

Genome Assembly pipeline

This pipeline is intended to be a standard procedure to assembly genomes. We developed this pipeline to be optimal while using two kinds of data:

- PacBio sequel V2 genomic data
- Illumina genomic data

Step 0: Genome size prediction

Kmergenie estimates the best k-mer length for genome de novo assembly and it predicts the genome size. This is a useful tool in de novo genome assemblies when the genome size is unknown.

Installation:

```
1 module load anaconda2
2 conda create -n kmergenie
3 source activate kmergenie
4 conda install r-base
5 conda install -c bioconda kmergenie
6 source deactivate
```

Kmergenie needs a list with the link to illumina fasta files:

```
1 cd /data/SBCS-MartinDuranLab/03-Giacomo/data/osedax/illumina/
2 readlink -f OSE_DAX-1_S55_L003_R* > list_for_kmergenie
```

Code to run the command:

```
1 #!/bin/bash
2 #$ -wd /data/SBCS-MartinDuranLab/03-Giacomo/logs/kmergenie/
3 #$ -j y
4 #$ -pe smp 6
5 #$ -l h_vmem=8G
```

```
6  # $ -l h_rt=240:0:0
7
8  module load anaconda2
9  source activate kmergenie #load the environment
10
11 kmergenie \
12  /data/SBCS-MartinDuranLab/03-Giacomo/data/osedax/illumina/list_
   for_kmergenie \
13  --diploid \
14  -o osedax_histograms
```

Results will be in `/data/SBCS-MartinDuranLab/03-Giacomo/logs/kmergenie/`.
There will be an `.html` file that will contain the genome size prediction.

Step 1: Assembly

Step 1.1: Assemblers

The core of this pipeline is a genome assembler software such as [Canu](#) or [Flye](#). We tried three different approaches and then we evaluated the best one. First we runned Flye, then we launched Canu and finally we runned Flye using the corrected reads obtained with Canu.

Flye installation:

Installing Flye de novo assembler:

```
1  cd
2  module load anaconda2
3  conda create -n Flye python=2.7 anaconda
4  conda install flye
```

Canu installation:

- [Download](#) the latest release on your computer

```
1  scp canu-1.8.Linux-amd64.tar.xz btx604@login.hpc.qmul.ac.uk:/data/
   SBCS-MartinDuranLab/03-Giacomo/src
```

```
2 cd /data/SBCS-MartinDuranLab/03-Giacomo/src
3 tar -xJf canu-1.8.Linux-amd64.tar.xz
```

Converting subreads.bam to fastq with [BAM2fastx tools](#):
installation:

```
1 module load anaconda2
2 source activate Flye
3 conda install bam2fastx
```

Code to run BAM2fastx tools:

```
1 #!/bin/bash
2 #$ -wd /data/home/btx604/scripts/envs_jobs/Flye/
3 #$ -j y
4 #$ -pe smp 4
5 #$ -l h_vmem=1G
6 #$ -l h_rt=72:0:0
7
8 module load anaconda2
9 source activate Flye #load the environment
10 bam2fastq -o /data/SBCS-MartinDuranLab/03-Giacomo/raw_data/PacBio/01-Oasisia/Oasisia_pb_raw /data/SBCS-MartinDuranLab/03-Giacomo/raw_data/PacBio/01-Oasisia/data2/pb/r64044_20190812_215729/1_A01/m64044_190812_220643.subreads.bam
```

Code to run Flye:

```
1 #!/bin/bash
2 #$ -wd /data/SBCS-MartinDuranLab/03-Giacomo/raw_data/PacBio/01-Oasisia/qsub_logs/
3 #$ -j y
4 #$ -pe smp 48
5 #$ -l h_vmem=10G
```

```
6  #$ -l h_rt=240:0:0
7  #$ -l highmem
8
9  module load anaconda2
10 source activate Flye #load the environment
11
12 flye \
13   --pacbio-raw /data/SBCS-MartinDuranLab/03-Giacomo/raw_data/PacB
io/01-Oasisia/Oasisia_pb_raw.fastq.gz \
14   --genome-size 1g \
15   --threads 48 \
16   --asm-coverage 40 \
17   --resume \
18   --out-dir /data/scratch/btx604/Oasisia/flye/
```

Code to run Canu (without using the grid):

```
1  #!/bin/bash
2  #$ -wd /data/SBCS-MartinDuranLab/03-Giacomo/raw_data/PacBio/01-O
asisia/qsub_logs/
3  #$ -j y
4  #$ -pe smp 48
5  #$ -l h_vmem=10G
6  #$ -l h_rt=240:0:0
7  #$ -l highmem
8
9  module load gnuplot
10 module load java/1.8.0_152-oracle
11
12 /data/SBCS-MartinDuranLab/03-Giacomo/bin/canu \
13   -p Oasisia -d /data/scratch/btx604/Oasisia/canu \
14   genomeSize=1g \
```

```

15 maxMemory=480g \
16 maxThreads=48 \
17 useGrid=false \
18 gridEngineResourceOption="-l h_vmem=MEMORY -pe smp THREADS" \
19 batOptions="-dg 3 -db 3 -dr 1 -ca 500 -cp 50" \
20 -pacbio-raw /data/SBCS-MartinDuranLab/03-Giacomo/raw_data/PacBio/01-Oasisia/Oasisia_pb_raw.fastq.gz

```

Code to run Canu (using the grid): ????????

- First a Canu script needs to be modified accordingly to the cluster requirements: in the script `/data/SBCS-MartinDuranLab/03-Giacomo/src/canu-1.8/Linux-amd64/lib/site_perl/canu/Grid_SGE.pm`: change `#"-pe threads THREADS");` with `#"-pe smp THREADS");`; and `#"-l mem=MEMORY");` with `#"-l h_vmem=MEMORY");`

```

1 cd
2
3 module load gnuplot
4 module load java/1.8.0_152-oracle
5
6 /data/SBCS-MartinDuranLab/03-Giacomo/bin/canu \
7 -p Oasisia -d /data/scratch/btx604/Oasisia/canu_V8 \
8 -pacbio-raw /data/SBCS-MartinDuranLab/03-Giacomo/raw_data/PacBio/01-Oasisia/Oasisia_pb_raw.fastq.gz \
9 genomeSize=1g \
10 useGrid=true \
11 gridEngineResourceOption="-l h_vmem=MEMORY -pe smp THREADS" \
12 gridOptions="-l h_rt=72:0:0 -j y -l highmem" \
13 "batOptions=-dg 3 -db 3 -dr 1 -ca 500 -cp 50" \
14 "gridOptionsGFA=-l h_vmem=300G" \
15 "gridOptionsovb=-l h_vmem=120G" \
16 "gridOptionsovs=-l h_vmem=120G" \

```

```

17 "gridOptionscorovl=-l h_vmem=30G" \
18 "gridOptionscor=-l h_vmem=30G" \
19 "gridOptionscns=l h_vmem=40G" \
20 gnuplot=$(which gnuplot) \
21 java=$(which java)
22 > /data/scratch/btx604/Oasisia/canu_V8/LOG/LOG 2>&1

