

# SquiggleKit: supplementary information

## Contents:

### [Contents:](#)

#### [1. Detailed description of finding 3' end motif example](#)

#### [2. Comparison to accessible methods with related functions](#)

##### [Table 1: Overview of tool comparisons](#)

#### [3. Extended comparison](#)

##### [BulkVis](#)

##### [Porettools](#)

##### [Picopore](#)

##### [PoRe](#)

##### [HDFView](#)

##### [Japsa](#)

##### [ReadUntil](#)

##### [Tombo/NanoRaw](#)

##### [Nanopolish](#)

#### [4. Standardisation of signal for comparison](#)

##### [Sup. Figure 1: Comparison of raw signal scaling strategies on MotifSeq scoring](#)

#### [5. Benchmarking](#)

##### [System specifications:](#)

##### [Fast5 fetcher speed](#)

##### [Table 2: Fast5 fetcher speed benchmark](#)

##### [Segmenter speed](#)

##### [Table 3: Segmenter speed \(default parameters on cDNA sequins\)](#)

##### [MotifSeq speed](#)

##### [Table 4: MotifSeq speed in function of k-mer length \(default parameters on .tsv files\)](#)

##### [Segmenter accuracy](#)

##### [Sup. Figure 2: Segmenter boundary assessment](#)

##### [Table 5: Segmenter Accuracy](#)

##### [MotifSeq accuracy](#)

##### [Sup. Figure 3: MotifSeq scoring benchmark](#)

#### [6. Supplementary References](#)

# 1. Detailed description of finding 3' end motif example

Note: The use of Parallel Alignment File (PAF) format over sam format is arbitrary. It does however only keep reads that aligned, whereas sam includes all reads, thus streamlining fast5 filtering.

1. Aligning reads with *minimap2* with paf output:

```
$ minimap2 -x ont-map sequin.fa reads.fastq > alignment.paf
```

2. Filtering paf file for reads of interest, where R2\_151 represents a unique identifier for the transcript group of interest.

```
$ grep R2_151 alignment.paf > filtered.paf
```

3. Using a filtered .paf file (generated in 2.) as input for *fast5\_fetcher*:

```
$ mkdir fast5  
$ python fast5_fetcher.py -p filtered.paf -s sequencing_summary.txt  
-i name.index -o ./fast5
```

4. Extract signals from fast5 files to a tab separated values (tsv) file:

```
$ python SquigglePull.py -rv -p ./fast5/ -f all > signals.tsv
```

5. Plot signals:

```
$ python SquigglePlot.py -s signals.tsv --save R2_151.pdf --  
save_path ./pics/ --no_show
```

6. Identify any segments in signals and visualise each one:

```
$ python segmenter.py -s signals.tsv -v > segments.tsv
```

7. Motif identification. Fasta format for creating signal from model (in a custom file named *motif.fa* for this example):

```
>my_kmer_name  
CATCTATCCAGGGTTAAATT
```

8. Find the best match to that k-mer in the signal, and visualise it:

```
$ python MotifSeq.py -s signals.tsv -i motif.fa -v >  
signals_kmer.tsv
```

## 2. Comparison to accessible methods with related functions

Table 1: Overview of tool comparisons

Tools	<i>SquiggleKit</i>				
	<i>fast5_fetcher</i>	<i>SquigglePull</i>	<i>SquigglePlot</i>	<i>Segmenter</i>	<i>MotifSeq</i>
<i>BulkVis</i>					
<i>Poretools</i>					
<i>Picopore</i>					
<i>poRe</i>					
<i>HDFView</i>					
<i>Japsa</i>					
<i>ReadUntil</i>					
<i>Tombo/NanoRaw</i>					
<i>Nanopolish</i>					

### Legend

Not the same	Similar functionality	Equal functionality

### Extraction and Visualisation:

- *Poretools* - extract events, not raw (Loman and Quinlan, 2014)
- *poRe* - extract events, not raw (Watson *et al.*, 2015)
- *Nanoraw* - raw signal visualisation and comparison, aligns signal to base
- *Tombo* - raw signal visualisation and comparison, aligns signal to base (Stoiber *et al.*, 2017)
- *HDFView* - single file, basic visualisation (<http://www.hdfgroup.org/HDF5/>)
- *BulkVis* - bulk fast5 files, real time (Payne *et al.*, 2018)

### File Management:

- *Picopore* - removes data from fast5 files to reduce size (Gigante, 2017)
- *Japsa* - manages files, streams data and performs base space analysis (Nguyen *et al.*, 2017)

**Analysis:**

- *ReadUntil* - local dynamic programming, real time (Loose *et al.*, 2016)
- *Nanopolish* - File management, segmentation, aligns signal to base (Loose *et al.*, 2016; Loman *et al.*, 2015; Loman and Quinlan, 2014)

