

TOAST2018

Day 2

Bioinformatics - Nanopore Sequence analysis

19th-20th March 2018
University of Technology Sydney

Tim Kahlke
tim.kahlke@uts.edu.au
Github: <https://github.com/timkahlke>
Twitter: @AdvancedTwigTec

Working environment

Login details

VirtualBox

If you are working with the provided VirtualBox image your login details are

User: toast2018

Password toast2018

Course data

Practicals

The directory for the course data can be found in your home directory in the folder *TOAST2018*. Precompiled data can additionally be re-downloaded from the google drive if needed.

Sequencing data *day 1*

The sequencing data from Day1 can be downloaded from the google drive link provided during the course.

Practical 1

Data analysis

&

Quality control

Data analysis and QC

The initial step of your data analysis is to get an overview of your data, e.g. what is the number of reads, their length, quality etc.

An increasing number of tools is available to analyse nanopore data. In this tutorial we will use two of the most common tools *nanoplot* and *poretools*.

Plot your base called reads with nanoplot

Nanoplot is a command-line tool written in python that can be used to visualise your nanopore reads and read alignments to a reference.

Analyse the sequencing data from day 1: find the directory that contains the fastq files that passed the quality control of the MinKNOW base call. For example: if you downloaded yesterdays sequencing data into the directory `~/TOAST2018/day1/MinKNOW` it will be located in `~/TOAST2018/day1/MinKNOW/reads/20180319_toast_GROUP/fastq/pass`.

Change into directory `TOAST/day2/prac1` and create a directory to store the nanoplot output in, e.g. `prac1_nanoplot1`, and use nanoplot to analyse your passed fastq files:

```
cd TOAST2018
mkdir prac1_nanoplot1
NanoPlot --fastq /DAY1_DIR/reads/20180319_toast_GROUP/fastq/pass/*
-o prac1_nanoplot --title toast2018_p1 --loglength
```

The options `--fastq` and `-o` set the input and output for the command. We also want the length of the sequences shown log-transformed (`--loglength`) and also set a title for the plots (`--title`).

Data analysis and QC

The above command will create several plots and summary files including a *.html file that can be opened in your browser.

Note

If you are working on your own computer or with the virtual machine you can just double-click the html file.

If you are working remotely, please first download the complete folder *prac1_nanoplot1* to your computer and then double-click the file html file.

Inspect the different plots and statistics:

- How many reads do you have in total?
- What is the average, minimum and maximum read length, what is the N50?
- What do the mean quality and the quality distribution of the run look like?

Plot the yield

Another way of analysing your sequencing run is plotting the read yield of your flow cell over time. Low yield or a drop off of sequencing yield can indicate contaminations in the library or problems with the flow cell itself.

Use the tool *poretools* and the fast5 files to create a yield plot

```
poretools yield_plot \  
/DAY1_DIR/reads/20180319_toast_GROUP/fast5/pass/0/
```

You will have to close the plot to exit the command

NOTE: The backslash “\” in the commands can be ignored throughout the tutorial if you type the command in one continuous line!

Data analysis and QC

Export fast5 to fasta

Use the tools *poretools* to extract fasta files from all fast5 files in directory */fast5/pass/0* of yesterday's sequencing run. Use the option *--min-length* to only extract reads that are at least 2,000 nucleotides long. Because *poretools* prints the sequences to the terminal you need to *redirect* the output into a file using the ">" operator.

```
poretools fasta --min-length 2000 \  
/DAY1_DIR/reads/20180319_toast_GROUP/fast5/pass/0 > day1.fasta
```

You should now have a fasta file called *day1.fasta* in your current directory.

What did you sequence?

