

HPC-getorganelle

```
1  #!/bin/bash
2
3  #PBS -S /bin/bash
4  #PBS -l walltime=06:00:00
5  #PBS -N JOBNAME
6  #PBS -o stdout.$PBS_JOBID
7  #PBS -e stderr.$PBS_JOBID
8  #PBS -l nodes=1:ppn=all:gpus=2
9  #PBS -m abe
10 #PBS -M pieter.asselman@ugent.be
11
12 #####
13 ##INFO##
14 #####
15
16 #DATArun on 2GPU/ 11Gbases - ~3h
17
18 ##BASECALL##
19 guppy_basecaller -i /scratch/gent/vo/000/gvo00058/vs
20
21 ##DEMULPLEX##
22 -guppy_barcode --i /scratch/gent/vo/000/gvo00058/vs
23
24 ##BASECALL and DEMULPLEX##
25 -guppy_basecaller -i /scratch/gent/vo/000/gvo00058/vs
26
```

```
vsc43352@gligar04: ~
Using username "vsc43352".
Authenticating with public key "rsa-key-HPC"
Passphrase for key "rsa-key-HPC":
Last login: Tue Aug  3 15:14:06 2021 from gligarha01.gastly.os

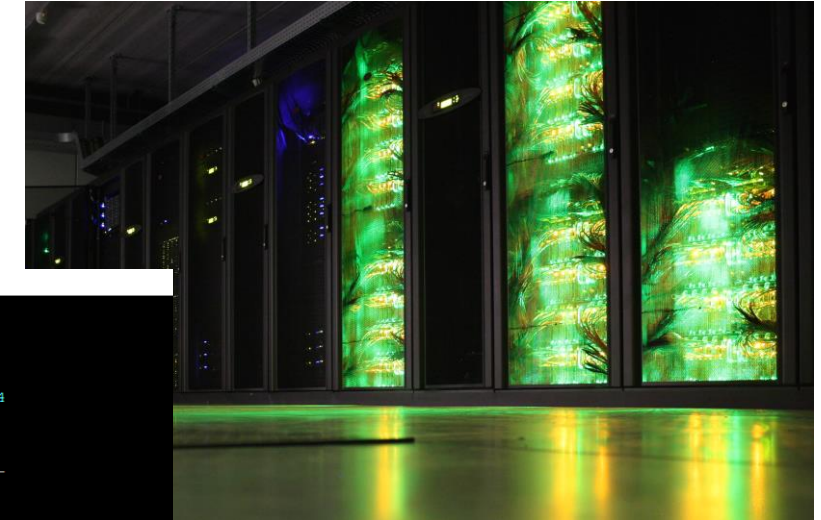
STEVIN HPC-UGent infrastructure status on Wed, 04 Aug 2021 07:25:04

cluster - full - free - part - total - running - queued
         nodes nodes free  nodes  jobs    jobs
-----
swalot   90    0   20   125   N/A    N/A
skitty   19    0   52    72   N/A    N/A
victini   8    1   86    96   N/A    N/A
joltik    9    0    0   10   N/A    N/A
kirlia   15    0    1   16   N/A    N/A
doduo   120    0    5   128   N/A    N/A

For a full view of the current loads and queues see:
https://hpc.ugent.be/clusterstate/
Updates on current system status and planned maintenance can be found on
https://www.ugent.be/hpc/en/infrastructure/status

We switched to new job command wrappers for all HPC-UGent Tier-2 clusters on Wed
June 9th 2021 at 17:22.
This switch should be transparant: you don't need to change your workflow or job
scripts.

If you notice any problems, or if any of these changes affect your work, please
contact hpc@ugent.be .
vsc43352@gligar04:~$
```



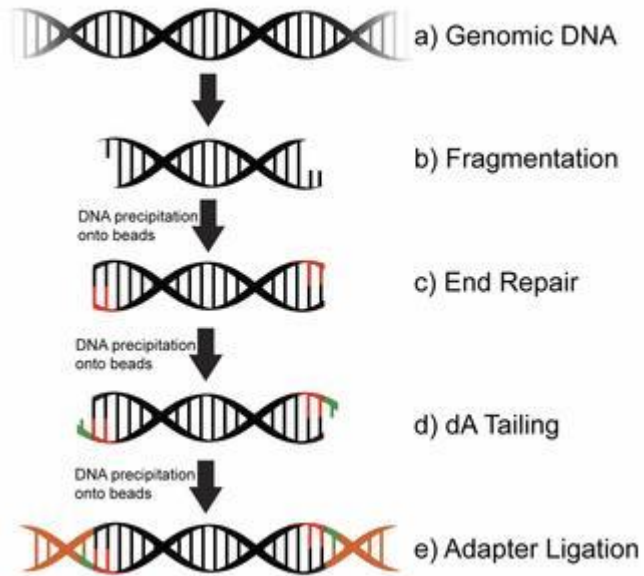
<https://ugent.be/hpc>
hpc@ugent.be

Pipeline: Cp assemblies

- I. Lib prep
- II. Sequencing - Demultiplexing
- III. QC
 - a) Sequencing Report
 - ✓ QC10: 1-10 error (90% accurate)
 - ✓ QC20: 1-100 error (99% accurate)
 - ✓ QC30: 1-1000 error (99,9% accurate)
 - b) Fastqc: <https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>
 - ✓ Quick overview of all reads
 - ✓ Summary graphs
 - ✓ HTML-based report
 - ✓ Multiqc! Overview of multiple HTML-reports from different analyses
 - c) Trim
 - ✓ Trimmomatic: <http://www.usadellab.org/cms/?page=trimmomatic>
 - ✓ Cutadapt: <https://cutadapt.readthedocs.io/en/stable/>
- IV. Cp – nr Assembly: Getorganelle
 - ✓ Toolkit for assembly of organelle genome: <https://github.com/Kinggerm/GetOrganelle>

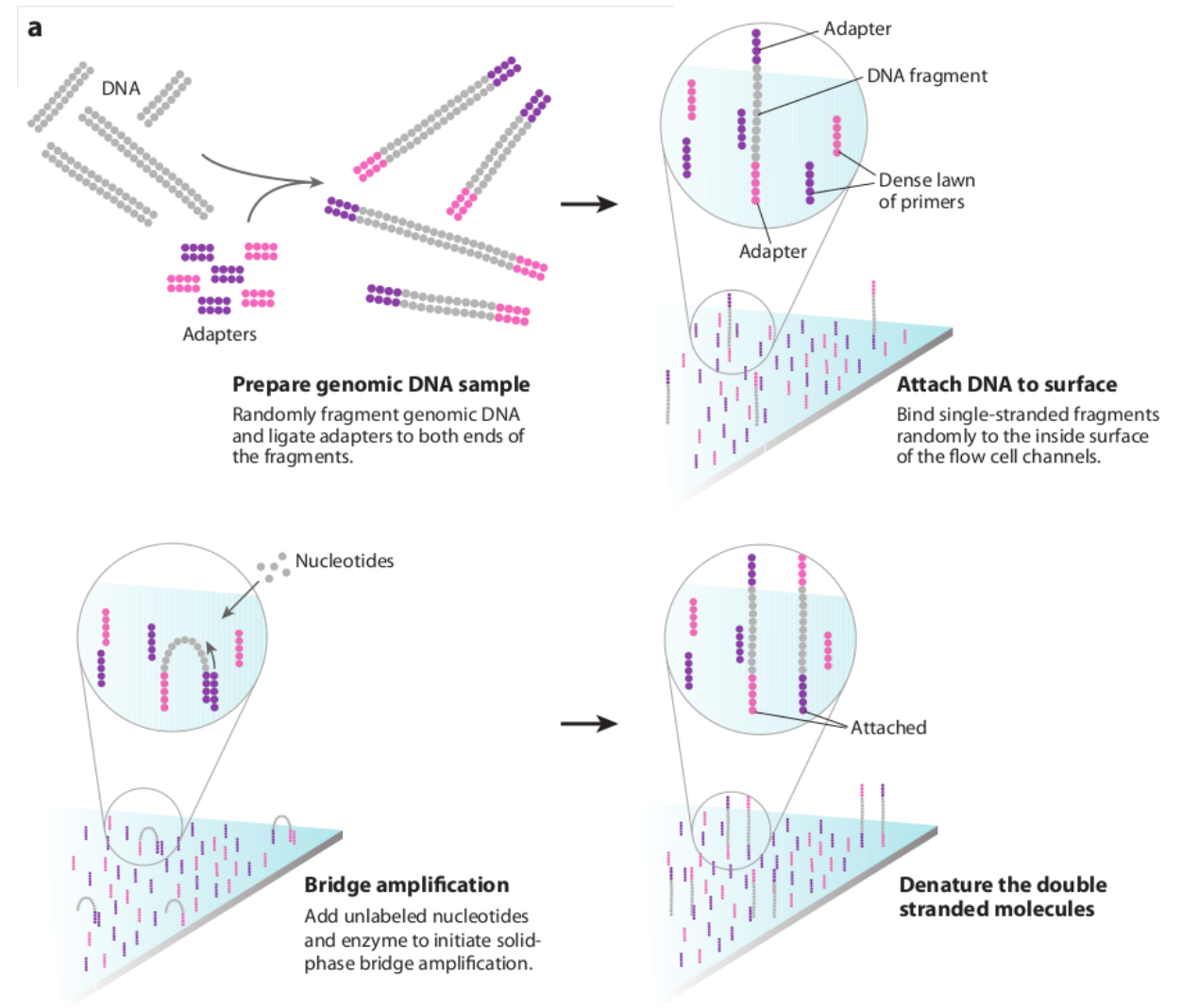
Pipeline: Library prep & Illumina Sequencing