**3. Extended comparison**

Comparison	<b>BulkVis</b>
Info	<i>BulkVis</i> enables visualization of bulk FAST5 files collected from Nanopore sequencers. <i>BulkVis</i> is provided for the visual inspection of challenging or difficult to sequence samples or where the user wishes to investigate specific events during a run. In these instances analysis of a bulk FAST5 file may provide some visual indication of the underlying issues
Last update	September 2018
Links	<a href="https://academic.oup.com/bioinformatics/advance-article/doi/10.1093/bioinformatics/bty841/5193712">https://academic.oup.com/bioinformatics/advance-article/doi/10.1093/bioinformatics/bty841/5193712</a> <a href="https://bulkvis.readthedocs.io/en/latest/index.html">https://bulkvis.readthedocs.io/en/latest/index.html</a>
Similarities	Parses sequencing_summary.txt, .paf, and bulk fast5 files to break down MinKNOW states: above, adapter, pore, transition, unblocking, and unclassified. Outputs in a .csv. This is similar to the adapter stall targeting by segmenter, however the segmenter algorithm is novel.  It allows for great visualisation and annotation of signals in a bulk fast5 file. This is similar to <i>SquigglePlot</i> , but are used for different applications
Differences	<i>SquigglePlot</i> and segmenter work with regular fast5 files, or any signal input, where <i>BulkVis</i> plots bulk fast5 files.
Comments and use	It has a specific use case around examining reads using MinKNOW's classifications and requires bulk fast5 file.

Comparison	<b>Poretools</b>
Info	<i>Poretools</i> provides a wealth of utilities and data exploration tools.

Last update	10 Jul 2017
Links	<a href="https://github.com/arq5x/porettools">https://github.com/arq5x/porettools</a>
Similarities	<i>Porettools</i> contains event data extraction similar to that of <i>SquigglePull</i>
Differences	<i>Porettools</i> extracts event data, which is no longer used as the basis of basecalling. Events changed to a static time for sampling and changed in format, breaking <i>Porettools</i> for the latest data.
Comments and use	<i>Porettools</i> was one of the first tools used for generating stats and getting a feel for nanopore data. It is no longer routinely used, and does not have support for raw signal data. If the user is after raw signal data and visualisation, picking <i>SquigglePull</i> and <i>SquigglePlot</i> would be the superior choice

<b>Comparison</b>	<b>Picopore</b>
Info	A tool for reducing the size of Oxford Nanopore Technologies' datasets without losing information.
Last update	31 Aug 2017
Links	<a href="http://dx.doi.org/10.12688/f1000research.11022.1">http://dx.doi.org/10.12688/f1000research.11022.1</a> <a href="https://github.com/scottgigante/picopore">https://github.com/scottgigante/picopore</a>
Similarities	Reduction in memory footprint, but by gutting fast5 files of various information. This is similar to extracting the raw signal with <i>SquigglePull</i> and gzipping. The reduction in memory footprint also overlaps with the reasoning behind <i>fast5_fetcher</i> , downsizing to only mapped fast5 files, or fast5 files of interest
Differences	The methodology of memory footprint for storage or use are quite different, where <i>picopore</i> guts and compresses data, <i>fast5_fetcher</i> filters and <i>SquigglePull</i> extracts into a different format.
Comments and use	<i>Picopore</i> is no longer under active development. Due to improvements in ONT's native HDF5 compression, lossless and deep-lossless compression no longer effectively reduce the size of nanopore files. With the increase in yields, and lack of methods to compress that data, both <i>fast5_fetcher</i> and <i>SquigglePull</i> can be of use to reduce

	what data you need depending on the kind of analysis you are trying to undertake.
--	---

<b>Comparison</b>	<b>PoRe</b>
Info	This package enables organisation and visualisation of MinION data on the MinION laptop for the novice user.
Last update	31 Aug 2017
Links	<a href="https://github.com/mw55309/poRe_docs">https://github.com/mw55309/poRe_docs</a> <a href="https://sourceforge.net/projects/rpore/files/">https://sourceforge.net/projects/rpore/files/</a> <a href="https://github.com/mw55309/poRe_scripts">https://github.com/mw55309/poRe_scripts</a>
Similarities	<i>poRe</i> contains event data extraction similar to that of <i>SquigglePull</i>
Differences	<i>poRe</i> does not extract or visualise raw signal
Comments and use	<i>poRe</i> is no longer under development. It works with legacy fast5 files and older basecaller data structures.

<b>Comparison</b>	<b>HDFView</b>
Info	<i>HDFView</i> is a visual tool written in Java for browsing and editing HDF5 files. View a file hierarchy in a tree structure. Create new files, add or delete groups and datasets. View and modify the content of a dataset. Add, delete and modify attributes.
Last update	updated current 2019
Links	<a href="http://www.hdfgroup.org/HDF5/">http://www.hdfgroup.org/HDF5/</a>
Similarities	<i>HDFView</i> allows for basic visualisation of raw signal data, with overlap with <i>SquigglePlot</i>
Differences	<i>SquigglePlot</i> has many more features than <i>HDFView</i> for visualisation, and can plot many files at once with custom settings, at high DPI for figures.
Comments and use	<i>HDFView</i> is excellent for exploring a fast5 file for information and structure or visualising one signal. When working with many signals, or just wanting to look at the signal directly, <i>SquigglePlot</i> is superior.