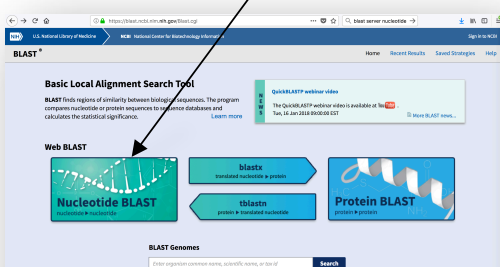
If you are sequencing your own projects you should have a pretty good understanding of where your DNA comes from. However, so far you don't know what your group actually sequenced yesterday.

Use the linux command *head* to list the first lines of the fasta file you just created and compare it to known genomes using the *Basic Local Alignment Tool* BLAST at the *Center for Biotechnology Information* (<https://blast.ncbi.nlm.nih.gov/Blast.cgi>)

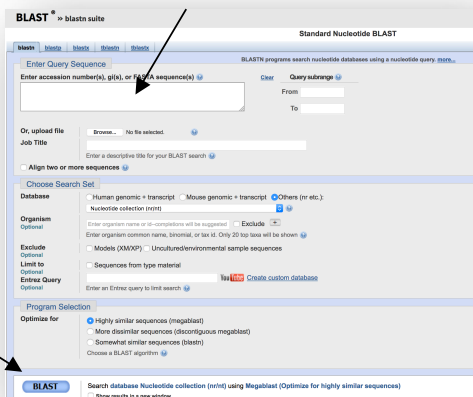
1. Mark and copy sequence (right-click -> copy)

[illegible]

2. Select *Nucleotide BLAST*



3. Paste your sequence



4. Start BLAST

Practical 2

Base calling with Albacore

Base calling with Albacore

First, change into the directory */TOAST2018/day2/prac2*. In the directory data you will find a *fast5* directory with fast5 files of a different MinION run. Use this directory as the input directory for the albacore run and call the output directory *albacore_out* (this directory will be created automatically by albacore).

```
read_fast5_basecaller.py -i ~/TOAST2018/day2/prac2/data/fast5 \  
-s ./albacore_out -f FLO-MIN106 -k SQK-LSK108 -o fastq -t 2
```

The command above will call albacore on the input fast5 directory option (-i), write the output to the directory given with option -s in fastq format (option -o). The options -f and -k defined the flow cell chemistry and the extraction kit, respectively. To see all possible attributes for those two parameters *call read_fast5_basecaller.py -i*.

While albacore is running you can see the estimated run-time in the lower right corner. For this particular data set it would take several hours to complete the base call. You can use the option -t *NUM* to specify the number of CPUs/cores to speed up the process. However, to save time we already ran alabacore on the data set. You can find the pre-computed data in the directory */TOAST2018/precomp/albacore_out*. Stop the albacore run by pressing *CTRL-C* and change into this directory.

Base calling with Albacore

Use `ls` to list the albacore output files and directories:

Configuration.cfg	The configuration file for the base call
Pipeline.log	A log file with base calling information per read
Sequencing_summary.txt	A summary of the sequencing run
workspace	A directory containing the fail and pass fastq directories for the base called reads

You will always find these files in your albacore output directory.

Compare *fail* and *pass* reads using nanoplot

Create a directories *alba_pass* in the course directory `~/TOAST2018/day2/prac2/` and use nanoplot on the fail and pass directories.

```
mkdir ~/TOAST2018/day2/prac2/alba_pass
```

```
NanoPlot --fastq ~/TOAST2018/precomp/albacore_out/workspace/pass/*  
-o /TOAST2018/day2/prac2/alba_pass --title alba_pass --loglength
```

Base calling with albacore

Do the same with the fastq files in the *fail* directory:

```
mkdir ~/TOAST2018/day2/prac2/alba_fail
```

```
NanoPlot --fastq ~/TOAST2018/precomp/alabcore_out/workspace/fail/*  
-o /TOAST/day2/prac2/alba_fail --title alba_fail --loglength
```

Open the different plots and compare them.

What are the main differences, e.g., wrt. sequence length and quality?