```

- it may happen that Canu crash before starting the assembly, in this case the assembly can be launched with a separated off grid job:

```

1 #!/bin/bash
2 #$ -wd /data/scratch/btx604/Oasisia/canu_V8
3 #$ -o /data/scratch/btx604/Oasisia/canu_V8
4 #$ -j y
5 #$ -pe smp 48
6 #$ -l h_vmem=10G
7 #$ -l h_rt=240:0:0
8 #$ -l highmem
9
10 module load gnuplot
11 module load java/1.8.0_152-oracle
12
13 /data/SBCS-MartinDuranLab/03-Giacomo/bin/canu \
14 -p Oasisia -d /data/scratch/btx604/Oasisia/canu_V8 \
15 genomeSize=750m \
16 maxMemory=480g \
17 maxThreads=48 \
18 useGrid=false \
19 gridEngineResourceOption="-l h_vmem=MEMORY -pe smp THREADS" \
20 batOptions="-dg 3 -db 3 -dr 1 -ca 500 -cp 50" \
21 -pacbio-raw /data/SBCS-MartinDuranLab/03-Giacomo/raw_data/PacBi
o/01-Oasisia/Oasisia_pb_raw.fastq.gz

```

Code to run Flye using Canu corrected reads:

```
1  #!/bin/bash
2  #$ -wd /data/SBCS-MartinDuranLab/03-Giacomo/raw_data/PacBio/01-0
   asisia/qsub_logs/
3  #$ -j y
4  #$ -pe smp 48
5  #$ -l h_vmem=10G
6  #$ -l h_rt=240:0:0
7  #$ -l highmem
8
9  module load anaconda2
10 source activate Flye #load the environment
11
12 flye \
13   --pacbio-corr /data/SBCS-MartinDuranLab/03-Giacomo/01-0asisiaDa
   ta/0asisia.correctedReads.fasta.gz \
14   --genome-size 1g \
15   --threads 48 \
16   --asm-coverage 40 \
17   --out-dir /data/scratch/btx604/0asisia/flye_corrected_reads_v2/
```

Step 1.2: Quality check

Busco provides quantitative measures for the assessment of genome assembly, gene set, and transcriptome completeness, based on evolutionarily-informed expectations of gene content from near-universal single-copy orthologs selected from OrthoDB v9

Code to run Busco on the Flye assembly:

```
1  #!/bin/bash
2  #$ -wd /data/SBCS-MartinDuranLab/03-Giacomo/logs/busco
```

```
3  #$ -pe smp 4
4  #$ -l h_vmem=2G
5  #$ -l h_rt=10:0:0
6  #$ -cwd
7  #$ -j y
8
9  cd /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia/busco/flye/
10
11  module load busco/3.0
12  module load augustus
13
14  export AUGUSTUS_CONFIG_PATH=/data/SBCS-MartinDuranLab/02-Chema/s
rc/Augustus/config
15
16  BUSCO.py -i /data/scratch/btx604/Oasisia/flye/assembly.fasta -m
genome -o Oasisia_Flye_Busco -c 4 -l /data/SBCS-MartinDuranLab/0
0-BlastDBs/metazoa_odb9
```

Code to run Busco on the Canu assembly:

```
1  #!/bin/bash
2  #$ -wd /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia/busco/c
anu/
3  #$ -pe smp 4
4  #$ -l h_vmem=2G
5  #$ -l h_rt=10:0:0
6  #$ -cwd
7  #$ -j y
8
9  cd /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia/busco/canu/
10
11  module load busco/3.0
```



```
12 module load augustus
13
14 export AUGUSTUS_CONFIG_PATH=/data/SBCS-MartinDuranLab/02-Chema/s
rc/Augustus/config
15
16 BUSCO.py -i /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia/ca
nu/Oasisia.contigs.fasta -m genome -o Oasisia_Canu_Busco -c 4 -l
/data/SBCS-MartinDuranLab/00-BlastDBs/metazoa_odb9
```

Code to run Busco on the Flye (corrected reads) assembly:

```
1 #!/bin/bash
2 #$ -wd /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia/busco/c
anu/
3 #$ -pe smp 4
4 #$ -l h_vmem=2G
5 #$ -l h_rt=10:0:0
6 #$ -cwd
7 #$ -j y
8
9 cd /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia/busco/canu/
10
11 module load busco/3.0
12 module load augustus
13
14 export AUGUSTUS_CONFIG_PATH=/data/SBCS-MartinDuranLab/02-Chema/s
rc/Augustus/config
15
16 BUSCO.py -i /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia/fl
ye_corrected_reads_v2/assembly.fasta -m genome -o Oasisia_Flye_c
orrected_reads_Busco -c 4 -l /data/SBCS-MartinDuranLab/00-BlastD
Bs/metazoa_odb9
```

Step 1.3: Assemblies comparison

QUAST quality assessment tool evaluates and compares genome assemblies. This tool improves on leading assembly comparison software with quality metrics. QUAST can evaluate assemblies both with a reference genome, as well as without a reference. QUAST produces many reports, summary tables and plots.

Installation:

```
1 cd
2 module load anaconda2
3 conda create -n quast
4 source activate quast
5 conda install -c bioconda quast
6 source deactivate
```

Code to run Quast:

```
1 #!/bin/bash
2 #$ -wd /data/SBCS-MartinDuranLab/03-Giacomo/logs/quast
3 #$ -pe smp 4
4 #$ -l h_vmem=2G
5 #$ -l h_rt=10:0:0
6 #$ -cwd
7 #$ -j y
8
9 module load anaconda2
10 source activate quast
11
12 quast \
13 /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia/canu/Oasisia.
14 contigs.fasta \
15 /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia/flye/assembly.
16 fasta \
17 /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia/flye_correcte
```

```
d_reads_v2/assembly.fasta \  
16 -o /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia/quast_canu  
_flye_flye \  
17 --eukaryote \  

```

Step 2: Polishing

Raw PacBio reads and Illumina paired-end data will be aligned to the initial assembly using **Bwa** and **Pbmm2** and casual errors occurred during the PacBio sequencing will be corrected with **Arrow** and **Racon**. This step includes many small tasks, therefore it could be easier to use a wrapper script for all the jobs.

Step 2.1: Preparing the PacBio data

Seqtk to convert PacBio subreads from fastq.gz to fasta

Code to run Seqtk:

```
1 #!/bin/bash  
2 #$ -wd /data/home/btx604/scripts  
3 #$ -j y  
4 #$ -o /data/SBCS-MartinDuranLab/03-Giacomo/logs/various  
5 #$ -pe smp 1  
6 #$ -l h_vmem=10G  
7 #$ -l h_rt=24:0:0  
8  
9 module load seqtk  
10  
11 seqtk seq -a /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia/0  
0-pacbio/oasisia_pb_raw.fastq.gz > /data/SBCS-MartinDuranLab/03-  
Giacomo/data/oasisia/00-pacbio/oasisia_pb_raw.fa
```

Samtools faidx to index the fasta assembly, therefore making it accessible to Arrow

Code to run Samtools:

```
1  #!/bin/bash
2  #$ -wd /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia/canu/
3  #$ -o /data/SBCS-MartinDuranLab/03-Giacomo/logs/samtools/
4  #$ -j y
5  #$ -pe smp 1
6  #$ -l h_vmem=12G
7  #$ -l h_rt=24:0:0
8
9  module load samtools/1.9
10
11 samtools faidx /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia/
   /canu/Oasisia.contigs.fasta
```

Step 2.2: Polishing with PacBio data

In order to polish our assembly using PacBio data we need to run two rounds of a combination of different softwares such as pbmn2 (that aligns the raw reads to the assembly), samtools and **pbindex** (that generate all the indexes that are needed to run Arrow), and Arrow (which compare the assembly with the raw reads and fix the errors occurred during the sequencing)

First we need to install pbmn2, pbindex (pbbam) and Arrow

```
1  module load anaconda2
2  conda create -n genomePolishing python=2.7 anaconda
3  source activate genomePolishing
4  conda install -c bioconda pbmm2
5  conda install -c bioconda pbgcpp
6  conda install pbbam
```

