albacore basecaller. Chiron is within 1% accuracy on bacterial DNA but only within 2% accuracy on Human DNA. More extensive training on a broader spectrum of species, including Human, can be expected to improve the performance of our model. There are also improvements in accuracy to be gained from better alignment of overlapping reads and consensus calling. Increasing the size of the sliding window will also improve accuracy but at the cost of increased memory and running time.

Bacterial and viral genome assemblies generated from Chiron basecalled reads all had less than 0.5% error, whereas those generated by Albacore had up to 0.8% error, Fig. 3. This marked reduction in error rate is essential for generating accurate single-nucleotide polymorphism genotypes, a prerequisite for many applications such as outbreak tracking. These results are consistent with those reported in a recent study of read and assembly level accuracy for *Klebsiella pneumoniae* [18].

Our model is substantially more computationally expensive than Albacore and somewhat more computationally expensive than BasecRAWler. This is to be expected given the extra depth in the neural network. Our model can be run in a GPU mode, which makes computation feasible on small- to medium-sized datasets on a modern desktop computer. Our method can be further sped up by increasing the step size of the sliding window, although this may impact accuracy. Also, there are several existing methods that can be used to accelerate NN-based base-

Table 4: Details on the number of reads and their median read length for data that was used to evaluate various basecallers

Sample	Number of reads	Median read length (bp)
Phage Lambda	34,383	5,720
<i>E. coli</i>	15,012	5,836
<i>M. tuberculosis</i>	147,594	3,423
Human	10,000	6,154

callers such as Chiron. One such example is Quantization, which reformats 32-bit float weights as 8-bit integers by binning the weight into a 256 linear set. As neural networks are robust to noise, this will likely have negligible impact on the performance. Weight pruning, which prunes the weights whose absolute value is under a certain threshold and then retrains the NN, is another method used to compress and accelerate NN [19].

Conclusion

We have presented a novel deep neural network approach for segmentation-free basecalling of raw nanopore signal. Our approach is the first method that can map the raw signal data directly to base sequence without segmentation. We trained our method on only 4,000 reads sequenced from the simple genome lambda virus and *E. coli*, but the method is sufficiently generalized to be able to basecall data from other species, including Human. Our method has state-of-art accuracy, outperforming the ONT cloud basecaller Metrichor as well as another third-party basecaller, BasecRAWler.

Methods

Deep neural network architecture

Our model combines a five-layer convolutional neural network (CNN) [20] with a three-layer recurrent neural network (RNN) and a fully connected (FC) layer in the last layer that calculates the probability for a CTC decoder to get the final output. This structure is similar to that used in speech recognition [21]. Both the CNN and RNN layers are found to be essential to the basecalling, as removing either would cause a dramatic drop in prediction accuracy, which is described more in the Training section.

Preliminaries. Let a raw signal input with T time points $\mathbf{s} = [s_1, s_2, \dots, s_T]$ and the corresponding DNA sequence label (with K bases) $\mathbf{y} = [y_1, y_2, \dots, y_K]$ with $y_i \in \{A, G, C, T\}$ be sampled from a training dataset $\chi = \{(\mathbf{s}^{(1)}, \mathbf{y}^{(1)}), (\mathbf{s}^{(2)}, \mathbf{y}^{(2)}), \dots\}$. Our network directly translates the input signal time series \mathbf{s} to the sequence \mathbf{y} without any segmentation steps.

The input signal is normalized by subtracting the mean of the whole read and dividing by the standard deviation. $\mathbf{s}' = (\mathbf{s} - \bar{s})/\text{std}(\mathbf{s})$.

Then the normalized signal is fed into a residual block [22] combined with global batch normalization [23] in the five convolution layers to extract the local pattern from the signal. The stride is set as 1 to ensure the output of the CNN has the same length as the input raw signal. The residual block is illustrated in Fig. 1. A convolution operation with a $l \times m$ filter, $n \times p$ stride, and s output channels on a k channels input is defined as:

$$\text{Output}(i, j, s) = \sum_{di < l, dj < m, q < k} \text{Input}(i \cdot n + di, j \cdot p + dj, q) \cdot \text{Filter}(di, dj, q, s).$$

An activation operation is performed after the convolution operation. Various kinds of activation functions can be chosen. However, in this model a rectified linear unit (ReLU) function is used as the activation operation, which has been reported to have a good performance in CNN, defined as :

$$\text{ReLU}(x) = \max(x, 0)$$

Following the convolution layers are multiple bidirectional RNN layers [24], a long short-term memory (LSTM) cell [25] is used as the RNN cell, with a separate batch normalization on the inside cell state and input term [26].