Additional exercise: Plotting with poretools

The data for this practical originated from a 24h sequencing run. If you are interested use the fast5 directory and poretools to create a yield_plot for this run. What do you see?

```
poretools yield_plot ~/TOAST2018/day2/prac2/data/fast5
```

What might be the explanations for this plot?

Practical 3

Read trimming and adapter removal

Read trimming & adapter removal

Adapter removal using porechop

Similar to other 2nd and 3rd generation sequencing platforms MinION library prep kits use ligate adapters to the ends of the DNA. Additionally, some library preparations kits, such as the older 2D kits, use additional adapters to link the two strands of the DNA. Also, in some instances chimeric reads may occur that include adapters in the middle of the sequence.

Porechop is a freely available open-source tool that among others can be used to find and remove adapters from Oxford nanopore reads and remove chimeric reads.

Create a directory *porechop* in the directory *TOAST2018/day2/prac3* and use *porechop* to remove adapters from the albacore reads that passed the QC.

```
mkdir ~/TOAST2018/day2/prac3/porechop

porechop -i ~/TOAST2018/precomp/albacore_out/workspace/pass/
-o ~/TOAST2018/day2/prac3/porechop/porechopped.fastq \
--discard_middle
```

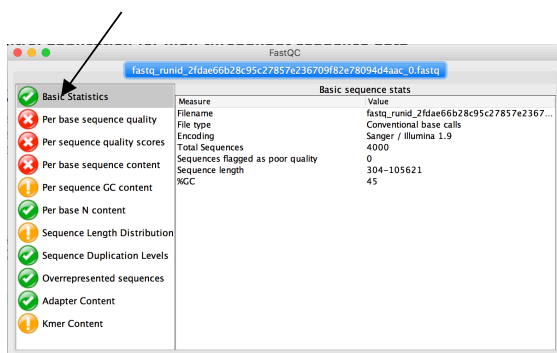
The above command will use the default values of *porechop* to search for adapters in all fastq files of the input directory, trim the reads and write them to file *trimmed.fastq* in the created *porechop* directory. The “*--discard_middle*” option will remove reads with internal adapters. As mentioned by R. Wick on the *porechop* github page (<https://github.com/rrwick/Porechop>), this is essential if you want to use nanopolish later in the analysis.

Read trimming & adapter removal

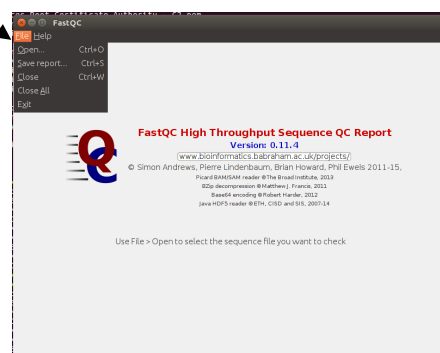
Plot reads using FastQC

In addition to poretools and nanoplot another tool can be used to analyse your sequences visually is FastQC. FastQC is a common tool for quality control of 2nd generation Illumina sequencing reads which can also be use to analyse 3rd generation long-read sequences, e.g. to identify low-quality regions in your read data.

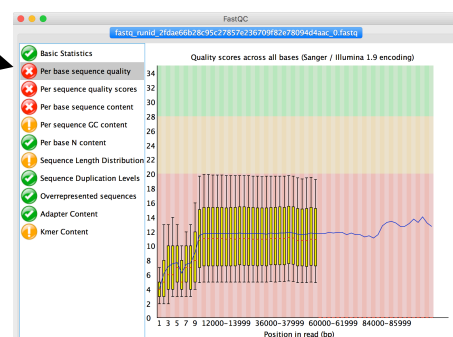
Type *fastqc* on the command-line to open its *graphical user interface* and load the fastq file of the reads that passed albacore



The second tab “Per Base sequence quality” shows the mean and standard deviation of the sequencing quality for each position in all reads of your data set.