First Round:

Align raw PB reads to the assembly in a bam file

Code to run pbmm2:

```
1  #!/bin/bash
2  #$ -wd /data/scratch/btx604/Oasisia/genomePolishing/pbalign/step
   1
3  #$ -j y
4  #$ -o /data/SBCS-MartinDuranLab/03-Giacomo/logs/genomePolishing
5  #$ -pe smp 5
6  #$ -l h_vmem=8G
7  #$ -l h_rt=72:0:0
8
9  cd /data/scratch/btx604/Oasisia/genomePolishing/pbalign/step1
10
11  module load anaconda2
12  source activate genomePolishing
13
14  pbmm2 align \
15  /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia/canu/Oasisia.
   contigs.fasta \
16  /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia/00-pacbio/dat
   a2/pb/r64044_20190812_215729/1_A01/m64044_190812_220643.subread
   s.bam \
17  /data/scratch/btx604/Oasisia/genomePolishing/pbalign/step1/oasi
   sia_pbalign_step1.bam
```

- Sorted output can be generated using `--sort`

Sorting and indexing the bam alignments

Code to run pbmm2:

```
1  #!/bin/bash
2  #$ -wd /data/scratch/btx604/Oasisia/genomePolishing/pbalign/step
   1/
3  #$ -o /data/SBCS-MartinDuranLab/03-Giacomo/logs/samtools/
4  #$ -j y
5  #$ -pe smp 1
```

```
6  #$ -l h_vmem=12G
7  #$ -l h_rt=24:0:0
8
9  module load samtools/1.9
10
11  samtools sort /data/scratch/btx604/Oasisia/genomePolishing/pbalign/step1/oasisia_pbalign_step1.bam -o /data/scratch/btx604/Oasisia/genomePolishing/pbalign/step1/oasisia_sorted_step1.bam
12  samtools index /data/scratch/btx604/Oasisia/genomePolishing/pbalign/step1/oasisia_sorted_step1.bam
```

Further indexing of the bam file

Code to run pbindex:

```
1  #!/bin/bash
2  #$ -wd /data/scratch/btx604/Oasisia/genomePolishing/pbalign/step1
3  #$ -j y
4  #$ -o /data/SBCS-MartinDuranLab/03-Giacomo/logs/genomePolishing
5  #$ -pe smp 1
6  #$ -l h_vmem=20G
7  #$ -l h_rt=120:0:0
8
9  module load anaconda2
10  source activate genomePolishing
11
12  cd /data/scratch/btx604/Oasisia/genomePolishing/pbalign/step1/
13
14  pbindex /data/scratch/btx604/Oasisia/genomePolishing/pbalign/step1/oasisia_sorted_step1.bam
```

Finally launch Arrow, the polishing tool

Code to run Arrow:

```
1  #!/bin/bash
2  #$ -wd /data/scratch/btx604/Oasisia/genomePolishing/pbalign/step
   1
3  #$ -j y
4  #$ -o /data/SBCS-MartinDuranLab/03-Giacomo/logs/genomePolishing
5  #$ -pe smp 4
6  #$ -l h_vmem=20G
7  #$ -l h_rt=120:0:0
8
9  module load anaconda2
10 source activate genomePolishing
11
12 cd /data/scratch/btx604/Oasisia/genomePolishing/pbalign/step1
13
14 arrow /data/scratch/btx604/Oasisia/genomePolishing/pbalign/step1
   /oasisia_sorted_step1.bam \
15 -r /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia/canu/Oasis
   ia.contigs.fasta \
16 -o oasisia_step1_variants.gff \
17 -o oasisia_step1_consensus.fasta \
18 -o oasisia_step1_consensus.fastq \
19 -j 4
```

Second Round:

First we need to generate another .fai index of the oasisia_step1_consensus.fasta file
Code to run Samtools:

Align raw PB reads to the fasta output of round1 Arrow

Code to run pbmm2:

Sorting and indexing the bam alignments

Code to run pbmm2:

Further indexing of the bam file

Code to run pbindex:

Launch Arrow for the second time

Code to run Arrow:

Preparing a wrapper script:

create folders: /data/scratch/btx604/Oasisia/genomePolishing/PB

```
1 /data/scratch/btx604/Oasisia/genomePolishing/PB/
2 /data/scratch/btx604/Oasisia/genomePolishing/PB/step1/
3 /data/scratch/btx604/Oasisia/genomePolishing/PB/step2/
4 /data/scratch/btx604/Oasisia/genomePolishing/PB/pacbio_reads/
5 /data/scratch/btx604/Oasisia/genomePolishing/PB/reference_genome
/
```