A) Workflow of the automated library preparation



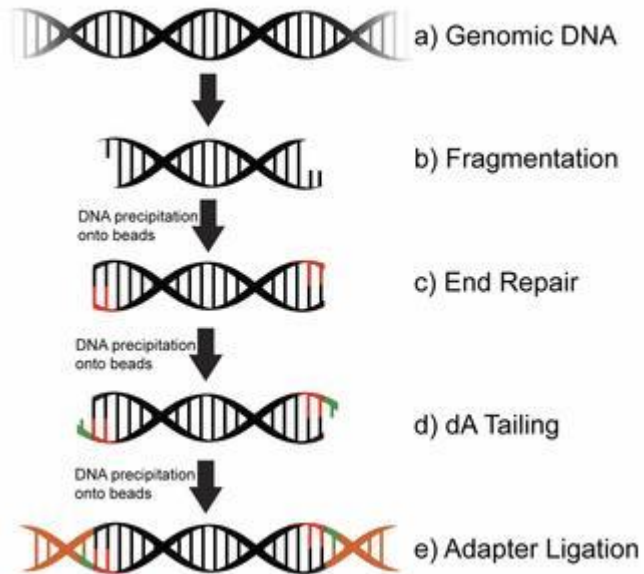
*Borgstrom E et.al.2011 Large Scale Library Generation for High Throughput Sequencing. PLoS ONE 6:e19119.

doi:10.1371/journal.pone.0019119 *



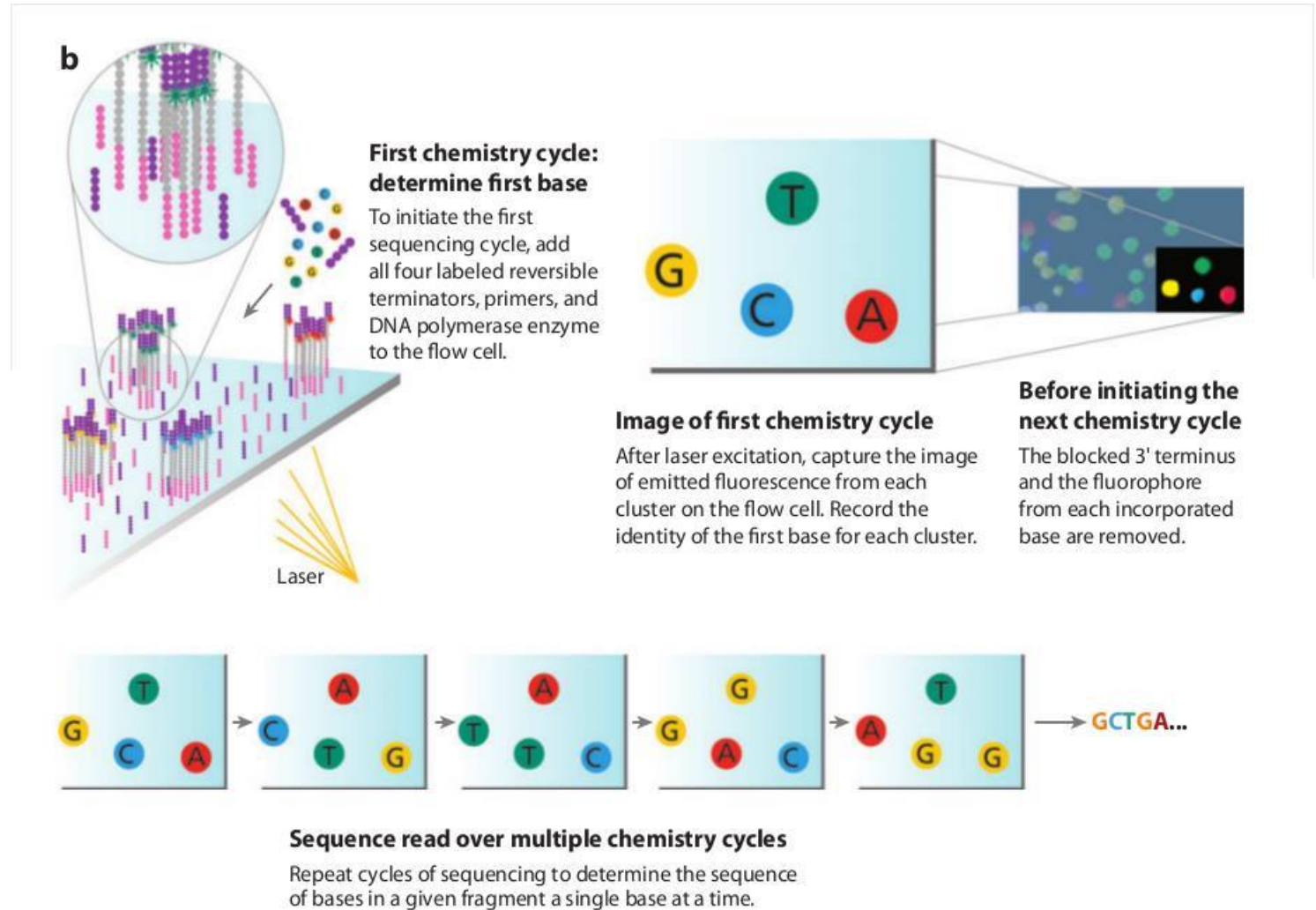
Pipeline: Library prep & Illumina Sequencing

A) Workflow of the automated library preparation

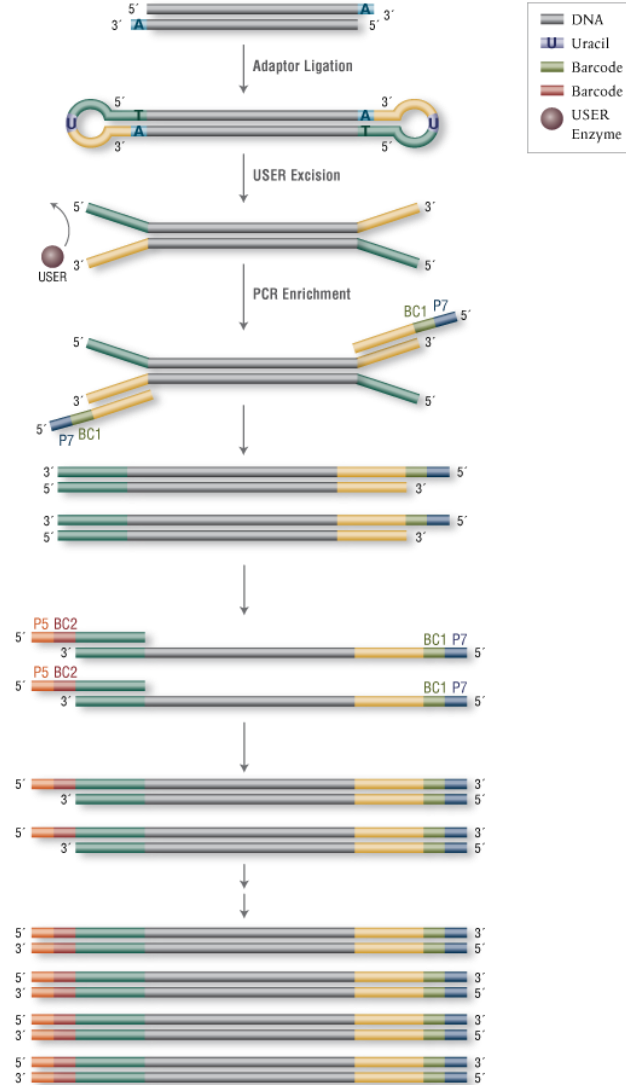
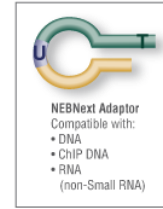


*Borgstrom E et.al.2011 Large Scale Library Generation for High Throughput Sequencing. PLoS ONE 6:e19119.