A typical batch normalization procedure [23] is

$$\text{BN}(\mathbf{x}; \gamma, \beta) = \beta + \gamma \odot \frac{\mathbf{x} - \hat{E}[\mathbf{x}]}{\sqrt{\hat{\text{Var}}[\mathbf{x}] + \epsilon}}, \quad (1)$$

where \mathbf{x} be a inactivation term.

Let \mathbf{h}_t^l be the output of l_{th} RNN layer at time t , the batch normalization for a LSTM cell is

$$(\mathbf{f}_t, \mathbf{i}_t, \mathbf{o}_t, \mathbf{g}_t) = \text{BN}(\mathbf{W}_h \mathbf{h}_{t-1}^l; \gamma_h, \beta_h) + \text{BN}(\mathbf{W}_x \mathbf{h}_{t-1}^{l-1}; \gamma_x, \beta_x) + \mathbf{b} \quad (2)$$

$$\mathbf{c}_t = \sigma(\mathbf{f}_t) \odot \mathbf{c}_{t-1} + \sigma(\mathbf{i}_t) \odot \tanh(\mathbf{g}_t) \quad (3)$$

$$\mathbf{h}_t = \sigma(\mathbf{o}_t) \odot \tanh(\text{BN}(\mathbf{c}_t; \gamma_c, \beta_c)) \quad (4)$$

The batch normalization is calculated separately in the recurrent term $\mathbf{W}_h \mathbf{h}_{t-1}^l$ as well as the input term $\mathbf{W}_x \mathbf{h}_{t-1}^{l-1}$. The parameters β_h and β_x are set to zero to avoid the redundancy with \mathbf{b} . The last forward layer \mathbf{h}_{tf}^l and the backward layer \mathbf{h}_{tb}^l are concatenated together as an input to a FC layer :

$$\mathbf{H}_i = [\mathbf{h}_{iw}^l, \mathbf{h}_{ib}^l]. \quad (5)$$

The final output is transferred through a FC layer followed by a softmax operation :

$$p(\mathbf{o}_i = j) = \frac{\exp \mathbf{W}_j \mathbf{H}_i}{\sum_j \exp \mathbf{W}_j \mathbf{H}_i} \quad (6)$$

The output \mathbf{o}_i , $i = 1, 2, \dots, T$ predicts the symbol given the input vector \mathbf{x} , $P(\mathbf{o}_i = j | \mathbf{x})$. If the read is a DNA sequence, then $j \in \{A, G, C, T, b\}$, where b represents a blank symbol (Fig. 1). During training, the CTC loss is calculated between the output sequence \mathbf{o} and label \mathbf{y} [13], and back-propagation is used to update the parameters. An Adam optimizer [26] with an initial learning rate of 0.001 is used to minimize the CTC loss.

During inference, the final sequence is constructed using either a greedy decoder [14] or a beam-search decoder [27]. The greedy decoder works by first getting the argument of maximum probability in each position of \mathbf{o} and then producing the sequence call by first removing the consecutive repeat, and then removing the blank symbols. For example, the greedy path of an output \mathbf{o} is A A - - A - - G -, here - represents the blank symbol, the consecutive repeat is removed first and leads to A - A - G -, and the blank is removed to get the final sequence AAG. The beam-search decoder, with beam width W , maintains a list of the W most probable sequences (after collapsing repeats and removing blanks) up to position i of \mathbf{o} . To obtain this list at position $i+1$, it constructs the probability of all possible extensions of the W most probable at position i based on adding each symbol according to $p(\mathbf{o}_i = j)$ and collapsing and summing up over repeated bases, or repeated blanks that are terminated by a nonblank. The greedy decoder is a special case of the beam-search decoder when the beam width is 1. It should be noted that the model can still call homopolymer repeats provided each repeated base is separated by a blank, which is typically the case.