Wrap script for the PacBio polishing phase (this is for osedax!)

```
1 #!/bin/bash
2 #$ -wd /data/scratch/btx604/osedax/genomePolishing/PB
3 #$ -j y
4 #$ -o /data/scratch/btx604/osedax/genomePolishing/PB
5 #$ -pe smp 8
6 #$ -l h_vmem=10G
7 #$ -l h_rt=240:0:0
8
9 #seqtk variables
10 raw_pb_reads_fq=/data/SBCS-MartinDuranLab/03-Giacomo/data/osedax
  /00-pacbio/osedax_pb_raw.fastq.gz
11 raw_pb_reads_fa=/data/scratch/btx604/osedax/genomePolishing/PB/p
  acbio_reads/osedax_pb_raw.fa
12 #faidx variables
13 reference_genome_step1=/data/scratch/btx604/osedax/genomePolishi
```



```
ng/PB/reference_genome/0sedax.contigs.fasta
14 #pbmm2 step1 variables
15 pb_subreads_bam=/data/SBCS-MartinDuranLab/03-Giacomo/data/osedax
    /00-pacbio/data2/pb/r64044_20191025_002245/2_B01/m64044_191025_1
    55828.subreads.bam
16 alignment_step1=/data/scratch/btx604/osedax/genomePolishing/PB/s
    tep1/osedax_pbmm2_step1.bam
17 #samtools step1 variables
18 alignment_step1_sorted=/data/scratch/btx604/osedax/genomePolishi
    ng/PB/step1/osedax_sorted_step1.bam
19 #arrow step1 variables
20 step1_variants=/data/scratch/btx604/osedax/genomePolishing/PB/st
    ep1/osedax_step1_variants.gff
21 step1_consensus_fa=/data/scratch/btx604/osedax/genomePolishing/P
    B/step1/osedax_step1_consensus.fasta
22 step1_consensus_fq=/data/scratch/btx604/osedax/genomePolishing/P
    B/step1/osedax_step1_consensus.fastq
23
24 #pbmm2 step2 variables
25 alignment_step2=/data/scratch/btx604/osedax/genomePolishing/PB/s
    tep2/osedax_pbmm2_step2.bam
26 #samtools step2 variables
27 alignment_step2_sorted=/data/scratch/btx604/osedax/genomePolishi
    ng/PB/step2/osedax_sorted_step2.bam
28 #arrow step2 variables
29 step2_variants=/data/scratch/btx604/osedax/genomePolishing/PB/st
    ep2/osedax_step2_variants.gff
30 step2_consensus_fa=/data/scratch/btx604/osedax/genomePolishing/P
    B/step2/osedax_step2_consensus.fasta
31 step2_consensus_fq=/data/scratch/btx604/osedax/genomePolishing/P
    B/step2/osedax_step2_consensus.fastq
32
33
```

```
34 module load seqtk
35 module load samtools/1.9
36 module load anaconda2
37 source activate genomePolishing
38
39
40 echo 'SEQTK _____
    _____'
41 if [ -e "$raw_pb_reads_fa" ]
42 then
43     echo "$raw_pb_reads_fa found."
44 else
45     seqtk seq -a $raw_pb_reads_fq > $raw_pb_reads_fa
46 fi
47
48 echo 'FAIDX STEP1_____
    _____'
49 if [ -e /data/scratch/btx604/osedax/genomePolishing/PB/reference
    _genome/*.fai ]
50 then
51     echo "/data/scratch/btx604/osedax/genomePolishing/PB/reference
    _genome/*.fai found."
52 else
53     samtools faidx $reference_genome_step1
54 fi
55
56 echo 'PBMM2 STEP1_____
    _____'
57 if [ -e "$alignment_step1" ]
58 then
59     echo "$alignment_step1 found."
```

```
60 else
61     pbmm2 align $reference_genome_step1 $pb_subreads_bam $alignmen
t_step1
62 fi
63
64 echo 'SAMTOOLS STEP1_____
        _____'
65 if [ -e "$alignment_step1_sorted" ]
66 then
67     echo "$alignment_step1_sorted found."
68 else
69     samtools sort $alignment_step1 -o $alignment_step1_sorted
70     samtools index $alignment_step1_sorted
71 fi
72
73 echo 'PBINDEX STEP1_____
        _____'
74 if [ -e /data/scratch/btx604/osedax/genomePolishing/PB/step1/*.p
bi ]
75 then
76     echo "/data/scratch/btx604/osedax/genomePolishing/PB/step1/*.p
bi found."
77 else
78     pbindex $alignment_step1_sorted
79 fi
80
81 echo 'ARROW STEP1_____
        _____'
82 if [ -e "$step1_consensus_fa" ]
83 then
84     echo "$step1_consensus_fa found."
```

```
85 else
86     arrow $alignment_step1_sorted -r $reference_genome_step1 -o $s
tep1_variants -o $step1_consensus_fa -o $step1_consensus_fq -j 8
87 fi
88
89 echo 'FAIDX STEP2_____
_____
'
90 if [ -e /data/scratch/btx604/osedax/genomePolishing/PB/step1/*.f
ai ]
91 then
92     echo "/data/scratch/btx604/osedax/genomePolishing/PB/step1/*.f
ai found."
93 else
94     samtools faidx $step1_consensus_fa
95 fi
96
97 echo 'PBMM2 STEP2_____
_____
'
98 if [ -e "$alignment_step2" ]
99 then
100     echo "$alignment_step2 found."
101 else
102     pbmm2 align $step1_consensus_fa $pb_subreads_bam $alignment_st
ep2
103 fi
104
105 echo 'SAMTOOLS STEP2_____
_____
'
106 if [ -e "$alignment_step2_sorted" ]
107 then
108     echo "$alignment_step2_sorted found."
109 else
```

```
110  samtools sort $alignment_step2 -o $alignment_step2_sorted
111  samtools index $alignment_step2_sorted
112  fi
113
114  echo 'PBINDEX STEP2_-----
-----'
115  if [ -e /data/scratch/btx604/osedax/genomePolishing/PB/step2/*.p
bi ]
116  then
117      echo "/data/scratch/btx604/osedax/genomePolishing/PB/step2/*.p
bi found."
118  else
119      pbindex $alignment_step2_sorted
120  fi
121
122  echo 'ARROW STEP2_-----
-----'
123  if [ -e "$step2_consensus_fa" ]
124  then
125      echo "$step2_consensus_fa found."
126  else
127      arrow $alignment_step2_sorted -r $step1_consensus_fa -o $step2
_variants -o $step2_consensus_fa -o $step2_consensus_fq -j 8
128  fi
```

Step 2.3: Preparing the Illumina data

Cutadapt is used to remove the adaptors employed during the Illumina sequencing from the raw data

Installation:

```
1  module load python/3.6.3
2  pip3 install --user --upgrade cutadapt
```

cutadapt will get installed in `~/.local` and can be run like this `~/.local/bin/cutadapt`

Code to run Cutadapt:

```
1  #!/bin/bash
2  #$ -cwd
3  #$ -j y
4  #$ -o /data/SBCS-MartinDuranLab/03-Giacomo/logs/cutadapt
5  #$ -pe smp 1
6  #$ -l h_vmem=12G
7  #$ -l h_rt=24:0:0
8
9  module load python/3.6.3
10 ~/.local/bin/cutadapt \
11   -a CTGTCTCTTATACACATCT \
12   -A CTGTCTCTTATACACATCT \
13   -m 30 \
14   -q 15,10 \
15   -o /data/scratch/btx604/Oasisia/cutadapt/02/Oasisia_illumina_R
16   1.fastq.gz \
17   -p /data/scratch/btx604/Oasisia/cutadapt/02/Oasisia_illumina_R
18   2.fastq.gz \
19   /data/SBCS-MartinDuranLab/03-Giacomo/01-OasisiaData/OT_S71_L002
20   _R1_001.fastq.gz \
21   /data/SBCS-MartinDuranLab/03-Giacomo/01-OasisiaData/OT_S71_L002
22   _R2_001.fastq.gz
```

- `-q 15,10 \` will quality-trim the 5' end with a cutoff of 15 and the 3' end with a cutoff of 10.
- `-m 30 \` Discard processed reads that are shorter than 30. Reads that are too short even before adapter removal are also discarded. Without this option, reads that have a length of zero (empty reads) are kept in the output.

Quality check of the cutadapted sequences with **FastQC**

Code to run fastQC:

```
1  #!/bin/bash
2  #$ -cwd
3  #$ -j y
4  #$ -o /data/SBCS-MartinDuranLab/03-Giacomo/logs/fastQC
5  #$ -pe smp 2
6  #$ -l h_vmem=1G
7  #$ -l h_rt=2:0:0
8
9  cd /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia/fastQC/
10 module load fastqc/0.11.5
11
12 fastqc /data/scratch/btx604/Oasisia/cutadapt/02/Oasisia_illumina
   _R1.fastq.gz
13 fastqc /data/scratch/btx604/Oasisia/cutadapt/02/Oasisia_illumina
   _R2.fastq.gz
```

Step 2.4: Polishing with Illumina data

Unique wrapper for illumina polishing. needed folders:

```
1  /data/scratch/btx604/$1/genomePolishing/illumina/raw_reads
2  /data/scratch/btx604/$1/genomePolishing/illumina/step1
3  /data/scratch/btx604/$1/genomePolishing/illumina/step2
```

- \$1 refers to the species name
- step2_consensus_fa (from PB wrapper) should be in step1
- illumina cutadapted reads named R1.fastq.gz and R2.fastq.gz should be in raw_reads

this script will take an input of the species name as the \$1 and should be launched like this:

```
qsub illumina_polishing_wrapper_v1.2_highmem.sh species
```

Before running the following script you need to ask the permission to use a high memory node to the IT service.