<b>Comparison</b>	<b>Japsa</b>
Info	<i>Japsa</i> has many packages that do a variety of nanopore data analysis and data streaming.
Last update	update current 2019 - multi package project
Links	<a href="https://japsa.readthedocs.io/en/latest/index.html">https://japsa.readthedocs.io/en/latest/index.html</a> <a href="https://www.ncbi.nlm.nih.gov/pubmed/28961965">https://www.ncbi.nlm.nih.gov/pubmed/28961965</a>
Similarities	<i>Japsa</i> has packages for real time file management and analysis during sequencing.
Differences	<i>Japsa</i> employs many base space sequence analysis methods like barcode demultiplexing. It does not however do any analysis, extraction, or visualisation of raw signal
Comments and use	<i>Japsa</i> seems to have been built as a pipeline for many real time analysis methods. It does not have any significant overlap other than file some file management.

<b>Comparison</b>	<b>ReadUntil</b>
Info	<i>ReadUntil</i> has a function, <code>squiggle_search2()</code> , to find if a particular sequencing matches a selected region or not, using the coordinates/position of the match. This is specifically designed for <i>ReadUntil</i> , and the scripts and methods can not be used for general exploration. No plotting of raw signal.
Last update	2016
Links	<a href="https://www.nature.com/articles/nmeth.3930">https://www.nature.com/articles/nmeth.3930</a> <a href="http://mattloose.github.io/RUscriptsdocs/">http://mattloose.github.io/RUscriptsdocs/</a> <a href="https://github.com/mattloose/RUscripts/tree/master/ReadUntil">https://github.com/mattloose/RUscripts/tree/master/ReadUntil</a>
Similarities	<i>ReadUntil</i> takes a reference sequence region and converts it into the event space using the models for basecalling. It then uses dynamic programming on the first portion of an incoming read to test if it matches within the selected region.
Differences	<i>MotifSeq</i> is a general method for finding the approximate position of where a particular sequence motif aligns with the raw signal data.

	<i>ReadUntil</i> and <i>MotifSeq</i> do this in similar ways, however while <i>ReadUntil</i> does the comparisons in event space with streaming data, <i>MotifSeq</i> does the comparisons in raw signal space.
Comments and use	<code>squiggle_search2()</code> from <i>ReadUntil</i> finds if a particular sequence matches a selected region or not, using the coordinates/position of the match. This is specifically designed for <i>ReadUntil</i> , and the scripts and methods can not be used for general exploration as is.

<b>Comparison</b>	<b>Tombo/NanoRaw</b>
Info	<i>Nanoraw</i> and <i>tombo</i> operate based on the <code>genome_resquiggle/resquiggle</code> method, mapping the raw signal with basecalls and reference alignment. All visualisation is centred around this method, and though it has the potential to be general use, it is not implemented or designed that way.
Last update	last updated current in 2019
Links	<a href="https://www.biorxiv.org/content/biorxiv/early/2017/04/10/094672.full.pdf">https://www.biorxiv.org/content/biorxiv/early/2017/04/10/094672.full.pdf</a> <a href="https://nanoporetech.github.io/tombo/tutorials.html">https://nanoporetech.github.io/tombo/tutorials.html</a> <a href="https://nanoraw.readthedocs.io/en/latest/resquiggle.html">https://nanoraw.readthedocs.io/en/latest/resquiggle.html</a>
Similarities	<p><i>Tombo</i> can visualise raw signal data similar to <i>SquigglePlot</i>.</p> <p><i>Tombo</i> can centre the visualisation around a particular motif, and extract that information with the API to find the signal associated with the basecalls, similar to <i>MotifSeq</i>.</p> <p><i>Tombo</i> can extract the signal data, similar to <i>SquigglePull</i>.</p>
Differences	<p><i>Tombo</i> can visualise multiple squiggles on one plot, and can associate them with the basecall as well creating a squiggle pileup plot.</p> <p>While <i>Tombo</i> focuses on multiple squiggles, <i>SquigglePlot</i> focuses on one squiggle at a time.</p> <p><i>Tombo</i> uses the <code>resquiggle</code> algorithm to match the aligned base sequence to the raw signal.</p> <p><i>MotifSeq</i> converts the sequence motif to a signal using a model, and uses a local dynamic programming method to identify the position of the motif.</p> <p>The signal extraction in <i>SquigglePull</i> is more general than <i>tombo</i>'s pipeline.</p>



Comments and use	<p><i>Tombo</i> is a fully fledged toolkit for finding differences in signal space, and exploring the raw signal of a dataset.</p> <p><i>Tombo</i> and <i>NanoRaw</i> were designed to find differences in the signal to identify modifications, and help train models to detect them.</p> <p><i>Fast5_fetcher</i> may be of use for filtering large datasets to be analysed with <i>tombo</i></p> <p>Use of either toolkit would be based on specific goals</p>
------------------	--