Convolutional network to extract local patterns. A total of 256 channel filters are used for all five convolutional layers. In each layer, there is a residual block [28] (Fig. 1) with two branches. A 1×1 filter is used for reshaping in the first branch. In the second branch, a 1×1 convolution filter is followed by a ReLU [29] activation function and a 1×3 filter with a ReLU activation function as well as a 1×1 filter. All filters have the same channel number of 256. An element-wise addition is performed on the two branches followed by a ReLU activation function. A global batch normalization operation is added after every convolution operation. A large kernel size (5,7,11) and different channel numbers (128,1024) are also tested. The above combination is found to yield the best performance.

Recurrent layers for unsegmented labeling. The local pattern extracted from the CNN described above is then fed to a three-layer RNN (Fig. 1). Under the current ONT sequencing settings, the DNA fragments translocate through the pore with a speed of roughly 250 or 450 bases per second, depending on the sequencing chemistry used, while the sampling rate is 4,000 samples per second. Because the sampling rate is higher than the translocation rate, each nucleotide usually stays in the current position for about 5 to 15 samplings, on average. Furthermore, as

a number of nearby nucleotides also influence the current, 40 to 100 samples (based on a 4- or 5-mer assumption) could contain information about a particular nucleotide. A three-layer bidirectional RNN is used for extracting this long range information. LSTM cells [26, 30] with 200 hidden units are used in every layer, and a FC layer is used to translate the output from the last RNN layer into a prediction. The output of the FC layer is then fed into a CTC decoder to obtain the predicted nucleotide sequence for the given raw signals.

Improving basecalling performance. To achieve better accuracy and less memory allocation, a sliding window is applied (default of 300 raw signals), with a preset sliding step size (default of 10% of window size), to the long raw signal. This gives a group of short reads with uniform length (window length) that overlap the original long read. Then, basecalling is run in parallel on these short reads, and the whole DNA sequence is reassembled by finding the maximum overlap between two adjacent short reads and read out of the consensus sequence. Note that here the reassembly is very easy because the order of the short reads is known. This procedure improves the accuracy of the basecalling and also enables parallel processing on one read.

Data preparation

Sequencing. The library preparations of the *E. coli* and *M. tuberculosis* samples were done using the *1D gDNA selecting for long reads using SQK-LSK108* (March 2017 version) protocol with the following modifications. Increase the incubation time to 20 minutes in each end-repair and ligation step; use 0.7x Agencourt[®] AMPure[®] XP beads (Beckman Coulter) immediately after the end-repair step and incubation of the eluted beads for 10 minutes; and use elution buffer (ELB) warmed up at 50°C with the incubation of the eluted bead at the same temperature. For the Lambda sample, the *1D Lambda Control Experiment for MinION device using SQK-LSK108* (January 2017 version) protocol was followed with the following changes: sheared the sample at 4000 rpm (2x1 minutes); 30 minutes of incubation in each end-repair step; and 20 minutes for adaptor ligation and elution of the library with 17 μ L of ELB. All samples were sequenced on new FLO-MIN106, version R9.4, flow cells with more than 1,100 active single pores, and the phage was sequenced in a MinION Mk1 (232 ng in 6-hour run) while the bacteria samples were sequenced in a MinION Mk1B (1 μ g *E. coli* and 595 ng *M. tuberculosis* in 22-hour and 44-hour runs, respectively). The *E. coli* sample was run on the MinKNOW, version 1.4.3, and the other samples in earlier versions of the software. The *E. coli* sample was also sequenced on Illumina MiSeq using paired-end 300x2 to 100-fold coverage. An assembly of the *E. coli* genome was constructed by running Spades [31] on the MiSeq sequencing data of the sample. The genome sequence of the Phage Lambda is NCBI reference sequence NC.001416.1.

Labeling of raw signal. Metrichor, the basecaller provided by ONT that runs as a cloud service, is used to basecall the MinION sequencing data first. Then, Nanoraw [32] is used for labeling the data. Briefly, the basecalled sequence data are aligned back to the genome of the sample; from the alignment, the errors introduced by Metrichor are corrected to avoid the bias from Metrichor being learned into Chiron. The corrected data are mapped back to the raw data. The resulting labeling consists of the raw signal data, as well as the boundaries of raw signals when the DNA fragment translocates to a new base. We use the base-level segmentation of the raw data to obtain matched pairs of signal segment (of lengths 200, 400, and 1000) together with the corresponding DNA base sequence. From this point onwards the

exact matching of the signal to each base within a segment is disregarded.