illumina_polishing_wrapper_v1.2_highmem.sh

```
1  #!/bin/bash
2  #$ -wd /data/scratch/btx604/
3  #$ -j y
4  #$ -o /data/scratch/btx604/
5  #$ -pe smp 24
6  #$ -l h_vmem=30G
7  #$ -l h_rt=240:0:0
8  #$ -l highmem
9
10 #bwa index step1 variables
11 reference_genome_step1_prefix="$1"_step2_consensus.fasta
12 reference_genome_step1=/data/scratch/btx604/$1/genomePolishing/i
    llumina/step1/$reference_genome_step1_prefix
13 #bwa mem step1 variables
14 R1=/data/scratch/btx604/$1/genomePolishing/illumina/raw_reads/R
    1.fastq.gz
15 R2=/data/scratch/btx604/$1/genomePolishing/illumina/raw_reads/R
    2.fastq.gz
16 alignment_step1_sam=/data/scratch/btx604/$1/genomePolishing/illu
    mina/step1/"$1"_alignment_step1.sam
17 #samtools view step1 variables
18 alignment_step1_bam=/data/scratch/btx604/$1/genomePolishing/illu
    mina/step1/"$1"_alignment_step1.bam
19 #samtools sort index step1 variables
20 alignment_step1_sorted=/data/scratch/btx604/$1/genomePolishing/i
    llumina/step1/"$1"_sorted_step1.bam
21 #pilon step1 variables
22 step1_pilon_prefix="$1"_pilon_step1
```



```
23 step1_pilon_fa=/data/scratch/btx604/$1/genomePolishing/illumina/
step1/"$1"_pilon_step1.fasta
24
25 #bwa index step2 variables
26 reference_genome_step2_prefix="$1"_pilon_step1.fasta
27 reference_genome_step2=/data/scratch/btx604/$1/genomePolishing/i
llumina/step2/$reference_genome_step2_prefix
28 #bwa mem step2 variables
29 alignment_step2_sam=/data/scratch/btx604/$1/genomePolishing/illu
mina/step2/"$1"_alignment_step2.sam
30 #samtools view step2 variables
31 alignment_step2_bam=/data/scratch/btx604/$1/genomePolishing/illu
mina/step2/"$1"_alignment_step2.bam
32 #samtools sort index step2 variables
33 alignment_step2_sorted=/data/scratch/btx604/$1/genomePolishing/i
llumina/step2/"$1"_sorted_step2.bam
34 #pilon step2 variables
35 step2_pilon_prefix="$1"_pilon_step2
36 step2_pilon_fa=/data/scratch/btx604/$1/genomePolishing/illumina/
step2/"$1"_pilon_step2.fasta
37
38 cd /data/scratch/btx604/$1/genomePolishing/illumina/step1
39
40 module load bwa
41 module load samtools/1.9
42 module load anaconda2
43 source activate genomePolishing
44 source activate pilon
45
46 echo 'working on '$1
47
```

```
48 echo 'BWA INDEX STEP1_____
   _____'
49 if [ -e /data/scratch/btx604/$1/genomePolishing/illumina/step1/
   *.ann ]
50 then
51     echo "/data/scratch/btx604/$1/genomePolishing/illumina/step1/
   *.ann found."
52 else
53     bwa index -p $reference_genome_step1_prefix -a bwtsv $referenc
   e_genome_step1
54 fi
55
56 echo 'BWA MEM STEP1_____
   _____'
57 if [ -e "$alignment_step1_sam" ]
58 then
59     echo "$alignment_step1_sam found."
60 else
61     bwa mem -t 24 -M $reference_genome_step1 $R1 $R2 > $alignment_
   step1_sam
62 fi
63
64 echo 'SAMTOOLS VIEW STEP1_____
   _____'
65 if [ -e "$alignment_step1_bam" ]
66 then
67     echo "$alignment_step1_bam found."
68 else
69     samtools view -S -b -h $alignment_step1_sam -o $alignment_step
   1_bam
70 fi
71
```

```
72 echo 'SAMTOOLS SORT INDEX STEP1_____
   _____'
73 if [ -e "$alignment_step1_sorted" ]
74 then
75     echo "$alignment_step1_sorted found."
76 else
77     samtools sort $alignment_step1_bam -o $alignment_step1_sorted
78     samtools index $alignment_step1_sorted
79 fi
80
81 echo 'PILON STEP1_____
   _____'
82 if [ -e "$step1_pilon_fa" ]
83 then
84     echo "$step1_pilon_fa found."
85 else
86     java -Xmx700G -jar /data/home/btx604/.conda/envs/pilon/share/p
   ilon-1.23-2/pilon-1.23.jar --genome $reference_genome_step1 --fr
   ags $alignment_step1_sorted --diploid --outdir /data/scratch/btx
   604/$1/genomePolishing/illumina/step1 --output $step1_pilon_pref
   ix --threads 4
87 fi
88
89 cp $step1_pilon_fa /data/scratch/btx604/$1/genomePolishing/illum
   ina/step2/
90 cd /data/scratch/btx604/$1/genomePolishing/illumina/step2
91
92 echo 'BWA INDEX STEP2_____
   _____'
93 if [ -e /data/scratch/btx604/$1/genomePolishing/illumina/step2/
   *.ann ]
94 then
```

```
95     echo "/data/scratch/btx604/$1/genomePolishing/illumina/step2/
*.ann found."
96 else
97     bwa index -p $reference_genome_step2_prefix -a bwtsw $referenc
e_genome_step2
98 fi
99
100 echo 'BWA MEM STEP2_____
_____ '
101 if [ -e "$alignment_step2_sam" ]
102 then
103     echo "$alignment_step2_sam found."
104 else
105     bwa mem -t 24 -M $reference_genome_step2 $R1 $R2 > $alignment_
step2_sam
106 fi
107
108 echo 'SAMTOOLS VIEW STEP2_____
_____ '
109 if [ -e "$alignment_step2_bam" ]
110 then
111     echo "$alignment_step2_bam found."
112 else
113     samtools view -S -b -h $alignment_step2_sam -o $alignment_step
2_bam
114 fi
115
116 echo 'SAMTOOLS SORT INDEX STEP2_____
_____ '
117 if [ -e "$alignment_step2_sorted" ]
118 then
119     echo "$alignment_step2_sorted found."
```

```
120 else
121     samtools sort $alignment_step2_bam -o $alignment_step2_sorted
122     samtools index $alignment_step2_sorted
123 fi
124
125 echo 'PILON STEP2-----
-----'
126 if [ -e "$step2_pilon_fa" ]
127 then
128     echo "$step2_pilon_fa found."
129 else
130     java -Xmx700G -jar /data/home/btx604/.conda/envs/pilon/share/p
ilon-1.23-2/pilon-1.23.jar --genome $reference_genome_step2 --fr
ags $alignment_step2_sorted --diploid --outdir /data/scratch/btx
604/$1/genomePolishing/illumina/step2 --output $step2_pilon_pref
ix --threads 4
131 fi
```

Step 2.5: Quality check

Busco and Quast will be used again here. The first software is used to obtain info about the polished version of our genome, while the second is used to compare the assembly at each stage of the polishing providing a detailed review about the changes occurred during this process.