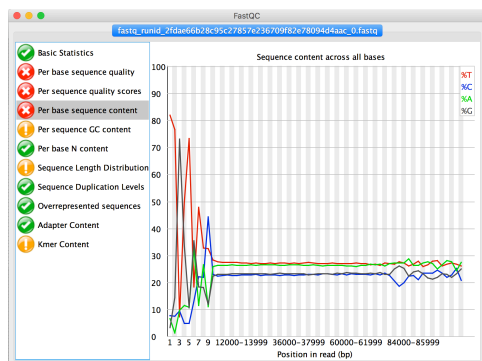
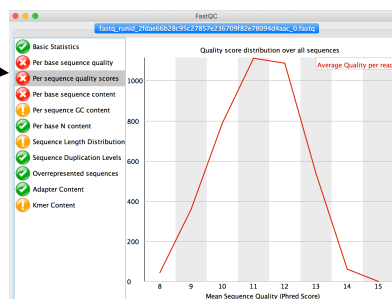


After loading the first tab will show you some basic statistics about your fastq file and sequences



Read trimming & adapter removal

The tab “Per Sequence Quality scores” shows the average quality score distribution of your nanopore reads



The tab “Per base sequence contents” show the average ratio of As, Ts, Cs and Gs in your data set. “Clean” data without sequencing should show almost parallel lines for all four nucleotides.

Using the FastQC output we can make a decision whether we want to trim our reads further, e.g., remove low-quality areas of the reads.

Which areas would you trim off?

Note: Exit FastQC to return to command-line prompt.

Read trimming & adapter removal

Trim reads using NanoFilt

NanoFilt is a python script that can be used to filter reads based on length and quality as well as trimming parts of the sequence.

Create a directory for your nanofilt output (e.g., `~/TOAST2018/day2/prac3/nanofilt`) and then use nanopolish to remove:

- all sequences shorter than 500 nucleotides (option `-l`)
- trim the first 10 nucleotides off all reads (option `--headcrop`)

```
mkdir ~/TOAST2018/day2/prac3/nanofilt
```

```
NanoFilt -l 500 --headcrop 10 \  
< ~/TOAST2018/day2/prac3/porechop/porechopped.fastq \  
> ~/TOAST2018/day2/prac3/nanofilt/nanofilt_trimmed.fastq
```

NanoFilt does not provide options for input or output files. Therefore we will use the two *redirect* operators “>” and “<” to

- redirect the file porechopped.fastq into NanoFilt (operator <)
- then redirect the output of NanoFilt into the file nanofilt_trimmed.fastq (>).

Use FastQC to check the result file and compare it to the original fastq.

Practical 4

Assembly using Minimap & miniasm

Assembly using minimap & miniasm

Assembling long-reads, PacBio and Oxford Nanopore, into contiguous sequences (contigs) has been challenging for common 2nd generation assemblers due to the high error rates of 3rd generation sequencing technologies. Recently, an increasing number of assemblers and assembly pipelines is available that take into account the specific characteristics of long-reads.

One way of assembling long-reads is the combination of *minimap* and *miniasm*.

Minimap is read mapper that can identify overlaps in reads.

Miniasm is a very fast overlap assembler that outputs unitigs, i.e., high confidence overlap sequences. In contrast to other long-read assemblers miniasm does not include a consensus step, i.e., miniasm unitigs have a similar error rate as the input reads.

First use minimap to map the filtered nanopore reads onto themselves

```
minimap2 -x ava-ont \  
~/TOAST2018/day2/prac3/nanofilt/nanofilt_trimmed.fastq \  
~/TOAST2018/day2/prac3/nanofilt/nanofilt_trimmed.fastq \  
| gzip -1 > ~/TOAST2018/day2/prac4/minimap.paf.gz
```

Minimap comes with several pre-configured parameter settings depending on the read data. The option *-x ava-ont* sets the default parameters for Oxford Nanopore reads.