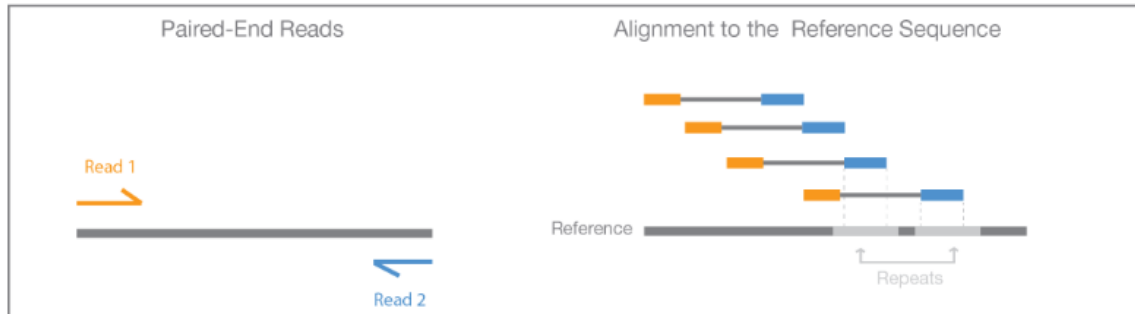
doi:10.1371/journal.pone.0019119 *



Pipeline: Library prep



Dual index barcoding P5/BC2 – BC1/P7

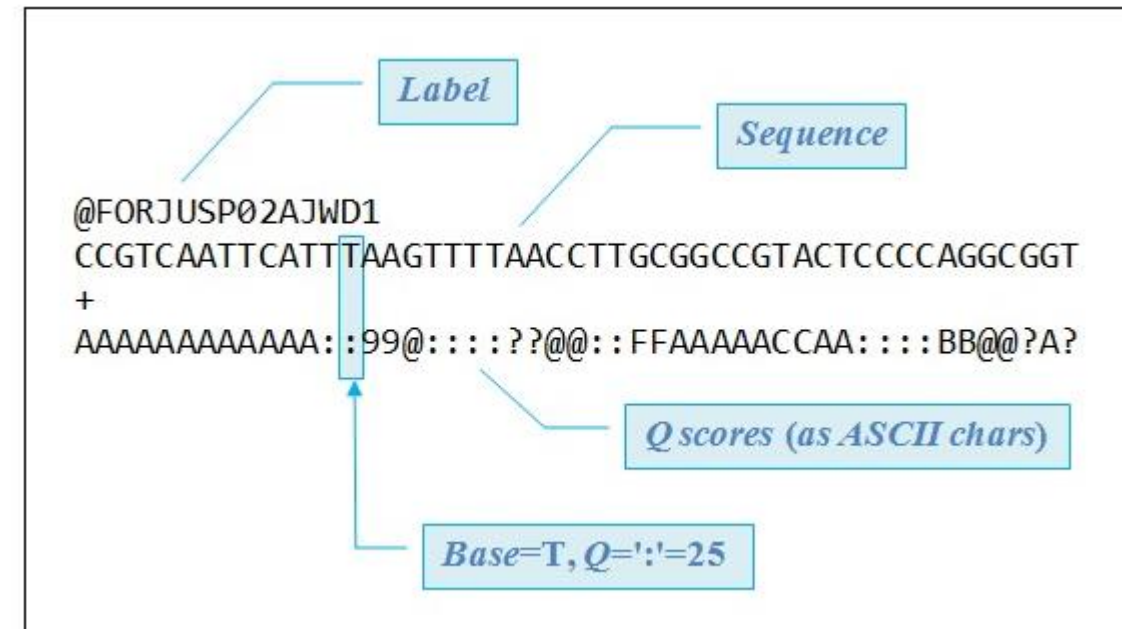


Pipeline: Data format

fastq

```
@de4d97b8-c1c5-4d48-83dd-cf35c8b1b262 runid=86431727e3e8469fcfb154f8fc6e278e2ca5217d read=4 ch=443
start_time=2017-10-30T10:11:28Z
TGATGCTTCGTTTCAGTTGCGTATTGCTCTGCCAGATACGCCCGCGCAGATGATTACGGCTCACCTGCTGGCTGCCAGTCTCAGGTCACCTGCGCCGAC
CGGTGCCGTGTTCCGGGTAGCCTGTTTCATCAGTAGTTTCAGTGTATGACGTGGCGGTTTATGCGAACTGGTCTGTGCTGTGTCAGGTCATCCTGCGCATC
CGGATGTAAACGGGGCGCGCGCGGGTATACCATCACCCCGGGCGCGGATGCGGGTGGTTTCTATAATGGGTGATGATGCTTTTCCGTCACCCAGCC
ACCGTTACGTTGTATCTAGACAGCAGCGGACTACGACGGTTCTATACCCCTTTGGGGTGGTTTCTGTACGAACGGAAGTAAAGCGCAGACGGTGCATG
TTTGACAGATTGTAATGGGTGCGGGTGTATGCGGTGTGCTATTTCGCTTGTACTTGGCGCGCACCGTCTCAGGCCAGCGGAGTTCAAATCCCTCCAGCGGAT
CTATCCTCTCACCAGCCCGGAACACCTGTGGTACGGATATATTAGTATTCCCTCAGTATCCAGCACGAGCTACTGTGTTTCAGCGCGCTTTTCCGTAT
CCACCGGACCTTCAATCCAGCCTTCGCTGGTGGCATCGATCCACACTCCAGCAGCGGTGGACTTGGAGTTGTCAGCCAGCTGCGCAGGGGTATGCCCTTAC
TGCTTCTTTACCCATTTCCTCGCTCCACCATTAAGCGAAACCGCAGCGGAGTTCCACATATAACATTTTGCATCTTGACCGAAATCACCTGGCAACACCC
GTCCCTTCGCTGCGCGTGCAGTCTCGCAGAGCCGCTATTTCGCGCATTCCTCATGACAGCAGGCCAGAACATTTGCTGGGACGCAATTATCAGTGGGG
AGAAATAGGTGTTCTGCTTACCGTTATCCGTTGCTGTATACGGGGTTCGGCACCAGGTGCCAGCATCCGCCACTACCGGGCGCCATACGCGACCCGAAAAAC
AGGATGCCAGTCATACACCGGGCCCAATGAAACTGCCCTCATGCTGCAAGGGTGGCTCCGGCGAGTAAAGAAATGATCCGGCGTGGCGGACGCCAGGAC
AATCTGGAATCGCCACCGGCTGACTGGCGACTCTGSGGAACAATGAATTACAGCGCCATCAGGCCAGAGTCTCATGCATAACACCGCGTTAACCGGACATT
TGGCGTCCGGCAATCCGTACCTGATACCGCAATCGCTCAGTTTCACAGCCGACGCCAGGCTATGTGGCAGTCCCGGATGGCTTACGCCCA
+
#####-7A83<D+5.).+731+)).( )-5./6<8+2*3,::744*322(+1*210265.7,,)5,'+:8789*/-./+17026289<:
4'2?+0<>:>4/,)6*2SS$'.*56/.3:<%(4',1:83(C8:3*+.17-,*-05162($(&'%+))-.3852-),*,-19-32=8?<=AC?45264?<
4=?>:*2(3,2-6AAAA</-50*4269%/( )-.<8,14/-736:*??,+,+.,:2>>230((, /30+6.3/-,-1387<8ED>0:169-298,+3
(, ('-($&*,0---+*'+&--.'(((*+/*4188/3:898=.80,)8:5>C79D2+-. (0/?4A(*'SS%&,+%57)6*9,.%) +/2E8A>@/CB96887
8C=7<120-*$(',-4:.(, &*&'(.02.83+3*+&/ (5C0+&'&#%&'&'%*,.)' **296--,,+7:8?*+&3-*+*2*&(*=-+-.%43:,0
10798*((,0((('S+)):/:352*/./0.+,+&.,+0*+.6C6(/4<D5CD.*53'66,11669%6+02334:8,-%5>)))*(+/,*+*(6(118:
9839@4C25++=3,1?%(1-1.-./03.,4),0*-;>((1.1+2'*&.,6.,(7-A/677+)%)/5*/06.-'%*,+*%45>6=45';D7)2+-
)) /21))8(-)($+88:32--22:C-39(3-+).(+),;E256754-617C0-*.) /+20+---A2,/2,('&'15/61052<9<90-)*(&(+*55
9=27:~09:;=C:66:/3/4*S/<<78S)&'3'/(, (**%&,-( ),+(-,--)+*%&'&+65/409(%&,50/17105A:;2*( /+:433:024'87/+9;
--(1A-,3-<5?61*)' '-0+6501=29:D223.*-,-':(346-7.,&&#%&/@<8*,,'54(, /+)&'..+**56*.,(S'5-(-:/6/3439
687>86.*6)1*&,-='&,' /+/-.*0:857/6*(,4BFD><:65,111/-8314<938,-/6-)&*8B(')116:CCB8=?,>0@'62+-94*$4620
<40:;+9S%327:'5$,&)%#(S)&-*.,-884,889996
+ )/-+)'',5<A==8'<+0>*( (,,:*3%6%&,6,16:<?3
@4ea9b632-cbca-455d-9061-88117b879114 run
start_time=2017-10-30T10:11:30Z
TGTTACTTTAGTTTCAGACCATTCACCATCAGATTATGTTT
ACACGTTTTGCCCGTACTCGGTACGAACAATCCTGATTAT
CCCTTATTGAGGGTAAAGACGCCAGAAATTGAAACGCGGC
```