<b>Comparison</b>	<b>Nanopolish</b>
Info	<i>Nanopolish</i> is a software package for signal-level analysis of Oxford Nanopore sequencing data. <i>Nanopolish</i> can calculate an improved consensus sequence for a draft genome assembly, detect base modifications, call SNPs and indels with respect to a reference genome and more.
Last update	current 2019
Links	<a href="https://github.com/jts/nanopolish">https://github.com/jts/nanopolish</a> <a href="https://nanopolish.readthedocs.io/en/latest/index.html">https://nanopolish.readthedocs.io/en/latest/index.html</a>
Similarities	File management, segmentation, aligns signal to base are similar to <i>fast5_fetcher</i> , <i>segmenter</i> , and <i>MotifSeq</i> respectively.
Differences	<i>Nanopolish</i> is aimed at specific problems and the tool is created as an integrated method of analysis, where <i>SquiggleKit</i> tools can all work standalone and in a more general way.
Comments and use	<i>Nanopolish</i> is best suited to standard methods of analysis. <i>SquiggleKit</i> is best used for developing tools like <i>nanopolish</i> , or exploring nanopore signal data.

## 4. Standardisation of signal for comparison

Many factors can contribute to non-sequence dependent variability in nanopore sequencing data (temperature, voltage, etc), which can lead to progressive distortion of raw current values for a single read. The three main types of raw signal distortion are shift, scale, and drift, which can be mitigated by normalising or standardising raw signal using global properties of the data. Standardisation is required to compare 2 raw signals against each other, or a simulated signal from the basecalling model and a raw signal, as implemented in *MotifSeq*.

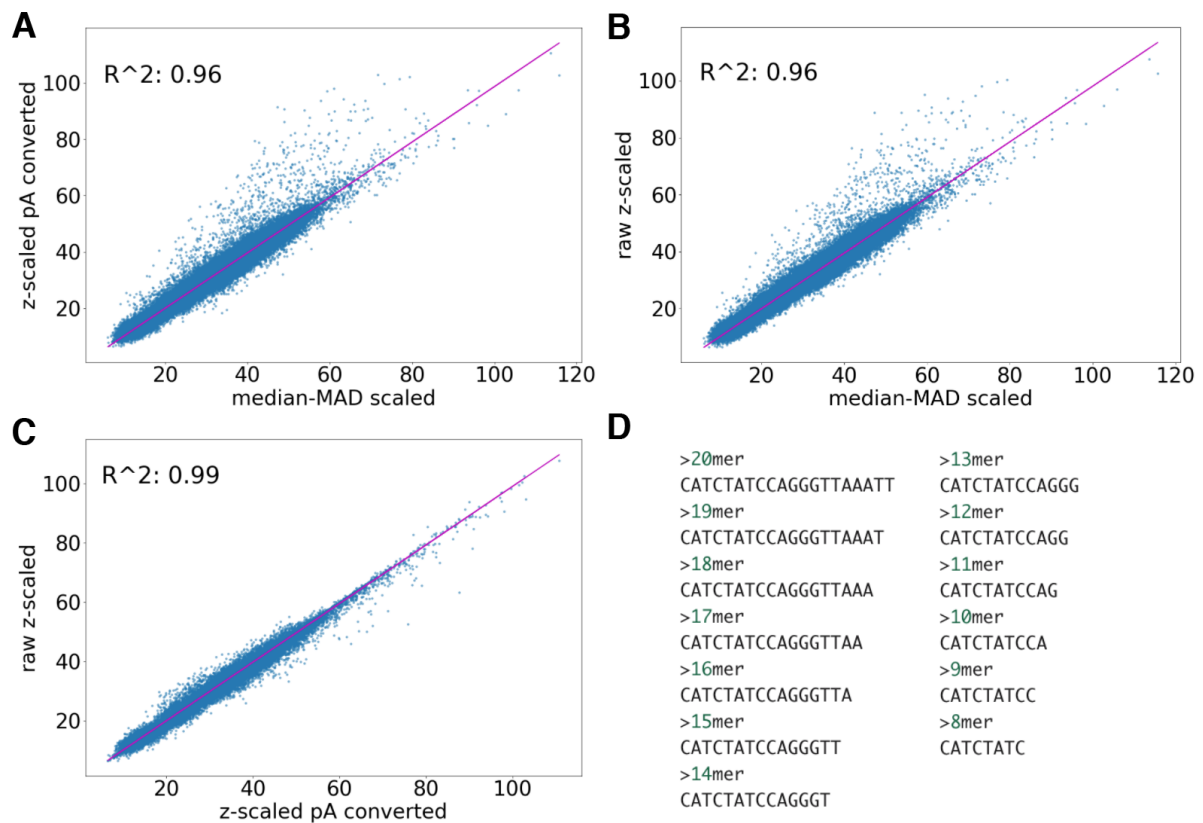
We compared 3 methods for raw signal normalisation: Z-score scaling, Median-Median-Absolute-Difference scaling (med-MAD), and pico-Ampere scaling.

Z-score scaling is performed using the function  $z = (x - u) / s$ , where **u** and **s** are the mean and standard deviation of the signal, respectively. It is more sensitive to any large spikes in the signal (which are relatively common), thus impacting the overall mean and standard deviation. One method of offsetting this is to filter large spikes above and below a threshold, as is done with the `--scale_hi/--scale_low` flags in *MotifSeq*, set at 900 and 0 respectively for DNA.