The additional *pipe* (|) into *gzip* is useful for large output files to compress them before writing the compressed output file.

Assembly using minimap & miniasm

Use the minimap output and the trimmed reads to assemble unitigs with miniasm.

```
miniasm -f ~/TOAST2018/day2/prac3/nanofilt/nanofilt_trimmed.fastq \  
~/TOAST2018/day2/prac4/minimap.paf.gz \  
> ~/TOAST2018/day2/prac4/miniasm.gfa
```

Miniamap and miniasm do not provide an option for output files but instead write the output directly to the terminal. Thus, the output has to be redirected using the “>” operator into the result [gfa file](#).

To convert the miniasm.gfa file into a fasta file of unitigs use the following [awk](#) command

```
awk '/^S/{print ">"$2"\n"$3}' ~/TOAST2018/day2/prac4/miniasm.gfa \  
> miniasm.fasta
```

Now use the tools *assembly-stats* to get some simple statistics about the assembly

```
assembly-stats ~/TOAST2018/day2/prac4/miniasm.fasta
```

Assembly-stats will tell you how many unitigs miniasm assembled, the total length of all

```
toast2018@toast2018-VirtualBox:~/TOAST2018/day2/prac4$ assembly-stats miniasm.fasta  
stats for miniasm.fasta  
sum = 859790, n = 7, ave = 122827.14, largest = 675939  
N50 = 675939, n = 1  
N60 = 675939, n = 1  
N70 = 675939, n = 1  
N80 = 54955, n = 2  
N90 = 34078, n = 4  
N100 = 12928, n = 7  
N_count = 0  
Gaps = 0
```

Assembly using minimap & miniasm

How to assess the assembly quality

The best way to analyse the quality of an alignment is to compare it to a published sequence. In directory `/TOAST2018/precomp/data` you will find the published sequence of chromosome_17 of the marine diatom *Thalassiosira pseudonana*, the source of most of the reads used in this tutorial.

First get some assembly statistics. Use *assembly-stats* to print report statistics about this chromosome, e.g., number of nucleotides and number of Ns (gaps or ambiguous sequences).

Compare the two statistics. Do some of the miniasm unitigs match the reference sequence in length?

To make downstream analysis easier first remove sequences that are much shorter than the reference sequence. Use the command *filter_fasta_by_seq_length.pl* from the *ampli-tools* package to remove all sequences shorter than 100,000 nucleotides from the assembly fasta

```
filter_fasta_by_seq_length.pl -i ~/TOAST2018/day2/prac4/miniasm.fasta  
-o ~/TOAST2018/day2/prac4/miniasm_filtered.fasta -a 100000
```

and count the resulting sequences with the linux command *grep*

```
grep -c ">" ~/TOAST2018/day2/prac4/miniasm_filtered.fasta
```

The above command will count the occurrence of the character ">" in the input file, which in a fasta file is equivalent to number of sequences.

Assembly using minimap & miniasm

How many sequences “survived” the filtering?

To compare the left unitig to the reference genome use the tool *dnadiff* from the mummer package. Mummer provides fast alignments of large highly similar sequences to each other.

```
dnadiff ~/TOAST2018/precomp/data/chr17.fasta -p dnadiff \  
~/TOAST2018/day2/prac4/miniasm_filtered.fasta
```

This command will create multiple output files of prefix *dnadiff* with alignment statistics and detailed coordinate information of the overlapping regions of both sequences.

Open the file *dnadiff.report* (e.g. double click) to see a general report of the alignment.

How much of both sequences was aligned?

What is the average percent identity of the alignment?

We can also visualise the alignment with a dot-plot using the command *mummerplot*.