\$
ASCII : 36
P: 0,5 (50% chance the basecall is wrong)
Q: 3



ASCII_BASE=33 Illumina, Ion Torrent, PacBio and Sanger

Q	P_error	ASCII	Q	P_error	ASCII	Q	P_error	ASCII	Q	P_error	ASCII	Q	P_error	ASCII
0	1.00000	33 !	11	0.07943	44 ,	22	0.00631	55 7	33	0.00050	66 B			
1	0.79433	34 "	12	0.06310	45 -	23	0.00501	56 8	34	0.00040	67 C			
2	0.63096	35 #	13	0.05012	46 .	24	0.00398	57 9	35	0.00032	68 D			
3	0.50119	36 \$	14	0.03981	47 /	25	0.00316	58 :	36	0.00025	69 E			
4	0.39811	37 %	15	0.03162	48 0	26	0.00251	59 ;	37	0.00020	70 F			
5	0.31623	38 &	16	0.02512	49 1	27	0.00200	60 <	38	0.00016	71 G			
6	0.25119	39 '	17	0.01995	50 2	28	0.00158	61 =	39	0.00013	72 H			
7	0.19953	40 (18	0.01585	51 3	29	0.00126	62 >	40	0.00010	73 I			
8	0.15849	41)	19	0.01259	52 4	30	0.00100	63 ?	41	0.00008	74 J			
9	0.12589	42 *	20	0.01000	53 5	31	0.00079	64 @	42	0.00006	75 K			
10	0.10000	43 +	21	0.00794	54 6	32	0.00063	65 A						

ASCII_BASE=64 Old Illumina

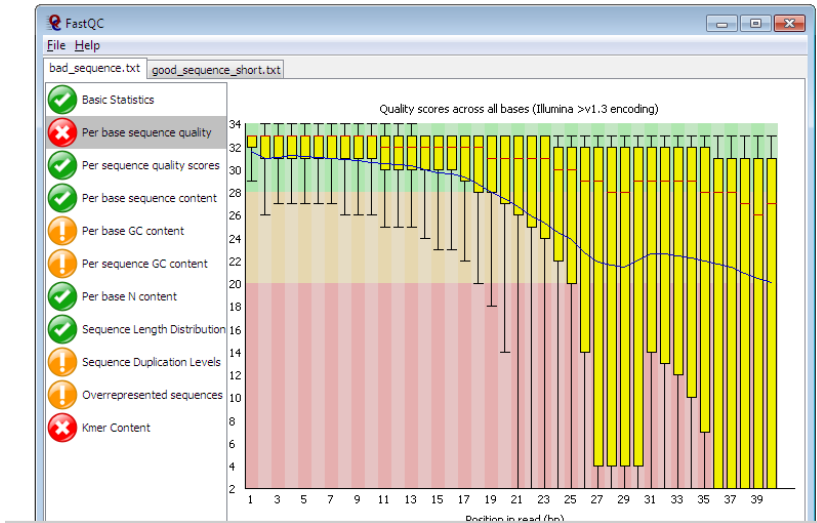
Q	P_error	ASCII	Q	P_error	ASCII	Q	P_error	ASCII	Q	P_error	ASCII
0	1.00000	64 @	11	0.07943	75 K	22	0.00631	86 V	33	0.00050	97 a
1	0.79433	65 A	12	0.06310	76 L	23	0.00501	87 W	34	0.00040	98 b
2	0.63096	66 B	13	0.05012	77 M	24	0.00398	88 X	35	0.00032	99 c
3	0.50119	67 C	14	0.03981	78 N	25	0.00316	89 Y	36	0.00025	100 d
4	0.39811	68 D	15	0.03162	79 O	26	0.00251	90 Z	37	0.00020	101 e
5	0.31623	69 E	16	0.02512	80 P	27	0.00200	91 [38	0.00016	102 f
6	0.25119	70 F	17	0.01995	81 Q	28	0.00158	92 \	39	0.00013	103 g
7	0.19953	71 G	18	0.01585	82 R	29	0.00126	93]	40	0.00010	104 h
8	0.15849	72 H	19	0.01259	83 S	30	0.00100	94 ^	41	0.00008	105 i
9	0.12589	73 I	20	0.01000	84 T	31	0.00079	95 _	42	0.00006	106 j
10	0.10000	74 J	21	0.00794	85 U	32	0.00063	96 `			

Pipeline: Cp assemblies

- I. Lib prep
- II. Sequencing - Demultiplexing
- III. QC
 - a) Sequencing Report
 - ✓ QC10: 1-10 error (90% accurate)
 - ✓ QC20: 1-100 error (99% accurate)
 - ✓ QC30: 1-1000 error (99,9% accurate)
 - b) Fastqc: <https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>
 - ✓ Quick overview of all reads
 - ✓ Summary graphs
 - ✓ HTML-based report
 - ✓ Multiqc! Overview of multiple HTMP-reports from different analyses
 - c) Trim
 - ✓ Trimmomatic: <http://www.usadellab.org/cms/?page=trimmomatic>
 - ✓ Cutadapt: <https://cutadapt.readthedocs.io/en/stable/>
- IV. Hybpiper
 - ✓ Toolkit designed for targeted sequence capture: <https://github.com/mossmatters/HybPiper>

Quality control – Fastqc

<https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>



Get QC info through fastqc:

- Download test data and scripts

```
$ git clone https://github.com/MycoMatics/intro-cpgenomes.git
```

```
$ cd intro-cpgenomes
```

```
$ kinit yourusername@UGENT.be
```

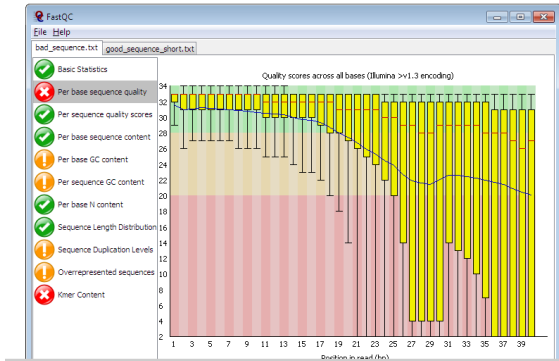
```
$ cp /UGent/yourusername/shares/data_hub_cemofe/Courses/illumina-QC/Illumina_data/OX1.tar.gz
```

- Take a look at the fastq.sh script

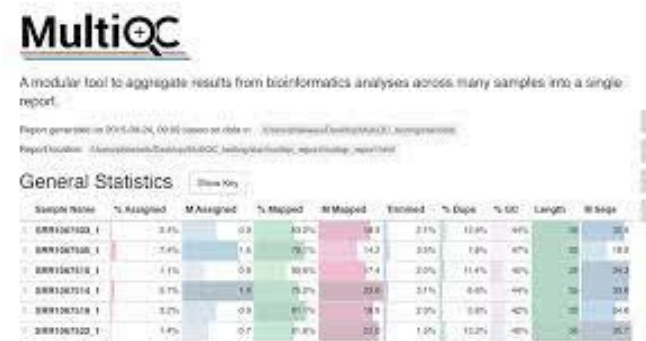
Quality control – Fastqc

```
1  #!/usr/bin/bash
2  # Line 1 is a she-bang that indicates that this is a Bash script.
3
4  #PBS -N fastqc.$PBS_JOBID          # Line 3-7 inform the scheduler about the resources required by this job:
5  #PBS -l nodes=1:ppn=all             # singe node (nodes=1) all core (ppn=all)
6  #PBS -o stdout.$PBS_JOBID          # redirect sterr stdout to separate files
7  #PBS -e stderr.$PBS_JOBID
8  #PBS -l walltime=01:00:00           # run for at most 2 minutes (walltime=00:02:00 max is 72hours)
9  #PBS -m abe                         # send mail when job (a)bort (b)egin (e)nd
10 #PBS -M <youremailaddresshere>      # specify your email address here
11
12
13 #Request software
14
15 ml FastQC/0.11.9-Java-11
16 ml MultiQC/1.14-foss-2022a
17
18 #Stage in data: Go to your current working directory and make sure both your data and scrit is there
19
20 cd $PBS_O_WORKDIR
21
22 #Make output directories
23 mkdir fastqc-reports
24
25 #Software commands
26 fastqc ./OX0001/*.fq.gz -o ./fastqc-reports && # run fastqc on all present fq.gz files
27 multiqc . # run multiqc on all fastqc output && in previous command prevents from starting multiqc early
28
```