Code to run Busco on the polished assembly:

```
1 #!/bin/bash
2 #$ -wd /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia/busco/g
enomePolishing
3 #$ -o /data/SBCS-MartinDuranLab/03-Giacomo/logs/busco
4 #$ -pe smp 4
5 #$ -l h_vmem=2G
6 #$ -l h_rt=10:0:0
```

```
7  # $ -j y
8
9  cd /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia/busco/genom
   ePolishing
10
11  module load busco/3.0
12  module load augustus
13
14  export AUGUSTUS_CONFIG_PATH=/data/SBCS-MartinDuranLab/02-Chema/s
   rc/Augustus/config
15
16  BUSCO.py -i /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia/ge
   nomePolishing/illumina/oasisia_pilon_step2.fasta -m genome -o oa
   sisia_polished_busco -c 4 -l /data/home/btx604/datasets/metazoa_
   odb9
```

Code to run Quast:

```
1  #!/bin/bash
2  # $ -wd /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia/quast/g
   enomePolishing
3  # $ -o /data/SBCS-MartinDuranLab/03-Giacomo/logs/quast
4  # $ -pe smp 8
5  # $ -l h_vmem=10G
6  # $ -l h_rt=10:0:0
7  # $ -cwd
8  # $ -j y
9
10  module load anaconda2
11  source activate quast
12
13  quast \
14  /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia/canu/Oasisia.
```

```
contigs.fasta \  
15 /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia/genomePolishi  
ng/PB/oasisia_step1_consensus.fasta \  
16 /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia/genomePolishi  
ng/PB/oasisia_step2_consensus.fasta \  
17 /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia/genomePolishi  
ng/illumina/oasisia_pilon_step1.fasta \  
18 /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia/genomePolishi  
ng/illumina/oasisia_pilon_step2.fasta \  
19 -o /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia/quast/geno  
mePolishing/ \  
20 --eukaryote \  

```

Step 3: Decontamination

BlobTools (1 not 2) is a tool able to determine the taxonomy of each contig in our assembled genome and flag those portions accordingly. This is very useful to identify which sequences truly derive from the target genome and which are derived from associated microbiome or contaminant organisms.

Step 3.1: Preparing the data

In order to run Blobtools we need to obtain 3 other files: 1. An alignment file, mapping Illumina cutadapted reads against the polished genome, 2. a BlastN dataset, 3. a BlastX dataset.

1-Alignment file

To do first (create a new working directory, copy and rename the polished assembly):

```
1 mkdir /data/scratch/btx604/oasisia/genomePolishing/illumina/blob  
tools  
2 cp oasisia_pilon_step2.fasta /data/scratch/btx604/oasisia/genome  
Polishing/illumina/blobtools  
3 mv oasisia_pilon_step2.fasta oasisia_step2_consensus.fasta  

```

Cutadapted R1 and R2 are in /data/scratch/btx604/\$1/genomePolishing/illumina

/raw_reads/

Now we can run the following steps all together in one script:

mapping_polished_oasisia_illumina_reads_v1.sh

```
1  #!/bin/bash
2  #$ -wd /data/scratch/btx604
3  #$ -j y
4  #$ -o /data/scratch/btx604
5  #$ -pe smp 24
6  #$ -l h_vmem=10G
7  #$ -l h_rt=24:0:0
8  #$ -l highmem
9
10 #bwa index step1 variables
11 reference_genome_step1_prefix="$1"_step2_consensus.fasta
12 reference_genome_step1=/data/scratch/btx604/$1/genomePolishing/i
    llumina/blobtools/$reference_genome_step1_prefix
13 #bwa mem step1 variables
14 R1=/data/scratch/btx604/$1/genomePolishing/illumina/raw_reads/R
    1.fastq.gz
15 R2=/data/scratch/btx604/$1/genomePolishing/illumina/raw_reads/R
    2.fastq.gz
16 alignment_step1_sam=/data/scratch/btx604/$1/genomePolishing/illu
    mina/blobtools/"$1"_alignment_step1.sam
17 #samtools view step1 variables
18 alignment_step1_bam=/data/scratch/btx604/$1/genomePolishing/illu
    mina/blobtools/"$1"_alignment_step1.bam
19 #samtools sort index step1 variables
20 alignment_step1_sorted=/data/scratch/btx604/$1/genomePolishing/i
    llumina/blobtools/"$1"_sorted_step1.bam
21
22 module load bwa
```



```
23 module load samtools/1.9
24
25 cd /data/scratch/btx604/oasisia/genomePolishing/illumina/blobtools
26
27 echo 'BWA INDEX STEP1_____
28 _____'
29 if [ -e /data/scratch/btx604/$1/genomePolishing/illumina/blobtools/*.ann ]
30 then
31     echo "/data/scratch/btx604/$1/genomePolishing/illumina/blobtools/*.ann found."
32 else
33     bwa index -p $reference_genome_step1_prefix -a bwtsv $reference_genome_step1
34 fi
35
36 echo 'BWA MEM STEP1_____
37 _____'
38 if [ -e "$alignment_step1_sam" ]
39 then
40     echo "$alignment_step1_sam found."
41 else
42     bwa mem -t 24 -M $reference_genome_step1 $R1 $R2 > $alignment_step1_sam
43 fi
44
45 echo 'SAMTOOLS VIEW STEP1_____
46 _____'
47 if [ -e "$alignment_step1_bam" ]
48 then
49     echo "$alignment_step1_bam found."
```

```
47 else
48     samtools view -S -b -h $alignment_step1_sam -o $alignment_step
1_bam
49 fi
50
51 echo 'SAMTOOLS SORT INDEX STEP1-----
-----'
52 if [ -e "$alignment_step1_sorted" ]
53 then
54     echo "$alignment_step1_sorted found."
55 else
56     samtools sort $alignment_step1_bam -o $alignment_step1_sorted
57     samtools index $alignment_step1_sorted
58 fi
```

2-BlastN dataset

Install the NCBI nucleotide database (nt). First we need to create a conda env in order to avoid problems with perl dependencies

```
1 conda create -n blast
2 source activate blast
3 conda install -c bioconda blast
```

Now we can run the code to download the last version of the NCBI nt database (V5-10/02/20)

update_blastdb_5_v1.sh

```
1 #!/bin/bash
2 #$ -wd /data/SBCS-MartinDuranLab/03-Giacomo/db/blast_nt_v5
3 #$ -j y
4 #$ -o /data/SBCS-MartinDuranLab/03-Giacomo/db/blast_nt_v5
5 #$ -pe smp 2
6 #$ -l h_vmem=10G
7 #$ -l h_rt=24:0:0
```

```
8
9 module load anaconda2
10 source activate blast
11
12 cd /data/SBCS-MartinDuranLab/03-Giacomo/db/blast_nt_v5
13 update_blastdb.pl --blastdb_version 5 nt --decompress
```

- the directory “blast_nt_v5” should be created previously

Finally we can run the code to obtain the blast nt dataset using the polished genome as query:

blastn_oasisia_v2.sh

```
1 #!/bin/bash
2 #$ -wd /data/SBCS-MartinDuranLab/03-Giacomo/db/blast_nt_v5/
3 #$ -o /data/scratch/btx604/oasisia/blobtools/nt_v5
4 #$ -j y
5 #$ -pe smp 16
6 #$ -l h_vmem=20G
7 #$ -l h_rt=120:0:0
8 #$ -l highmem
9
10 module load anaconda2
11 source activate blobtools2
12
13 cd /data/SBCS-MartinDuranLab/03-Giacomo/db/blast_nt_v5/
14
15 export BLASTDB=/data/SBCS-MartinDuranLab/03-Giacomo/db/blast_nt_
16 v5/
17
18 blastn -db nt \
19 -query /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia/genome
20 Polishing/illumina/oasisia_pilon_step2.fasta \
```

```
19 -outfmt "6 qseqid staxids bitscore std" \  
20 -max_target_seqs 10 \  
21 -max_hsps 1 \  
22 -evaluate 1e-25 \  
23 -num_threads 16 \  
24 -out /data/scratch/btx604/oasisia/blobtools/nt_v5/blast.out
```

3-BlastX dataset

Create conda environment

```
1 module load anaconda2  
2 conda config --add channels bioconda  
3 conda config --add channels conda-forge  
4 conda config --add channels r  
5 cd /data/SBCS-MartinDuranLab/03-Giacomo/src/blobtools2  
6 conda create -n blobtools2 -y python=3.6 docopt pyyaml ujson pysam  
  tqdm nodejs seqtk  
7 source activate blobtools2  
8 conda install -c bioconda -y pysam seqtk  
9 conda install -c bioconda -y blast=2.9 busco diamond minimap2
```