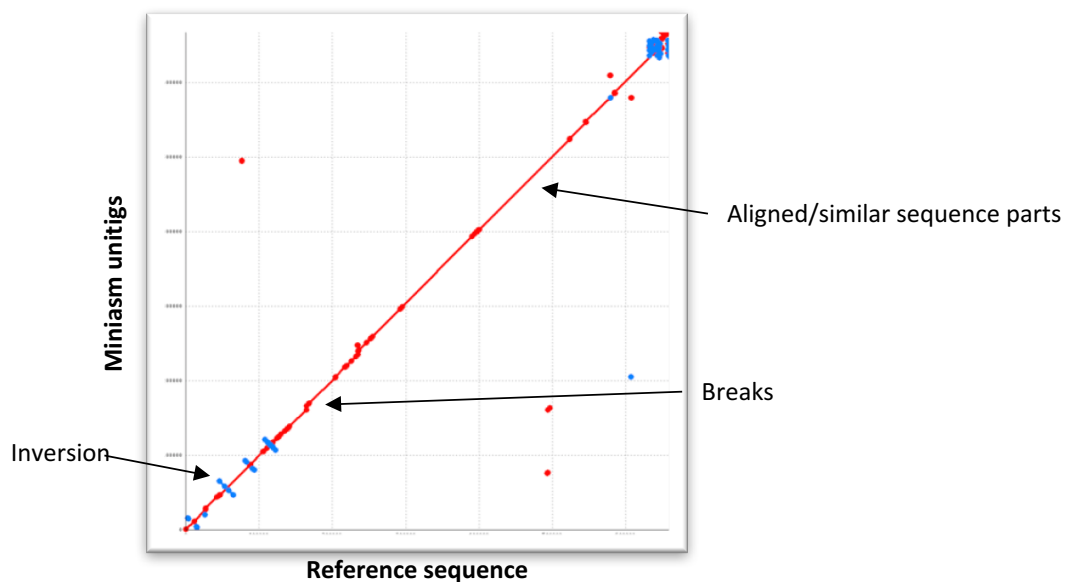
```
mummerplot --png -p miniasm ~/TOAST2018/day2/prac4/dnadiff.delta \  
-R ~/TOAST2018/precomp/data/chr17.fasta \  
-Q ~/TOAST2018/day2/prac4/miniasm_filtered.fasta
```

The above command will create an image file in png format in directory ~/TOAST2018/day2/prac4. Open it by double clicking the file and analyse the plot.

Assembly using minimap & miniasm

Dot plots

A dot-plot is a visual representation of the similarity of two sequences. In this case the resulting dot-plot shows the miniasm unitig on the y-axis and the reference sequence of chr17 on the x-axis. Sequence parts that overlap are shown as a diagonal line from the lower left to upper right corner. Similarly, orthogonal lines (upper left to lower right) indicate inversions in one of the sequences. Breaks and gaps in the line indicate deletions or insertions in either of the sequences.



What do you think about this initial assembly?

Practical 5

Create consensus sequences using racon

Consensus sequences using racon

Due to consistently high error rates of Oxford Nanopore sequencing data downstream analysis has to include some sort of error correction especially if no high quality short read data is available.

Consensus assemblies try to reduce error rates by choosing the *most likely* sequence of a given assembly and a set of raw reads. Although this does not incorporate the raw signal information of the flow cell to correct individual reads it can significantly improve the quality of an assembly.

The software package racon has been developed to complement the minimap and miniasm tools. It provides a fasta consensus algorithm that uses either 2nd generation short reads or raw noisy long-reads to correct draft assemblies.

To improve a draft assembly with racon map the reads that should be used for error correction against the assembly.