Training and testing datasets. A dataset using 2,000 reads from *E. coli* and 2,000 reads from Phage Lambda is created for training Chiron. In every start of the training epoch, the dataset is shuffled first and then fed into the model by batch. Training on this mixture dataset gave the model better performance both on generality and accuracy on not only the *E. coli* and Phage Lambda but also on *M. tuberculosis* and Human data. The testing dataset is shown in Table 4.

Training

The labeling from Metrichor described previously is used to train Chiron. Although the neural network architecture is translation invariant and not restricted by the sequence length, a uniform length of sequences is suited for batch feeding and thus can accelerate the training process. From this view, the original reads were cut into short segments with a uniform length of 200, 400, and 1,000 and trained on these batches in alternation. Several different architectures of the neural network were tested (Table 5), with the CNN-RNN network architecture having the best accuracy compared to a CNN- or RNN-only network. Also, using more layers seems to increase the performance of the model; however, the time consumed for training and basecalling is also increased. In the final structure, an NN with five convolution layers and three recurrent layers is adopted, as adding layers above this structure gave negligible performance improvement but required more calculation and also increased the risk of overfitting (Table 5).

Parameters for basecalling

All basecallers were invoked on the same set of reads for each sample. When using Chiron to basecall, the raw signal was first sliced by a 300 length window, the window was slid by 30, and the sliced segments were fed into the basecaller with a batch size equal to 1,100. Then, the output short reads were simply assembled by a pair-wise alignment between neighboring reads, and the consensus sequence was output from this alignment. All basecalling with Albacore (v1.1.1 and v2.0.1) and BasecRAWller [13] (version 0.1) was done with default parameters. For the configuration setting in Albacore, `r94.450bps.linear.cfg` was used for all samples, as this matches the flowcell and kit used for each sample. The data were basecalled on Metrichor on 3 June 2017 (Lambda), 18 May 2017 (*E. coli*), 4 June 2017 (*M. tuberculosis*), and 20 June 2017 (NA12878-Human).

Quality score

The quality score is calculated using the following algorithm: $qs = 10 * \log_{10}(\frac{P_1}{P_2})$, where P_1 is the probability of the most probable base in the current position, and P_2 is the probability of the second probable base in the current position.

Comparison of raw read accuracy

To assess the performance of each program, the resulting FASTA/FASTQ file from basecalling was aligned to the reference genome using graphmap [33] with the default parameters. The resulting BAM file was then assessed using the japsa error analysis tool (`jsa.hts.errorAnalysis`), which looks at the deletion, insertion, and mismatch rates; the number of unaligned and aligned reads; and the identification rate compared

Table 5: Comparison of normalized edit distance with different neural network architectures.

Architecture	Normalized edit distance
3 convolutional layers	0.4007 ± 0.0277
5 convolutional layers	0.3903 ± 0.0230
10 convolutional layers	0.3874 ± 0.0186
3 bidirectional recurrent layers	0.2987 ± 0.0221
5 bidirectional recurrent layers	0.2930 ± 0.0215
3 convolutional layers + 3 bidirectional recurrent layers	0.2011 ± 0.0252
5 convolutional layers + 5 bidirectional recurrent layers	0.2001 ± 0.0177

The normalized edit distance is the edit distance between predicted reads and labeled reads and normalized by segment length.

to the reference genome. The identity rate was calculated as $\frac{\text{number of matched bases}}{\text{number of bases in reference}}$ and is the marker used here for base-calling accuracy.

Assembly identity rate comparison

We assessed the quality of assemblies generated from reads produced by different basecallers. For each basecaller, a *de novo* assembly was generated using only Nanopore reads for the *M. tuberculosis*, *E. coli*, and Lambda Phage genomes. We used Minimap2 [34] and Miniasm [35] to generate a draft genome, then Racon [36] was used to polish on the draft genome for 10 rounds.

Availability of supporting data

The *M. tuberculosis* sequencing data have been deposited in Genbank under project number PRJNA386696. The Human nanopore data were downloaded from <https://github.com/nanopore-wgs-consortium/NA12878>. Supporting data, including training and testing datasets, are available via GigaDB [37].