Alternatively, raw signal can be shifted and scaled by converting to pico ampres, then scaled. This is done by extracting the digitisation, offset, and range values from the fast5 files, and converting the raw data points using the function  $pA = (x + offset) * (range/digitisation)$ . This can be done automatically when the raw signal is extracted using the `--convert` flag in *SquigglePull*. There is no significant difference in *MotifSeq* scores when using Z-score or Z-score pA scaling (**Supplementary Figure 1C**).

*MotifSeq* can do Z-score or median-absolute-difference scaling, in a similar way, but which are more or less sensitive to large spikes in the signal. The med-MAD scaling uses the following function:  $z = (x - med) / MAD$ , where **med** and **MAD** are the median and Median-Absolute-Difference of the signal respectively. Med-MAD is more robust at handling noise in nanopore signals than z-score scaling, where the mean and standard deviation can be significantly impacted by the occasional spike in current. Med-MAD is used by nanopore basecalling software for this reason. This is evidenced in **Supplementary Figure 1A-B** by slightly lower *MotifSeq* scores, indicative of lower penalties for outlier spikes in current.

All 3 normalisation methods have been implemented in *MotifSeq*, which uses med-MAD by default given the benchmarking results described below (see section 5).



Sup. Figure 1: Comparison of raw signal scaling strategies on MotifSeq scoring

**(A-C)** Pairwise comparison of MotifSeq scores for all combinations of raw Z-score scaling (raw z-scaled), raw Z-score scaling with conversion to picoAmperes (z-scaled pA converted), and median-median-absolute-difference (median-MAD scaled). Scores were generated by running MotifSeq with 12 different k-mers **(D)** corresponding to the 3' end used in the main example against the raw signals of 10,000 randomly selected cDNA reads, which were scaled using either of the 3 methods described above.

## 5. Benchmarking

Benchmarking was carried out to assess speed and accuracy for certain *SquiggleKit* tools. Below is a comparison of speed for *fast5\_fetcher*, *segmenter*, and *MotifSeq*. Accuracy benchmarking was performed for *Segmenter* and *MotifSeq*. All local benchmarks were carried out on the same hardware, as described below.

### System specifications:

Linux omega-1-31.local 2.6.32-642.3.1.el6.x86\_64  
x86\_64 x86\_64 x86\_64 GNU/Linux  
CPU: AMD Opteron(tm) Processor 6282 SE (only 1 CPU used)  
CPU MHz: 2600  
Architecture: x86\_64  
RAM: 32GB  
Disk: 400GB NVME  
Python version: 2.7.14

### Fast5\_fetcher speed

By taking a dataset with 7,332,202 single fast5 files, tarballed into groups of 4000, the average runtime (after 5 replicates) of *fast5\_fetcher* extraction was measured at different magnitudes (i.e. number of files to extract).

Table 2: Fast5\_fetcher speed benchmark

Number of files	Size of files	Average runtime (s)	Reads per Second
1,000	32MB	204.2	4.9
10,000	697MB	1673.8	6.0
100,000	9.7GB	20133.8	5.0
1,000,000	102GB	259546	3.9

### Segmenter speed

*Segmenter* identifies regions of low complexity using a greedy algorithm in a manner similar to a Markov chain. It processes raw signal as follows:

1. Calculate the median current intensity of the full read;
2. Set thresholds about the median using a fraction of the standard deviation (parameter `--std_scale`, default 0.75);
3. For each data point, calculate the difference to the median. Stretches of signal are returned if at least  $w$  (parameter `--window`, default 150) consecutive data points are within the standard deviation threshold;
4. Outliers are tolerated (parameter `--error`, default 5) as are consecutive low complexity regions within a given boundary (parameter `--seg_dist`, default 50).

Table 3: Segmenter speed (default parameters on cDNA sequins)

Number of files	time (s)	Reads per second	Time for gzip input (s)	Reads per second (gzip)
103,526	916	~113	3,430	~30

## MotifSeq speed

Table 4: MotifSeq speed in function of k-mer length (default parameters on .tsv files)