Use minimap to map the trimmed reads from directory `~/TOAST2018/prac3/nanofilt` against the miniasm assembly and subsequently use the filtered reads and the mapping to build the consensus assembly.

```
minimap2 ~/TOAST2018/day2/prac4/miniasm_filtered.fasta \  
~/TOAST2018/day2/prac3/nanofilt/nanofilt_trimmed.fastq \  
> ~/TOAST2018/day2/prac5/minimap.racon.paf  
  
racon ~/TOAST2018/day2/prac3/nanofilt/nanofilt_trimmed.fastq \  
~/TOAST2018/day2/prac5/minimap.racon.paf \  
~/TOAST2018/day2/prac4/miniasm_filtered.fasta \  
> ~/TOAST2018/day2/prac5/consensus_assembly.fasta
```


Consensus sequences using racon

Analyse the assembly quality by comparing the consensus assembly to the published sequence using *dnadiff* and *mummerplot*.

```
dnadiff ~/TOAST2018/precomp/data/chr17.fasta -p dd_con \  
~/TOAST2018/day2/prac5/consensus_assembly.fasta  
  
mummerplot --png -p miniasm ~/TOAST2018/day2/prac5/dd_con.delta \  
-R ~/TOAST/precomp/data/chr17.fasta \  
-Q ~/TOAST/day2/prac5/consensus_assembly.fasta
```

Did the assembly improve?

Which parts did not improve?

Run racon a second time but this time on the consensus assembly

```
minimap2 ~/TOAST2018/day2/prac5/consensus_assembly.fasta \  
~/TOAST2018/day2/prac3/nanofilt/nanofilt_trimmed.fastq \  
> ~/TOAST2018/day2/prac5/minimap_2.racon.paf  
  
racon ~/TOAST2018/day2/prac3/nanofilt/nanofilt_trimmed.fastq \  
~/TOAST2018/day2/prac5/minimap_2.racon.paf \  
~/TOAST2018/day2/prac5/consensus_assembly.fasta \  
> ~/TOAST2018/day2/prac5/consensus_assembly2.fas
```

Consensus sequences using racon

Analyse the second consensus

```
dnadiff ~/TOAST2018/precomp/data/chr17.fasta -p dd_con2 \  
~/TOAST2018/day2/prac5/consensus_assembly2.fasta  
  
mummerplot --png -p miniasm2 ~/TOAST2018/day2/prac5/dd_con2.delta \  
-R ~/TOAST/precomp/data/chr17.fasta \  
-Q ~/TOAST/day2/prac5/consensus_assembly2.fasta
```

What are the improvements?

Did some statistics decrease? If so which and why?

Practical 6

Error correction using nanopolish

Error correction using nanopolish

Nanopolish is another free open-source tool for nanopore data analysis. In contrast to tools such as racon it uses not only the sequence information from fasta/fastq files but also utilises the raw signal of each read that is stored in the fast5 files: it compares the raw signal of each nucleotide in a dataset and tries to identify incorrect base calls as well as DNA modifications such as methylation.

The recommended workflow to polish and existing draft assembly includes the following steps:

1. Nanopolish index	Use nanopolish to link the raw fast5 signal with base called fastq sequences
2. BWA index	Prepare your draft genome for use with read aligner BWA
3. Read mapping	Use BWA together with samtools to map the base called reads back to your draft genome
4. Nanopolish consensus	Build an error corrected consensus

For large assemblies (>50K) nanopolish provides two scripts *nanopolish_makerange.py* and *nanopolish_merge.py* to split the sequence into multiple segments before the error correction steps and merging them again afterwards (see http://nanopolish.readthedocs.io/en/latest/quickstart_consensus.html for more details).

Error correction using nanopolish

1. Nanopolish index

To index the nanopore reads, i.e., to connect the fastq sequences with the corresponding raw signal first concatenate all fastq files in the `albacore_pass` directory into one fastq file using the linux command `cat`:

```
cat ~/TOAST2018/precomp/albacore_out/workspace/pass/* \  
> ~/TOAST2018/day2/prac6/all_pass.fastq
```

This will write one fastq file into the `prac6` folder containing all sequences that passed the initial albacore quality control.

Use nanopolish to index your fast5 and fastq reads

```
nanopolish index -d ~/TOAST2018/day2/prac2/data/fast5 \  
~/TOAST2018/day2/prac6/all_pass.fastq
```

This command will create several files index files in the directory for practical 6.