Availability of supporting source code and requirements

Program and code are available at <https://github.com/haotianteng/chiron> pypi package index 0.3 at <https://pypi.python.org/pypi/chiron>. Chiron is registered in SciCrunch with RRID:SCR_015950. Chiron is available under a Mozilla Public License v2.0. Chiron is built with Tensorflow and requires python 2.7

Abbreviations

CNN: convolutional neural network; CPU: central processing unit; CTC: connectionist temporal classification; ELB: elution buffer; FC: fully connected; GPU:; LSTM: long short-term memory; NCBI: National Center for Biotechnology Information; NN: neural network; ONT: Oxford Nanopore Technologies; ReLU: rectified linear unit; RNN: recurrent neural network.

Competing interests

L.C. is a participant of Oxford Nanopore's MinION Access Programme and received the MinION device, MinION flow cells, and Oxford Nanopore sequencing kits in return for an early access fee deposit. L.C. and M.D.C. received travel and accommodation expenses to speak at an Oxford Nanopore-organized conference.

None of the authors have any commercial or financial interest in Oxford Nanopore Technologies Ltd.

Funding

LC is supported by an NHMRC career development fellowship (GNT1130084). The research is supported by an ARC research grant (DP170102626). MH is supported by a Westpac Future Leaders Scholarship (2016) awarded by the Westpac Bicentennial Foundation.

Author contributions

M.H., M.D.C., and L.C. conceived the study and designed the experimental framework. H.T. designed and implemented the Chiron algorithm. M.D.C., L.C., and T.D. designed and performed the MinION sequencing. H.T. and M.D.C. labeled the training data. H.T. and M.H. ran the performance comparison. H.T. and M.D.C. wrote the initial draft. H.T., M.H., and L.C. refined the manuscript. All authors contributed to editing the final manuscript.

Acknowledgements

We thank Jianhua Guo for contributing the DNA for the *E. coli* sample. We thank Arnold Bainomugisa for extracting DNA for the *M. tuberculosis* sample. We thank Sheng Wang and Han Qiao for the helpful discussion. We thank Jain et al. [14] for the open Human nanopore dataset.

References

1. Kasianowicz JJ, Brandin E, Branton D, et al. Characterization of individual polynucleotide molecules using a membrane channel. *Proc Nat Acad of Sci* 1996;**93**(24):13770–3.
2. Branton D, Deamer DW, Marziali A, et al. The potential and challenges of nanopore sequencing. *Nature Biotechnology* 2008;**26**(10):1146–53.
3. Stoddart D, Heron AJ, Mikhailova E, et al. Single-nucleotide discrimination in immobilized DNA oligonucleotides with a biological nanopore. *Proc Nat Acad of Sci U S A* 2009;**106**(19):7702–7.
4. Ashton PM, Nair S, Dallman T, et al. MinION nanopore sequencing identifies the position and structure of a bacterial antibiotic resistance island. *Nature Biotechnology* 2014;**33**(3):296–300.
5. Cao MD, Ganesamoorthy D, Elliott AG, et al. Streaming algorithms for identification of pathogens and antibiotic resistance.