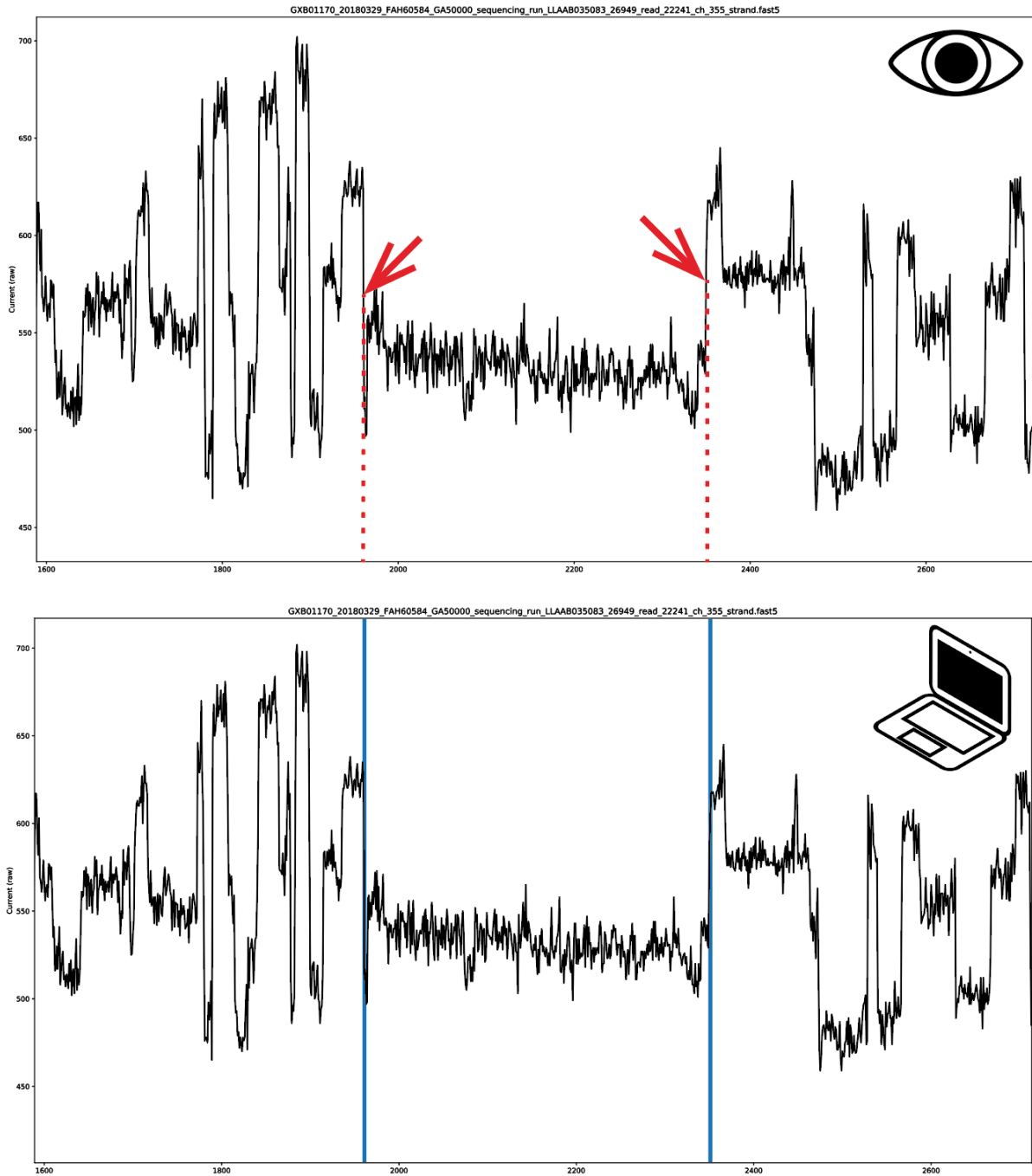
K-mer length	Total time (s)	Reads per second	Time per read (s)
12	92.13	54.27	0.0184
30	143.82	34.77	0.0288
50	232.05	21.55	0.0464
75	338.47	14.77	0.067
100	402.98	12.41	0.080
256	963.55	5.19	0.193

This is comparable to the speed attained by (Loose *et al.*, 2016) in their speed testing on “Intel(R) Xeon(R) E5-2690 version 3 central processing units (CPUs) running at 2.60 GHz (server)” attaining ~0.3 seconds per read for a k-mer length of 256.

## Segmenter accuracy

Segmenter is designed to identify low complexity regions in raw nanopore signal above a given window length. It was designed as a qualitative annotation tool to identify homopolymer boundaries and polymer translocation stalling. Although it could be used to find associations between poly-A tail lengths and signal length, this is not the subject of this toolkit.

Identifying low complexity regions in raw signal is not straightforward to benchmark given the stochastic nature of single molecule sensing. Consequently, we evaluated the accuracy of *Segmenter* by manually inspecting 100 raw cDNA signals using *Segmenter*'s “-v” argument (**Supplementary Figure 2**). Homopolymer boundaries corresponding to poly-A regions in cDNA reads were visually identified, their position recorded and compared to the those automatically identified by *Segmenter* (**Table 5**).



## Sup. Figure 2: Segmenter boundary assessment

Low complexity regions identified by *Segmenter* were manually inspected (top) to identify precise boundaries visually using *matplotlib* functionality. The positions on the x-axis corresponding to boundary signals (red lines) were recorded and compared to those output by *Segmenter* (bottom, blue lines),