Quality control – Fastqc + MultiQC



<https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>



<https://multiqc.info/>

Get QC info through fastqc:

- Download test data and scripts

```
$ git clone https://github.com/MycoMatics/intro-cpgenomes.git
```

```
$ cd intro-cpgenomes
```

```
$ kinit yourusername@UGENT.be
```

```
$ cp /UGent/yourusername/shares/data_hub_cemofe/Courses/illumina-QC/Illumina_data/OX1.tar.gz
```

- Run the fastq.sh script

```
$ qsub fastq.sh
```

Quality controle – Fastqc

- **Basic statistics**
- **Per base sequence quality**
 - => Range of quality values across all bases in each position
- **Per sequence quality scores**
 - => Quality distribution over all sequences
- **Per base sequence content**
 - => IF GC content ~50% => 25% chance of either A,T,G or C
 - => Bias at start: PCR duplicates, adapter/primer contamination
- **Per sequence GC content**
 - => GC distribution over all sequences => contamination indication other organisms
- **Per base N content**
- **Sequence length distribution**
- **Sequence duplication levels**
 - => Most sequences occur once (WGS)
 - Low level duplication: high level coverage throughout whole genome
 - High level duplication: enrichment bias
- **Overrepresented sequences**
- **Adapter content**

Pipeline: Cp assemblies

- I. Lib prep
- II. Sequencing - Demultiplexing
- III. QC
 - a) Sequencing Report
 - ✓ QC10: 1-10 error (90% accurate)
 - ✓ QC20: 1-100 error (99% accurate)
 - ✓ QC30: 1-1000 error (99,9% accurate)
 - b) Fastqc: <https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>
 - ✓ Quick overview of all reads
 - ✓ Summary graphs
 - ✓ HTML-based report
 - ✓ Multiqc! Overview of multiple HTML-reports from different analyses
 - c) Trim
 - ✓ Trimmomatic: <http://www.usadellab.org/cms/?page=trimmomatic>
 - ✓ Cutadapt: <https://cutadapt.readthedocs.io/en/stable/>
- IV. Cp – nr Assembly: Getorganelle
 - ✓ Toolkit for assemble of organelle genome: <https://github.com/Kinggerm/GetOrganelle>

Quality control: Trimming sequences

Simple Mode

The dark blue and red are the initial raw sequence that goes into Trimmomatic

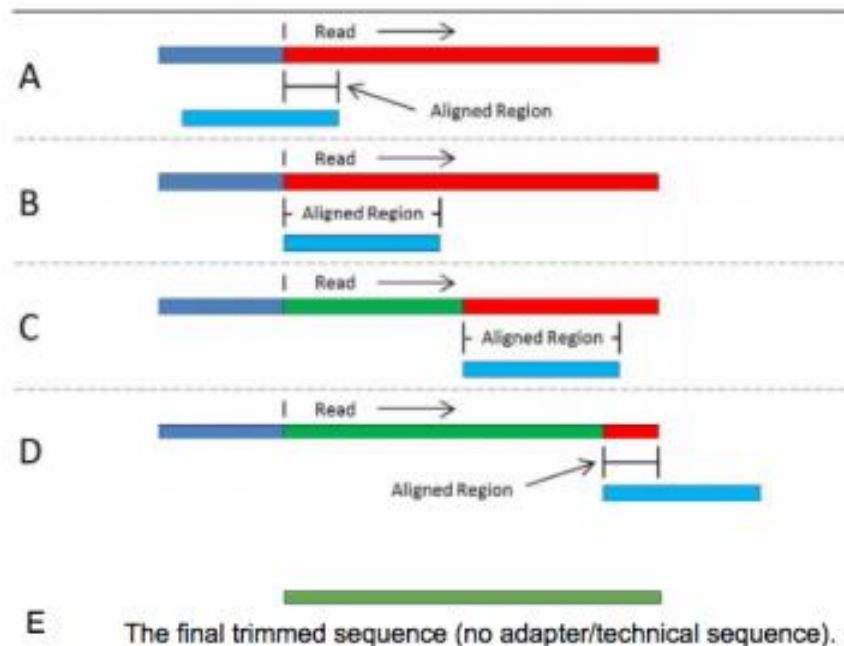


Figure modified from Bolger et al. 2014 (see link below). Caption from Bolger et al. 2014.

Fig. 1. Putative sequence alignments as tested in Simple Mode. The alignment process begins with partial overlap at the 5' end of the read (A), increasing to a full length 5' overlap (B), followed by full overlaps at all positions (C) and finishes with partial overlap at the 3' end of the read (D). Note that the upstream 'adapter' sequence is for illustration only, and is not part of the read or the aligned region.

Here the program "missed" removing the little bits that aligned well in A and D.



Quality control: Trimming sequences

Palindrome Mode

Figure modified from Bolger et al. 2014 (see link below). Caption from Bolger et al. 2014.

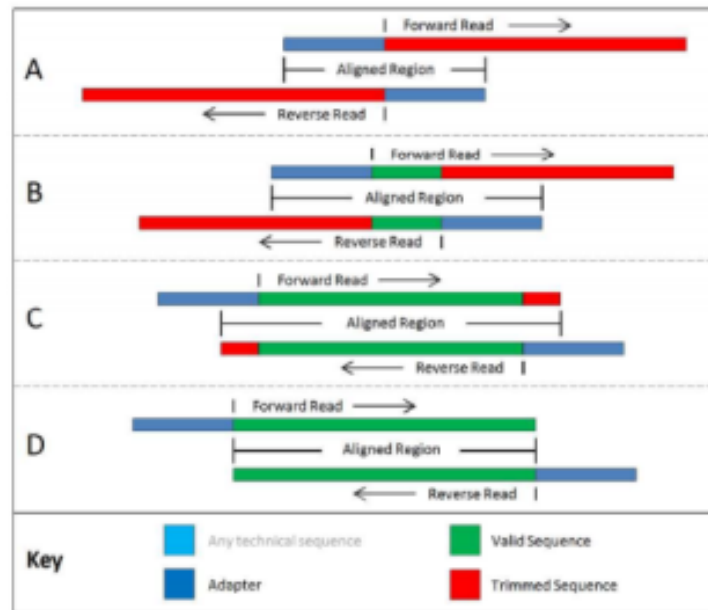


Fig. 2. Putative sequence alignments as tested in Palindrome Mode. The alignment process begins with the adapters completely overlapping the reads (A) testing for immediate 'read-through', then proceeds by checking for later overlap (B), including partial adapter read-through (C), finishing when the overlap indicates no read-through into the adapters (D).

<https://academic.oup.com/bioinformatics/article/30/15/2114/2390096/Trimmomatic-a-flexible-trimmer-for-Illumina>

Quality controle – Trimmomatic

<http://www.usadellab.org/cms/?page=trimmomatic>

Use trimmomatic to clip adapter sequences from your reads:

- Data structure folder intro-cp genomes
 - `adapters`
 - `alladapterstrimmomatic.fa`
 - `OX0001`
 - `_1.fq.gz` AND `_2.fq.gz`
 - `trimmomatic_v04.sh`
- take a look the trimomatic_v04.sh script

```

1  #!/bin/bash
2  #PBS -N trimmomatic
3  #PBS -l nodes=1:ppn=8
4  #PBS -o stdout.$PBS_JOBID
5  #PBS -e stderr.$PBS_JOBID
6  #PBS -l walltime=0:30:00
7  #PBS -m abe
8  #PBS -M pieter.asselman@ugent.be
9
10 # Data Paths
11 ILLUMINA_RAWDATA=/yourpath/to/OX0001 #take notice of the file extension names different options '.fq.gz' '.fastq.gz'
12 ILLUMINA_ADAPTERS=/yourpath/to/adapters
13
14 # make trimmed data directory
15 mkdir $ILLUMINA_RAWDATA/trimmed-data
16
17 # Create sampleslist to iterate, serves as input for iteration process in trimmomatic
18 for file in $ILLUMINA_RAWDATA/*_1.fq.gz
19 do
20     sample_name=$(basename "$file" _1.fq.gz) # extract sample name
21     echo "${sample_name}" >> $ILLUMINA_RAWDATA/samples.txt
22 done &&
23
24 # Load modules
25 module load Trimmomatic/0.39-Java-11
26
27 ### Run trimmomatic SE OR PE reads #-out what you don't need!
28
29 while read p
30 do
31     ### ACTIVATE FOR SE data
32     #java -jar $EBROOTTRIMMOMATIC/trimmomatic-0.39.jar SE -phred33 \
33     #-trimlog $TRIMMOMATIC_OUT/$p".log" \
34     #$ILLUMINA_RAWDATA/$p"_R1.fastq.gz" \
35     #$TRIMMOMATIC_OUT/$p"_R1_trimmed.fastq" \
36     #ILLUMINACLIP:$ILLUMINA_ADAPTERS/alladapterstrimmomatic.fa:2:30:10:1:TRUE SLIDINGWINDOW:5:20
37
38     ### ACTIVATE FOR PE data
39     java -jar $EBROOTTRIMMOMATIC/trimmomatic-0.39.jar PE -phred33 -trimlog $ILLUMINA_RAWDATA/trimmed-data/$p".log" \
40     $ILLUMINA_RAWDATA/$p"_1.fq.gz" \
41     $ILLUMINA_RAWDATA/$p"_2.fq.gz" \
42     $ILLUMINA_RAWDATA/trimmed-data/$p"_1_trimmedpaired.fastq" \
43     $ILLUMINA_RAWDATA/trimmed-data/$p"_1_trimmed_unpaired.fastq" \
44     $ILLUMINA_RAWDATA/trimmed-data/$p"_2_trimmedpaired.fastq" \
45     $ILLUMINA_RAWDATA/trimmed-data/$p"_2_trimmed_unpaired.fastq" \
46     ILLUMINACLIP:$ILLUMINA_ADAPTERS/alladapterstrimmomatic.fa:2:30:10:1:TRUE SLIDINGWINDOW:5:20
47
48 done < $ILLUMINA_RAWDATA/samples.txt
49
50
51

```

Quality controle – Trimmomatic

<http://www.usadellab.org/cms/?page=trimmomatic>

Use trimmomatic to clip adaptersequences from your reads:

- Data structure folder intro-cpgenomes
 - `adapters`
 - `alladapterstrimmomatic.fa`
 - `OX0001`
 - `_1.fq.gz` AND `_2.fq.gz`
 - `trimmomatic_v04.sh`
- Run the `trimomatic_v04.sh` script (make necessary adjustments to fit to your environment)
`$ qsub trimmomatic_v04.sh`

Quality control: Trimming sequences

```
module load Trimmomatic/0.39-Java-11
java -jar $EBROOTTRIMMOMATIC/trimmomatic-0.39.jar PE \
  input_forward.fq.gz \
  input_reverse.fq.gz \
  output_forward_paired.fq.gz \
  output_forward_unpaired.fq.gz \
  output_reverse_paired.fq.gz \
  output_reverse_unpaired.fq.gz \
  ILLUMINACLIP:<path>\adapters.fa:2:30:10:2:keepBothReads LEADING:3 TRAILING:3 MINLEN:36 SLIDINGWINDOW:4:15
```

This will perform the following:

- **Remove adapters** (ILLUMINACLIP:TruSeq3-PE.fa:2:30:10); adapter sequences can be found [here](#)
 - Allow maximally 2 mismatches.
 - extended seeds and clipped if score of 30 (about 50 bases) for PE (SE=10 ~17 bases)
 - Remove leading low quality or N bases (below quality 3) (LEADING:3)
- **Remove trailing low quality or N bases** (below quality 3) (TRAILING:3)
- **Scan the read with a 4-base wide sliding window**, cutting when the average quality per base drops below 15 (SLIDINGWINDOW:4:15)
- **Drop reads** below the 36 bases long (MINLEN:36)

Quality control: Trimming sequences

OUTPUT TRIMMOMATIC

Fastq files

- ✓ output_forward_paired.fq.gz
- ✓ output_**forward_unpaired**.fq.gz
- ✓ output_reverse_paired.fq.gz
- ✓ output_**reverse_unpaired**.fq.gz

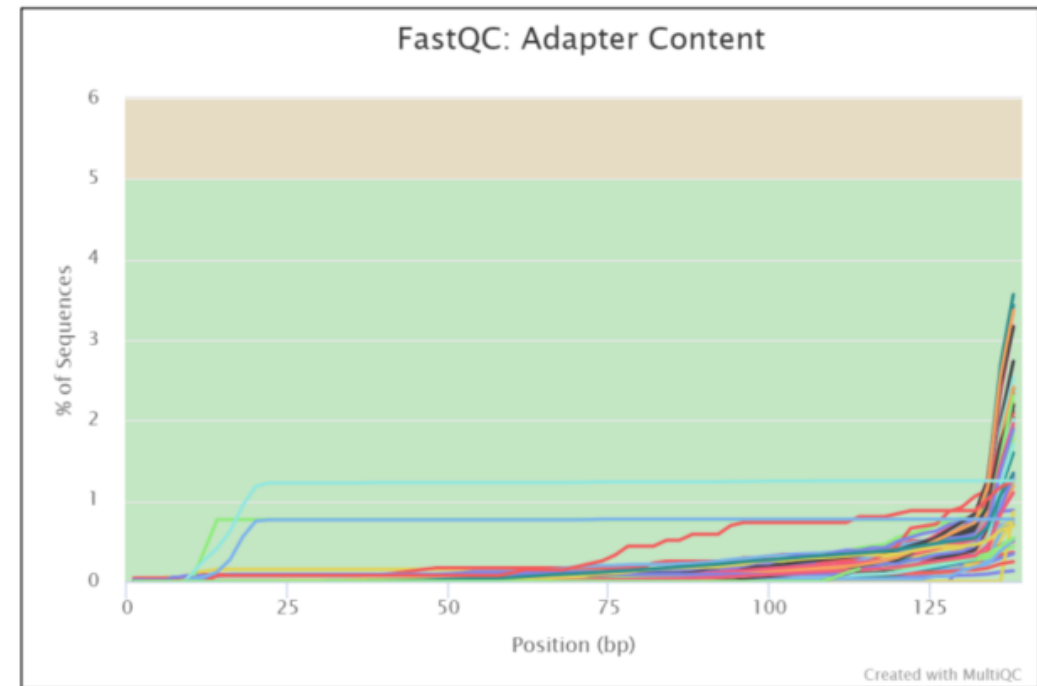
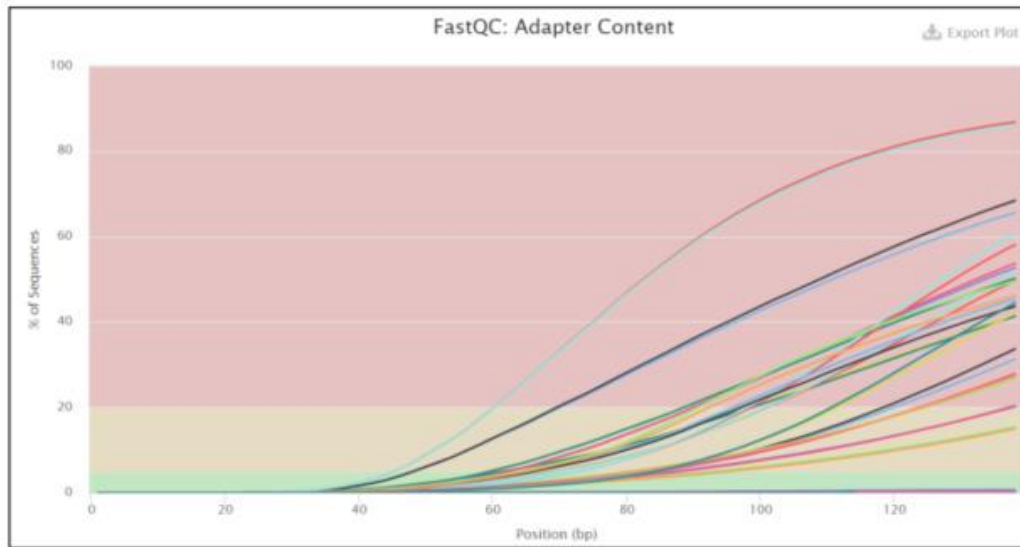