2. BWA index

BWA is a widely used free open-source read mapper developed to quickly map sequences back to a reference genome. To speed up mapping of large read files *bwa* indexes the reference genome

Index the draft assembly using

```
bwa index ~/TOAST2018/day2/prac4/miniasm_filtered.fasta
```

Error correction using nanopolish

3. Map reads to draft assembly

As input for the mapped reads nanopolish expects a sorted and indexed *bam* file.

Use *bwa* to create an output file in bam format and the free open-source package *samtools* to subsequently sort and index it.

```
bwa mem -x ont2d ~/TOAST2018/day2/prac4/miniasm_filtered.fasta \  
~/TOAST2018/day2/prac6/all_pass.fastq | samtools sort \  
-o ~/TOAST2018/day2/prac6/reads.sorted.bam  
  
samtools index reads.sorted.bam
```

4. Run nanopolish

Use the mapping information together with the linked fastq and fast5 files to create a polished assembly

```
nanopolish variants --consensus ~/TOAST2018/day2/prac6/nano.fasta \  
-r ~/TOAST2018/day2/prac6/all_pass.fastq \  
-b ~/TOAST2018/day2/prac6/reads.sorted.bam  
-g ~/TOAST2018/day2/prac4/miniasm_filtered.fasta  
-q dcm,dam  
--min-candidate-frequency 0.1
```

Note

This nanopolish call will run for several hours. You can find the final file *nano.fasta* in directory *~/TOAST2018/precomp/data*.

Error correction using nanopolish

The above call will call nanopolish in *methylation aware* mode. DNA modifications can lead to errors in the base calling due to the fact that modified nucleotides show a different raw signal than their unmodified counterparts. By using option *-q dcm,dam* nanopolish is trying to determine the differences in the signal of 5-methyle-cytosine and unmodified cytosines and use this information to correct potential base calling errors.

As before, use dnadiff to assess the nanopolish results and compare them with racon:

```
dnadiff ~/TOAST2018/precomp/data/chr17.fasta -p dd_con3 \  
~/TOAST2018/day2/prac5/consensus_assembly2.fasta
```

What are the differences?

If you run racon on the polished data does it improve?

Note

Despite the systematic errors of nanopore base callers sequencing depth as well as organism can affect nanopolish results. The current eukaryotic test data set has a mean sequencing depth of <10. Nanopolish results of >99% accuracy have been reported for bacterial data set with >15 x coverage.

Practical 7

Canu assembly quality

Canu assembly quality

In contrast to miniasm the nanopore assembler *Canu* combines unitig assembler and contig assembly and also provides error correction. Recent assembler comparisons on bacterial data indicated higher accuracy and lower error rates for canu assemblies. However, canu's accuracy comes the cost of significant increases in run time and computational resources.

In directory `~/TOAST2018/precomp/data` you can find the canu assembled contigs of the tutorial data in file *canu.contigs.fasta*.

Additionally, the same directory includes the the largest if the canu assembled contigs polished with nanopolish in file *canu.contigs.polished.fasta*.

The commands used to assemble and polish the the sequences are:

```
Canu -p OUTPUT_PREFIX -d output_directory genomeSize=0.6m \  
--nanopore-raw READS.FASTQ useGrid=false maxThreads=NO_CPUS \  
maxMemory=RAM  
  
bwa index CANU_ASSEMBLY  
  
bwa mem -x ont2d -t THREADS READS.FASTQ |samtools sort -o \  
OUTPUT_BAM  
samtools index OUTPUT_BAM  
  
nanopolish variants --consensus OUTPUT_FASTA -r READS.FASTQ \  
-g CANU_ASSEMBLY -b READS.BAM -min-candidate-frequency 0.1\  
-q dcm,dam
```

Hope you enjoyed these practical.

For questions and improvement suggestions
please contact me at tim.kahlke@uts.edu.au