Table 5: Segmenter Accuracy

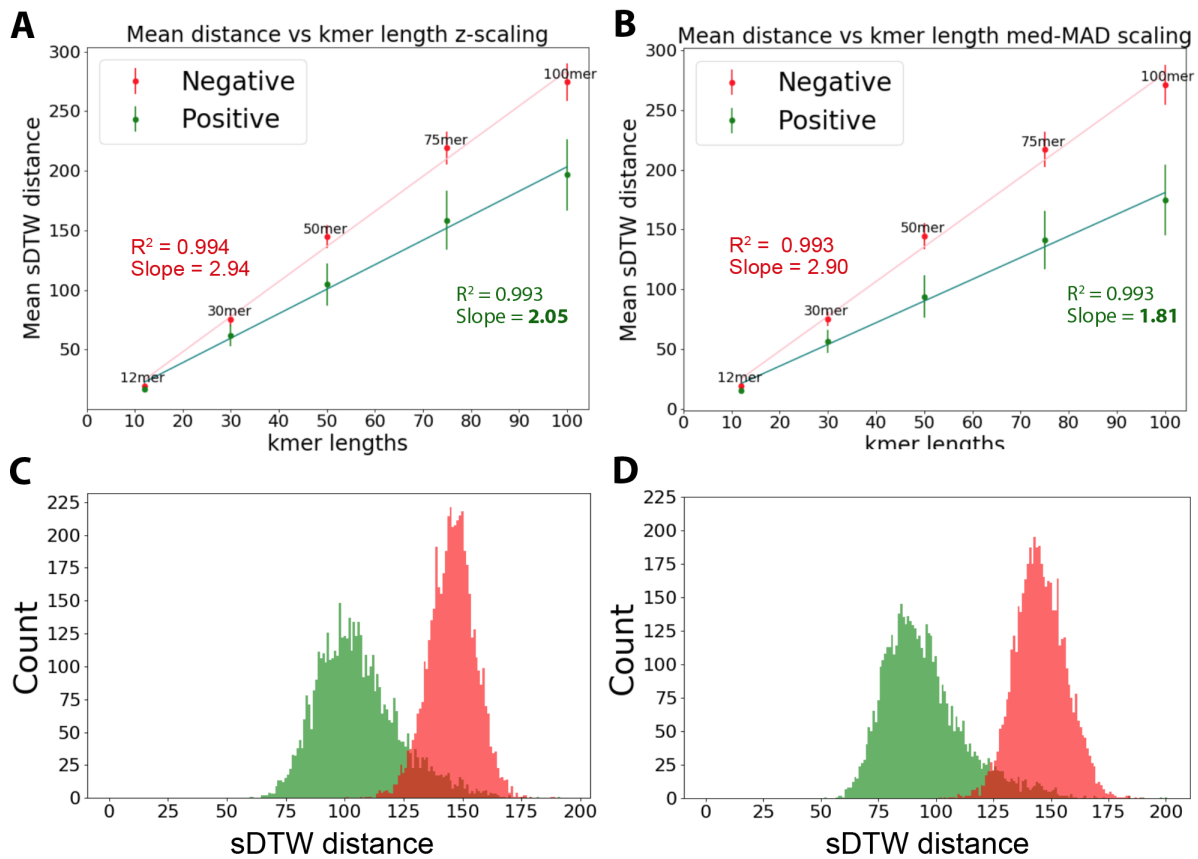
Differences	Mean difference	Median difference	Stdev of difference
Start	2.32	1.0	5.4
End	9.25	3.0	13.7
Length	9.75	3.0	14.7

As can be seen, start accuracy is quite good, with the largest variation in the end position targeting. A mean length difference of 9.75 data points approximately corresponds to 1 nucleotide worth of signal (~8-11 data points per nucleotide for DNA pore models).

### MotifSeq accuracy

We benchmarked *MotifSeq* accuracy by comparing minimum dynamic programming scores across positive and negative control datasets, which were generated as follows:

- (i) Using a sequin synthetic RNA spike-in control cDNA sequencing run (ENA accession number PRJEB33439), we randomly selected a sequin isoform from all isoforms with at least 10,000 uniquely mapping reads (minimap2 using the preset “-x map-ont”) and a MAPQ score of 60 via the resulting .paf file.
- (ii) The corresponding reference transcript (available at [www.sequin.xyz](http://www.sequin.xyz)) was used to randomly extract k-mers (within 15nt of the transcript extremities) of lengths 12, 30, 50, 75, and 100 nucleotides.
- (iii) For each k-mer, positive and negative control sets were generated by identifying all reference sequin isoforms with or without the selected k-mer, respectively, using direct sequence matching.
- (iv) A random selection of 5000 cDNA reads were selected that uniquely map to the reference isoforms for both control sets defined in (iii). We ensure the mapped sequences overlapped the position of the k-mer identified in the reference sequences. Reads were selected with *fast5\_fetcher* and raw signals extracted using *SquigglePull*.
- (v) Each k-mer motif was submitted to *MotifSeq* comparisons across both positive and negative control datasets (**Supplementary Figure 3**).



### Sup. Figure 3: MotifSeq scoring benchmark

The *MotifSeq* mean score (sDTW distance) and standard deviation (error bars) for 5000 comparisons against positive (red, top) and negative (green, bottom) controls is plotted for 5 different k-mer sizes in conjunction with **(A)** z-scaling and **(B)** med-MAD scaling of raw signal. A linear regression trend line is fitted to 5 different k-mer sizes with associated slopes and correlation of determination ( $R^2$ ). **(C)** Distribution of *MotifSeq* scores for 50-mers using z-scaling and **(D)** med-MAD scaling.