- tance potential from real-time MinIONTM sequencing. *GigaScience* 2016;5(1):32, 10.1186/s13742-016-0137-2.
6. Cao MD, Nguyen SH, Ganesamoorthy D, et al. Scaffolding and completing genome assemblies in real-time with nanopore sequencing. *Nature Communications* 2017;8:14515, 10.1038/ncomms14515.
 7. Cao MD, Ganesamoorthy D, Cooper MA, et al. Realtime analysis and visualization of MinION sequencing data with npReader. *Bioinformatics* 2016;32(5):764–6.
 8. Quick J, Loman NJ, Duraffour S, et al. Real-time, portable genome sequencing for Ebola surveillance. *Nature* 2016;530(7589):228–32.
 9. Faria NR, Sabino EC, Nunes MR, et al. Mobile real-time surveillance of Zika virus in Brazil. *Genome Medicine* 2016;8(1):97.
 10. McIntyre AB, Rizzardi L, Angela MY, et al. Nanopore sequencing in microgravity. *npj Microgravity* 2016;2:16035.
 11. Castro-Wallace SL, Chiu CY, John KK, et al. Nanopore DNA sequencing and genome assembly on the International Space Station. *Scientific Reports* 2017;p. 18022.
 12. Graves A, Fernández S, Gomez F, et al. Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks. In: *Proceedings of the 23rd International Conference on Machine Learning ACM*; 2006. p. 369–376.
 13. Stobier M, Brown J, BasecRAWller: Streaming Nanopore Basecalling Directly from Raw Signal. *bioRxiv* 2017;:133058.
 14. Jain M, Koren S, Quick J, et al. Nanopore sequencing and assembly of a human genome with ultra-long reads. *Nature Biotechnology* 2018, 338–345. .
 15. Laszlo AH, Derrington IM, Ross BC, et al. Decoding long nanopore sequencing reads of natural DNA. *Nature Biotechnology* 2014;32(8):829–833.
 16. David M, Dursi LJ, Yao D, et al. Nanocall: an open source basecaller for Oxford Nanopore sequencing data. *Bioinformatics* 2016;33(1):49–55.
 17. Boža V, Brejová B, Vinař T. DeepNano: deep recurrent neural networks for base calling in MinION nanopore reads. *PLoS ONE* 2017, 12, 6(0): e0178751 .
 18. Wick RR, Judd LM, Holt KE. Comparison of Oxford Nanopore Basecalling Tools; 2017. Available from: <https://doi.org/10.5281/zenodo.1082696>.
 19. Han S, Mao H, Dally WJ. Deep compression: compressing deep neural networks with pruning, trained quantization and Huffman coding arXiv preprint arXiv:151000149. 2015.
 20. LeCun Y, Bengio Y, Hinton G. Deep learning. *Nature* 2015;521(7553):436–44.
 21. Amodei D, Anubhai R, Battenberg E, et al. Deep Speech 2: end-to-end speech recognition in English and Mandarin. In: *International Conference on Machine Learning*; 2016. p. 173–182.
 22. Ioffe S, Szegedy C. Batch normalization: accelerating deep network training by reducing internal covariate shift. arXiv preprint arXiv:150203167. 2015.
 23. Schuster M, Paliwal KK. Bidirectional recurrent neural networks. *IEEE Transactions on Signal Processing* 1997;45(11):2673–81.
 24. Hochreiter S, Schmidhuber J. Long short-term memory. *Neural Computation* 1997;9(8):1735–80.
 25. Cooijmans T, Ballas N, Laurent C et al. Recurrent batch normalization arXiv preprint arXiv:160309025. 2016.
 26. Kingma D, Ba J. Adam: a method for stochastic optimization arXiv preprint arXiv:1412.6980. 2014.
 27. Graves A, Jaitly N. Towards end-to-end speech recognition with recurrent neural networks. In: *Proceedings of the 31st International Conference on Machine Learning (ICML-14)*; 2014. p. 1764–1772.
 28. He K, Zhang X, Ren S, et al. Deep residual learning for image recognition. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*; 2016. p. 770–778.
 29. Nair V, Hinton GE. Rectified linear units improve restricted boltzmann machines. In: *Proceedings of the 27th International Conference on Machine Learning (ICML-10)*; 2010. p. 807–814.
 30. Gers FA, Schmidhuber J, Cummins F. Learning to forget: continual prediction with LSTM. *Neural Computation* 2000;12(10):2451–71.
 31. Bankevich A, Nurk S, Antipov D, et al. SPAdes: a new genome assembly algorithm and its applications to single-cell sequencing. *Journal of Computational Biology* 2012;19(5):455–77.
 32. Stoiber MH, Quick J, Egan R, et al. De novo identification of DNA modifications enabled by genome-guided nanopore signal processing. *bioRxiv* 2017;p. 094672.
 33. Sović I, Šikić M, Wilm A, et al. Fast and sensitive mapping of nanopore sequencing reads with GraphMap. *Nature Communications* 2016;7:11307.
 34. Li H. Minimap2: versatile pairwise alignment for nucleotide sequences arXiv. 2017;1708.
 35. Li H. Minimap and Miniasm: fast mapping and de novo assembly for noisy long sequences. *Bioinformatics* 2016;32(14):2103–10.
 36. Vaser R, Sović I, Nagarajan N, et al. Fast and accurate de novo genome assembly from long uncorrected reads. *Genome Research* 2017;27(5):737–746.
 37. Teng H, Cao MD, Hall MB, et al. Supporting data for “Chiron: translating nanopore raw signal directly into nucleotide sequence using deep learning.” *GigaScience Database* 2018, <http://dx.doi.org/10.5524/100424>.