Unpaired reads can be mapped like
SE reads if necessary

Output logs

- ✓ Input Read Pairs: 16923155
- ✓ Both Surviving Reads: 13396826
- ✓ Both Surviving Read Percent: **79.16**
- ✓ Forward Only Surviving Reads: 3199726
- ✓ Forward Only Surviving Read Percent:
18.91
- ✓ Reverse Only Surviving Reads: 126909
- ✓ Reverse Only Surviving Read Percent:
0.75
- ✓ Dropped Reads: 199694
- ✓ Dropped Read Percent: 1.18

Quality control: Trimmomatic (before and after)

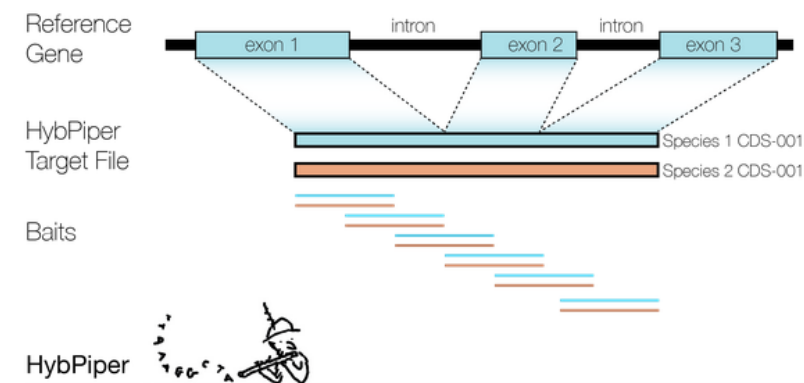


Pipeline: Cp assemblies

- I. Lib prep
- II. Sequencing - Demultiplexing
- III. QC
 - a) Sequencing Report
 - ✓ QC10: 1-10 error (90% accurate)
 - ✓ QC20: 1-100 error (99% accurate)
 - ✓ QC30: 1-1000 error (99,9% accurate)
 - b) Fastqc: <https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>
 - ✓ Quick overview of all reads
 - ✓ Summary graphs
 - ✓ HTML-based report
 - ✓ Multiqc! Overview of multiple HTML-reports from different analyses
 - c) Trim
 - ✓ Trimmomatic: <http://www.usadellab.org/cms/?page=trimmomatic>
 - ✓ Cutadapt: <https://cutadapt.readthedocs.io/en/stable/>
- IV. Cp – nr Assembly: Getorganelle
 - ✓ Toolkit for assemble of organelle genome: <https://github.com/Kinggerm/GetOrganelle>

Hybpiper 2

Target file and bait design (pre-HybPiper)

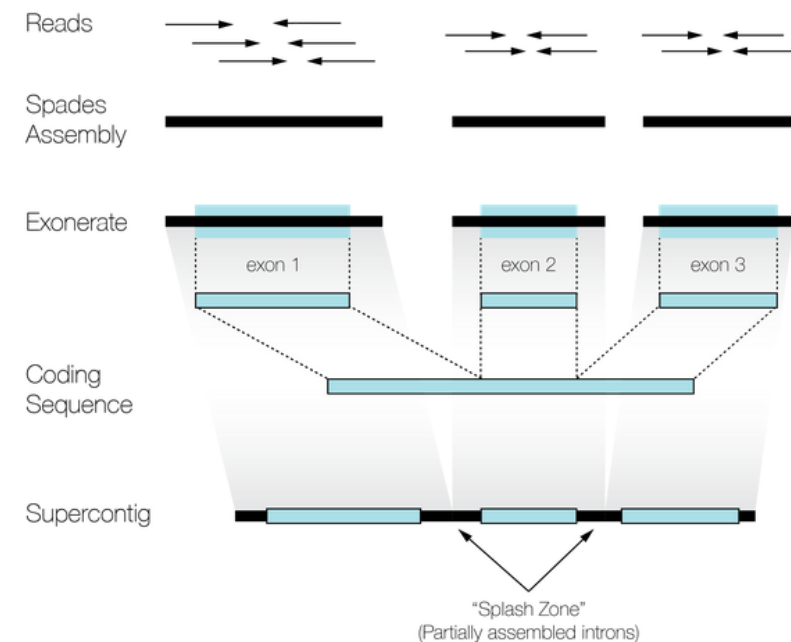


Target File

Targets are the complete sequence to be recovered. The target file includes targeted genes, and may include multiple orthologous sequences (e.g. multiple taxa) per gene.

Baits

Baits are short RNA sequences that hybridize with DNA library fragments during target enrichment. They may be designed from mature transcripts (as shown here) or from individual exons.



hybpiper assemble

Reads are searched against the target file and sorted according to the target gene:

`distribute_reads_to_targets.py`

The appropriate target gene is identified as the reference gene:

`distribute_targets.py`

Reads into contigs with SPAdes, optimized for single-gene assembly:

`spades_runner.py`

After assembly, SPAdes contigs are aligned to the reference, scaffolded, and translated. Intron sequences and supercontigs (scaffolded/merged SPAdes contigs) are generated:

`exonerate_hits.py`

hybpiper check_targets

Checks target file for proper formatting and flags low complexity sequences.

hybpiper stats

Summarizes gene recovery from multiple samples including sequence lengths and number of paralogs

hybpiper retrieve_sequences

Retrieve sequences generated from multiple runs of HybPiper

hybpiper recovery_heatmap

generate gene recovery heatmap

hybpiper paralog_retriever

Retrieve paralog sequences for a given gene, for all samples



Appl Plant Sci. 2016 Jul; 4(7): apps.1600016.

Published online 2016 Jul 12. doi: [10.3732/apps.1600016](https://doi.org/10.3732/apps.1600016)

PMCID: PMC4948903

PMID: [27437175](https://pubmed.ncbi.nlm.nih.gov/27437175/)

HybPiper: Extracting coding sequence and introns for phylogenetics from high-throughput sequencing reads using target enrichment¹

Matthew G. Johnson,^{2,6} Elliot M. Gardner,^{2,3} Yang Liu,⁴ Rafael Medina,⁴ Bernard Goffinet,⁴ A. Jonathan Shaw,⁵ Nyree J. C. Zerega,^{2,3} and Norman J. Wickett^{2,3}

► Author information ► Article notes ► Copyright and License information ► [PMC Disclaimer](#)

Read the paper: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4948903/>

Hybpiper

<https://github.com/mossmatters/HybPiper>

Run hybpiper on the test data.

The tool is extremely well documented. This is an exercise for home.

- <https://github.com/mossmatters/HybPiper>
- Download the test data
- Run on the HPC
- Let's take a look at the script


```

1  #!/usr/bin/bash
2  #PBS -N hybpiper_conda_tutorialdata
3  #PBS -l nodes=1:ppn=8
4  #PBS -o hybpiper.$PBS_JOBID.stdout
5  #PBS -e hybpiper.$PBS_JOBID.Stderr
6  #PBS -l walltime=01:00:00
7  #PBS -m abe
8  #PBS -M pieter.asselman@ugent.be
9
10 cd $PBS_O_WORKDIR
11
12 source ~/.bashrc
13 conda activate hybpiper &&
14
15 # Unpack the test dataset
16 tar -zxvf test_reads.fastq.tar.gz
17
18 # Remove any previous runs
19 parallel rm -r {} ::: namelist.txt
20
21
22 # Run main HybPiper command with all available CPUs
23 while read sample_name
24 do
25     hybpiper assemble -r ${sample_name}*.fastq -t_dna test_targets.fasta --prefix ${sample_name} --bwa --run_intronerate
26 done < namelist.txt
27
28
29 # Get runs statistics
30 hybpiper stats -t_dna test_targets.fasta gene namelist.txt
31
32
33 # Get heatmap of length recovery
34 hybpiper recovery_heatmap seq_lengths.tsv
35
36 # Recover DNA and amino-acid sequences
37 hybpiper retrieve_sequences -t_dna test_targets.fasta dna --sample_names namelist.txt --fasta_dir 01_dna_seqs
38 hybpiper retrieve_sequences -t_dna test_targets.fasta aa --sample_names namelist.txt --fasta_dir 02_aa_seqs
39
40
41 # Recover paralog sequences
42 hybpiper paralog_retriever namelist.txt -t_dna test_targets.fasta
43
44
45 echo "DONE!"

```

Hybpiper is installed in a conda env.
You first need to activate the env.

Download and prep data

Run hybpiper tools

Hybpiper

<https://github.com/mossmatters/HybPiper>

Run hybpiper on the test data.