As expected, the separation of true positive scores from true negative scores increases with k-mer length. The separation further increases when using med-MAD scaling vs z-scaling, highlighting the improved tolerance of med-MAD scaling for outlier current spikes through slightly lower scores overall. *MotifSeq* scores against the negative control appear to be normally distributed (**Supplementary Figure 3C-D**, in red or on the right) suggesting that they can be interpreted as Z-scores to assess the confidence of hits via the background mean and standard deviation, which can be converted to a probability using the formula:

$$P(X < \mu) = P\left(\frac{X - \mu}{\sigma} < \frac{\mu_q - \mu_b}{\sigma_b}\right)$$

Where  $\mu$  is the mean,  $\mu_q$  is the query mean,  $\mu_b$  is the background mean,  $\sigma_b$  is background standard deviation.



Furthermore, the seemingly linear correlation between *MotifSeq* scores and k-mer lengths suggests that linear regression can be used to extrapolate background models to a given motif size (**Table 6**). Consequently, we incorporated this additional scoring metric into *MotifSeq*, which now

Table 6: Expected background mean and standard deviation by k-mer length

<b>k-mer</b>	<b>Mean</b>	<b>stdev</b>	<b>k-mer</b>	<b>Mean</b>	<b>stdev</b>
<b>10</b>	19.44	1.65	<b>160</b>	455.16	38.54
<b>15</b>	33.96	2.88	<b>165</b>	469.68	39.77
<b>20</b>	48.49	4.11	<b>170</b>	484.21	41
<b>25</b>	63.01	5.34	<b>175</b>	498.73	42.23
<b>30</b>	77.54	6.57	<b>180</b>	513.26	43.46
<b>35</b>	92.06	7.8	<b>185</b>	527.78	44.69
<b>40</b>	106.58	9.03	<b>190</b>	542.3	45.92
<b>45</b>	121.11	10.26	<b>195</b>	556.83	47.15
<b>50</b>	135.63	11.49	<b>200</b>	571.35	48.38
<b>55</b>	150.16	12.72	<b>205</b>	585.87	49.61
<b>60</b>	164.68	13.95	<b>210</b>	600.4	50.84
<b>65</b>	179.2	15.18	<b>215</b>	614.92	52.07
<b>70</b>	193.73	16.4	<b>220</b>	629.45	53.3
<b>75</b>	208.25	17.63	<b>225</b>	643.97	54.53
<b>80</b>	222.78	18.86	<b>230</b>	658.49	55.76
<b>85</b>	237.3	20.09	<b>235</b>	673.02	56.99
<b>90</b>	251.82	21.32	<b>240</b>	687.54	58.22
<b>95</b>	266.35	22.55	<b>245</b>	702.07	59.45
<b>100</b>	280.87	23.78	<b>250</b>	716.59	60.68
<b>105</b>	295.4	25.01	<b>255</b>	731.11	61.91
<b>110</b>	309.92	26.24	<b>260</b>	745.64	63.14
<b>115</b>	324.44	27.47	<b>265</b>	760.16	64.37
<b>120</b>	338.97	28.7	<b>270</b>	774.69	65.6
<b>125</b>	353.49	29.93	<b>275</b>	789.21	66.83
<b>130</b>	368.02	31.16	<b>280</b>	803.73	68.06
<b>135</b>	382.54	32.39	<b>285</b>	818.26	69.29
<b>140</b>	397.06	33.62	<b>290</b>	832.78	70.52
<b>145</b>	411.59	34.85	<b>295</b>	847.31	71.75
<b>150</b>	426.11	36.08	<b>300</b>	861.83	72.98

155	440.64	37.31			
-----	--------	-------	--	--	--

\*Median-MAD scaling was used for modelling.

## 6. Supplementary References

- Gigante,S. (2017) Picopore: A tool for reducing the storage size of Oxford Nanopore Technologies datasets without loss of functionality. *F1000Res.*, **6**, 227.
- Loman,N.J. *et al.* (2015) A complete bacterial genome assembled *de novo* using only nanopore sequencing data. *Nat. Methods*, **12**, 733.
- Loman,N.J. and Quinlan,A.R. (2014) Poretools: a toolkit for analyzing nanopore sequence data. *Bioinformatics*, **30**, 3399–3401.
- Loose,M. *et al.* (2016) Real-time selective sequencing using nanopore technology. *Nat. Methods*, **13**, 751.
- Nguyen,S.H. *et al.* (2017) Real-time demultiplexing Nanopore barcoded sequencing data with npBarcode. *Bioinformatics*, **33**, 3988–3990.
- Payne,A. *et al.* (2018) BulkVis: a graphical viewer for Oxford nanopore bulk FAST5 files. *Bioinformatics*.
- Stoiber, M.H. (2017) De novo Identification of DNA Modifications Enabled by Genome-Guided Nanopore Signal Processing, bioRxiv, <https://doi.org/10.1101/094672> .
- Watson,M. *et al.* (2015) poRe: an R package for the visualization and analysis of nanopore sequencing data. *Bioinformatics*, **31**, 114–115.