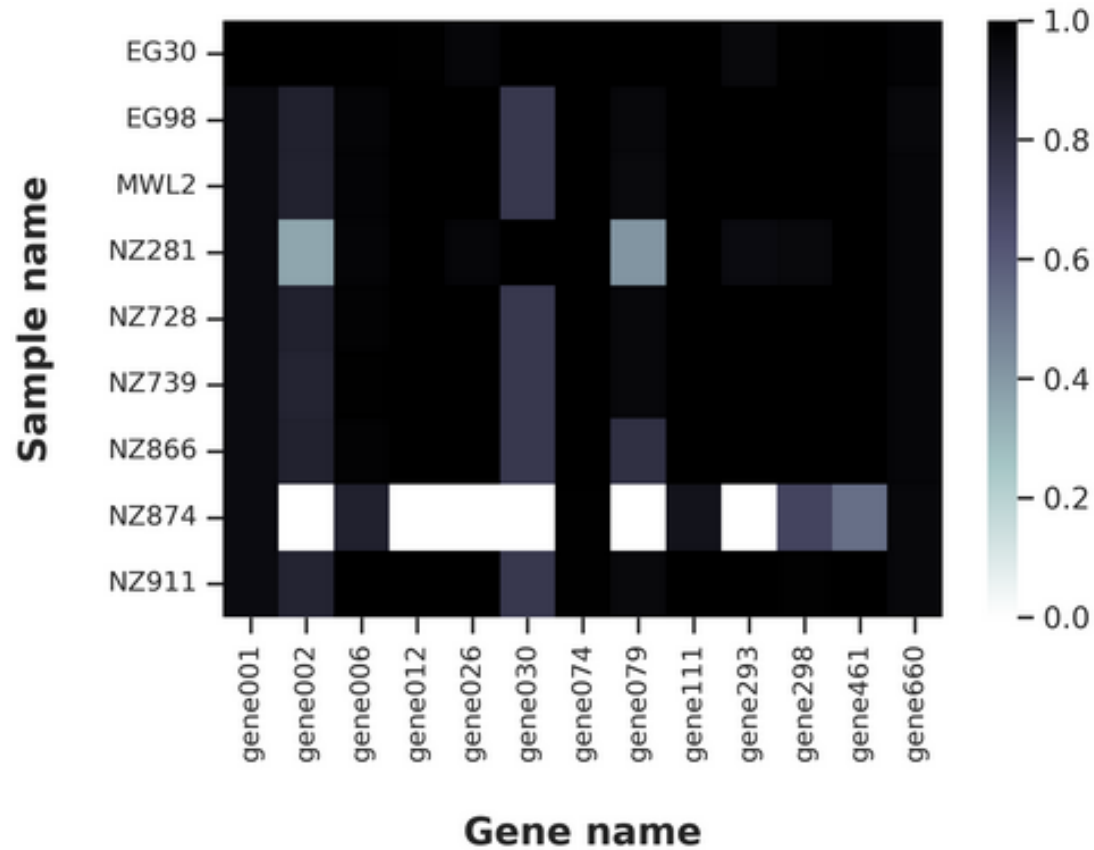
The tool is extremely well documented. This is an exercise for home.

- <https://github.com/mossmatters/HybPiper>
- Download the test data
- Run on the HPC
- Let's take a look at the script
- Output (more to discover than just a fancy heatmap, but just to get the gist of it)

Hybpiper

<https://github.com/mossmatters/HybPiper>

Percentage length recovery for each gene, relative to mean of targetfile references



Documentation

- Website: <https://www.ugent.be/hpc/en>
- HPC & linux documentation: <https://www.ugent.be/hpc/en/support/documentation.htm>
- Open Stack Dashboard: <https://login.hpc.ugent.be>
- Cluster state info: <https://shieldon.ugent.be:8083/pbsmon-web-users/>
- Software installation: <https://www.ugent.be/hpc/en/support/software-installation-request>

Tip: DIY installations 'easybuild' see Chp 28

- Flemish compute center: <https://www.vscentrum.be/>

[Need help with issues: hpc@ugent.be](mailto:hpc@ugent.be)

Introduction HPC-UGent

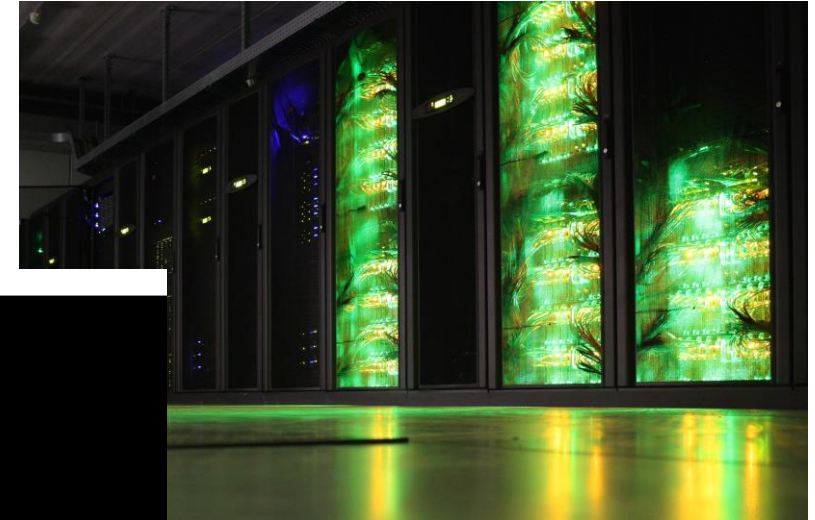
Questions?
Comments?
Suggestions?
Concerns?

```
1  #!/bin/bash
2
3  #PBS -S /bin/bash
4  #PBS -l walltime=06:00:00
5  #PBS -N JOBNAME
6  #PBS -o stdout.$PBS_JOBID
7  #PBS -e stderr.$PBS_JOBID
8  #PBS -l nodes=1:ppn=all:gpus=2
9  #PBS -m abe
10 #PBS -M pieter.asselman@ugent.be
11
12 #####
13 ##INFO##
14 #####
15
16 #DATArun on 2GPU/ 11Gbases - ~3h
17
18 ##BASECALL##
19 guppy_basecaller -i /scratch/gent/vo/000/gv
20
21 ##DEMULPLEX##
22 -guppy_barcode --i /scratch/gent/vo/000/gvo00058/v
23
24 ##BASECALL and DEMULPLEX##
25 -guppy_basecaller -i /scratch/gent/vo/000/gvo00058/v
26
```

```
https://www.ugent.be/hpc/en/infrastructure/status

We switched to new job command wrappers for all HPC-UGent Tier-2 clusters on Wed
June 9th 2021 at 17:22.
This switch should be transparent: you don't need to change your workflow or job
scripts.

If you notice any problems, or if any of these changes affect your work, please
contact hpc@ugent.be .
vsc43352@gligar04:~$
```



<https://ugent.be/hpc>
hpc@ugent.be

DEMO – file permissions

```
# ls -l file
-rw-r--r-- 1 root root 0 Nov 19 23:49 file
```

Diagram illustrating the breakdown of the permissions `-rw-r--r--`:

- File type:** `-` (regular file)
- Owner (rw-):** `rw` (Read, Write)
- Group (r--):** `r--` (Read)
- Other (r--):** `r--` (Read)

Legend:

- `r` = Readable
- `w` = Writeable
- `x` = Executable
- `-` = Denied

UGO

`/user/gent/433/vsc43352/data_gent_vo/data_share_group`

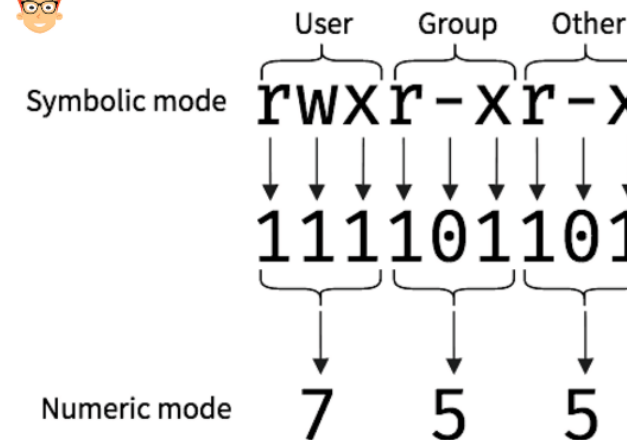
Examples:

`chmod u+rwx,g+rwx,o-rwx`

`chmod u=rwx,g=rwx`

`chgrp -R groupname file(or folder)`

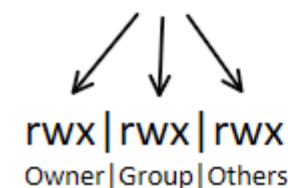
=> check for groups: `$ groups`



`drwxrwxrwx`

`d` = Directory
`r` = Read
`w` = Write
`x` = Execute

`chmod 777`



7	rwX	111
6	rw-	110
5	r-X	101
4	r--	100
3	-wX	011
2	-w-	010
1	--X	001
0	---	000