- I think this can be avoided, what I need here is just diamond that can also be loaded with `module load diamond`

Download and format the UniProt reference_proteomes as a diamond database:

```
1 #!/bin/bash  
2 #$ -wd /data/SBCS-MartinDuranLab/03-Giacomo/src/blobtools2  
3 #$ -o /data/SBCS-MartinDuranLab/03-Giacomo/src/blobtools2  
4 #$ -j y  
5 #$ -pe smp 16  
6 #$ -l h_vmem=3G  
7 #$ -l h_rt=24:0:0  
8
```

```
9 module load anaconda2
10 source activate blobtools2
11
12 cd /data/SBCS-MartinDuranLab/03-Giacomo/src/blobtools2
13
14 mkdir -p uniprot
15
16 # fetch latest reference proteome database
17 wget -q -O uniprot/reference_proteomes.tar.gz \
18     ftp.ebi.ac.uk/pub/databases/uniprot/current_release/knowledgebase/reference_proteomes/$(curl \
19     -vs ftp.ebi.ac.uk/pub/databases/uniprot/current_release/knowledgebase/reference_proteomes/ 2>&1 | \
20     awk '/tar.gz/ {print $9}')
21
22 cd uniprot
23
24 tar xf reference_proteomes.tar.gz
25
26 # extract and concatenate protein FASTA files
27 touch reference_proteomes.fasta.gz
28 find . -mindepth 2 | grep "fasta.gz" | grep -v 'DNA' | grep -v 'additional' | xargs cat >> reference_proteomes.fasta.gz
29
30 # extract and concatenate taxid mapping files
31 echo "accession\taccession.version\ttaxid\ttgi" > reference_proteomes.taxid_map
32 zcat */*.idmapping.gz | grep "NCBI_TaxID" | awk '{print $1 "\t" $1 "\t" $3 "\t" 0}' >> reference_proteomes.taxid_map
33
34 # make diamond blast db with taxonomic information included
```

```
35 diamond makedb -p 16 --in reference_proteomes.fasta.gz --taxonmap  
p reference_proteomes.taxid_map --taxonnodes ../taxdump/nodes.dm  
p -d reference_proteomes.dmnd
```

Code to obtain the blastX dataset using the polished genome as query:

blastx_oasisia_v1.sh

```
1  #!/bin/bash
2  #$ -wd /data/scratch/btx604/oasisia/blobtools
3  #$ -o /data/scratch/btx604/oasisia/blobtools
4  #$ -j y
5  #$ -pe smp 16
6  #$ -l h_vmem=20G
7  #$ -l h_rt=120:0:0
8  #$ -l highmem
9
10 module load anaconda2
11 source activate blobtools2
12
13 cd /data/scratch/btx604/oasisia/blobtools
14
15 diamond blastx \
16 --query /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia/genom  
ePolishing/illumina/oasisia_pilon_step2.fasta \
17 --db /data/SBCS-MartinDuranLab/03-Giacomo/src/blobtools2/unipro  
t/reference_proteomes.dmnd \
18 --outfmt 6 qseqid staxids bitscore qseqid sseqid pident length  
mismatch gapopen qstart qend sstart send evalue bitscore \
19 --sensitive \
20 --max-target-seqs 1 \
21 --evaluate 1e-25 \
22 --threads 16 \
23 > /data/scratch/btx604/oasisia/blobtools/diamond.out
```

Step 3.2: Blobtools

Blobtools is a command line tool designed to aid genome assembly QC and contaminant/cobiont detection and filtering.

First we need to create a BlobDB data structure based on input files that we have obtained in 3.1

blobtools1_create_oasisia_v2.sh

```
1  #!/bin/bash
2  #$ -wd /data/scratch/btx604/oasisia/blobtools/nt_v5
3  #$ -o /data/scratch/btx604/oasisia/blobtools/nt_v5
4  #$ -j y
5  #$ -pe smp 4
6  #$ -l h_vmem=5G
7  #$ -l h_rt=24:0:0
8
9  module load anaconda2
10 source activate blobtools
11
12 cd /data/scratch/btx604/oasisia/blobtools/nt_v5
13
14 blobtools create \
15   -i /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia/genomePolishing/illumina/oasisia_pilon_step2.fasta \
16   -b /data/scratch/btx604/oasisia/genomePolishing/illumina/blobtools/oasisia_sorted_step1.bam \
17   -t /data/scratch/btx604/oasisia/blobtools/nt_v5/blast.out \
18   -t /data/scratch/btx604/oasisia/blobtools/diamond.out \
19   -o /data/scratch/btx604/oasisia/blobtools/nt_v5/oasisia_blobplot
```

Then we need to generate a tabular output containing the taxonomic rank

informations. This table contains information about the sequences in the order they appear in the assembly file and it will be used in 3.3. The following command is also generating two plots that can help to better understand your assembly.

`blobtools1_view_plot_oasisia_v2.sh`

```
1  #!/bin/bash
2  #$ -wd /data/scratch/btx604/oasisia/blobtools/nt_v5
3  #$ -o /data/scratch/btx604/oasisia/blobtools/nt_v5
4  #$ -j y
5  #$ -pe smp 4
6  #$ -l h_vmem=5G
7  #$ -l h_rt=24:0:0
8
9  module load anaconda2
10 source activate blobtools
11
12 cd /data/scratch/btx604/oasisia/blobtools/nt_v5
13
14 blobtools view \
15   -i /data/scratch/btx604/oasisia/blobtools/nt_v5/oasisia_blobplo
   t.blobDB.json \
16   -o /data/scratch/btx604/oasisia/blobtools/nt_v5/
17
18 blobtools plot \
19   -i /data/scratch/btx604/oasisia/blobtools/nt_v5/oasisia_blobplo
   t.blobDB.json \
20   -o /data/scratch/btx604/oasisia/blobtools/nt_v5/
```

Step 3.3: Filter the assembly

Now we need to extract fasta sequences from the polished assembly based on taxonomy. In order to accomplish this, 3 steps will be required:

1-Divide blobtools table

This will generate five tables containing only the sequences from a particular group of organisms

PAY ATTENTION! the last line of oasisia_blobplot.blobDB.bestsum.table.txt is not well formatted, in fact `wc -l` is giving one line more, the best is to delete the last empty line. Furthermore, the first lines of host table containing comments and info can be removed.

To get all the phyla present in my assembly do this:

```
awk -F '\t' '{print $6}' oasisia_blobplot.blobDB.bestsum.table.txt | sort | uniq
```

And then grep phyla by phyla with the following:

`new_filter_blobtools_v1.sh`

```
1  #!/bin/bash
2  #$ -wd /data/scratch/btx604/
3  #$ -o /data/scratch/btx604/
4  #$ -j y
5  #$ -pe smp 2
6  #$ -l h_vmem=5G
7  #$ -l h_rt=24:0:0
8
9  #variables:
10 blobplot_original=/data/SBCS-MartinDuranLab/03-Giacomo/data/$1/blobtools/"$1"_blobplot.blobDB.bestsum.table.txt
11 blobplot=/data/SBCS-MartinDuranLab/03-Giacomo/data/$1/blobtools/"$1"_blobplot_tail_13.blobDB.bestsum.table.txt
12 host_table=/data/SBCS-MartinDuranLab/03-Giacomo/data/$1/blobtools/new_filter/new_filter_"$1"_blobplot_host.table.txt
13 proteobacteria_table=/data/SBCS-MartinDuranLab/03-Giacomo/data/$1/blobtools/new_filter/new_filter_"$1"_blobplot_proteobacteria.table.txt
14 actinobacteria_table=/data/SBCS-MartinDuranLab/03-Giacomo/data/$1/blobtools/new_filter/new_filter_"$1"_blobplot_actinobacteria.table.txt
```

```
a.table.txt
15 other_bacteria_table=/data/SBCS-MartinDuranLab/03-Giacomo/data
   /$1/blobtools/new_filter/new_filter_"$1"_blobplot_other_bacteri
   a.table.txt
16 virus_table=/data/SBCS-MartinDuranLab/03-Giacomo/data/$1/blobtoo
   ls/new_filter/new_filter_"$1"_blobplot_virus.table.txt
17
18 echo 'working on '$1
19
20 cd /data/SBCS-MartinDuranLab/03-Giacomo/data/$1/blobtools
21 mkdir new_filter
22
23 tail -n +13 $blobplot_original > $blobplot
24
25 grep -wi Arthropoda $blobplot >> $host_table
26 grep -wi Annelida $blobplot >> $host_table
27 grep -wi Ascomycota $blobplot >> $host_table
28 grep -wi Brachiopoda $blobplot >> $host_table
29 grep -wi Chordata $blobplot >> $host_table
30 grep -wi Cnidaria $blobplot >> $host_table
31 grep -wi Echinodermata $blobplot >> $host_table
32 grep -wi Hemichordata $blobplot >> $host_table
33 grep -wi Mollusca $blobplot >> $host_table
34 grep -wi Nematoda $blobplot >> $host_table
35 grep -wi Nemertea $blobplot >> $host_table
36 grep -wi Platyhelminthes $blobplot >> $host_table
37 grep -wi Porifera $blobplot >> $host_table
38 grep -wi Priapulida $blobplot >> $host_table
39 grep -wi Streptophyta $blobplot >> $host_table
40 grep -wi Chytridiomycota $blobplot >> $host_table
41 grep -wi Ciliophora $blobplot >> $host_table
```

```
42 grep -wi Euglenozoa $blobplot >> $host_table
43 grep -wi Apicomplexa $blobplot >> $host_table
44 grep -wi Basidiomycota $blobplot >> $host_table
45 grep -wi Chlorophyta $blobplot >> $host_table
46 grep -wi Cryptophyta $blobplot >> $host_table
47 grep -wi Eukaryota-undef $blobplot >> $host_table
48 grep -wi Evosea $blobplot >> $host_table
49 grep -wi Haptista $blobplot >> $host_table
50 grep -wi Mucoromycota $blobplot >> $host_table
51 grep -wi Placozoa $blobplot >> $host_table
52 grep -wi Rotifera $blobplot >> $host_table
53
54 grep -wi unresolved $blobplot >> $host_table
55 grep -wi phylum.t.6%s $blobplot >> $host_table
56 grep -wi no-hit $blobplot >> $host_table
57
58 grep -wi Proteobacteria $blobplot >> $proteobacteria_table
59
60 grep -wi Actinobacteria $blobplot >> $actinobacteria_table
61
62
63 grep -wi Bacteria-undef $blobplot >> $other_bacteria_table
64 grep -wi Bacteroidetes $blobplot >> $other_bacteria_table
65 grep -wi Calditrichaeota $blobplot >> $other_bacteria_table
66 grep -wi Chlorobi $blobplot >> $other_bacteria_table
67 grep -wi Cyanobacteria $blobplot >> $other_bacteria_table
68 grep -wi Firmicutes $blobplot >> $other_bacteria_table
69 grep -wi Fusobacteria $blobplot >> $other_bacteria_table
70 grep -wi Ignavibacteriae $blobplot >> $other_bacteria_table
71 grep -wi Kiritimatiellaeota $blobplot >> $other_bacteria_table
```

```

72 grep -wi Lentisphaerae $blobplot >> $other_bacteria_table
73 grep -wi Nitrospirae $blobplot >> $other_bacteria_table
74 grep -wi Planctomycetes $blobplot >> $other_bacteria_table
75 grep -wi Verrucomicrobia $blobplot >> $other_bacteria_table
76 grep -wi Chlamydiae $blobplot >> $other_bacteria_table
77 grep -wi Nitrospinae $blobplot >> $other_bacteria_table
78 grep -wi Spirochaetes $blobplot >> $other_bacteria_table
79 grep -wi Thermodesulfobacteria $blobplot >> $other_bacteria_table
80
81 grep -wi Viruses-undef $blobplot >> $virus_table

```

```

1 qsub new_filter_blobtools_v1.sh oasisia
2 qsub new_filter_blobtools_v1.sh riftia
3 qsub new_filter_blobtools_v1.sh osedax

```

Quick check:

```

1 wc -l $blobplot #this is the total
2 wc -l $host_table #sum
3 wc -l $proteobacteria_table #sum
4 wc -l $actinobacteria_table #sum
5 wc -l $other_bacteria_table #sum
6 wc -l $virus_table #sum

```

2-Obtain headers

The next script will create a list of contig names (one per line) that will then be used to extract fasta sequences. It requires 3 inputs from the command line: \$1 is the name to one of the tables we obtained before, \$2 is the taxa of the organisms contained in the table (e.g. Actinobacteria, host...) and \$3 is the path of the directory containing the input file and where you want to place the output.

create_headers_list_v1.sh

```

1 #!/bin/bash

```

```
2  #$ -wd /data/scratch/btx604/
3  #$ -o /data/scratch/btx604/
4  #$ -j y
5  #$ -pe smp 1
6  #$ -l h_vmem=1G
7  #$ -l h_rt=1:0:0
8
9  #This script is working on subsets depending on taxa of blobtool
  s tables (e.g. oasisia_blobplot_host.table.txt)
10
11 #variables
12 input=$1
13 taxa=$2
14 working_directory=$3
15 output=headers_"$2".txt
16
17 cd $working_directory
18 touch $output
19
20 while IFS= read -r line
21 do
22   echo $line | awk '{print $1;}' >> $output
23 done < $input
```

```
qsub create_headers_list_v1.sh oasisia_blobplot_Proteobacteria.t
able.txt proteobacteria /data/scratch/btx604/oasisia/blobtools/n
t_v5/
```

- This example is showing how to launch the script from the command line

Now you have four/five files containing the headers divided by taxa

3-Obtain sequences

Let's create a directory "fasta" where we want to collect our fasta files, one for each taxa of interest.

```
mkdir fasta
```

Samtools as a function (faidx) that can extract subsequence from indexed reference sequence when they are matching with a specific pattern, or as in this case with the headers.

samtools_faidx_oasisia_blob_fastas_v1.sh

```
1  #!/bin/bash
2  #$ -wd /data/scratch/btx604/oasisia/blobtools/nt_v5/fasta/
3  #$ -o /data/scratch/btx604/oasisia/blobtools/nt_v5/fasta/
4  #$ -j y
5  #$ -pe smp 1
6  #$ -l h_vmem=5G
7  #$ -l h_rt=3:0:0
8
9  module load samtools/1.9
10
11  samtools faidx /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia
    /genomePolishing/illumina/oasisia_pilon_step2.fasta -r /data/scr
    atch/btx604/oasisia/blobtools/nt_v5/headers_proteobacteria.txt >
    /data/scratch/btx604/oasisia/blobtools/nt_v5/fasta/oasisia_prote
    obacteria.fasta
12
13  samtools faidx /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia
    /genomePolishing/illumina/oasisia_pilon_step2.fasta -r /data/scr
    atch/btx604/oasisia/blobtools/nt_v5/headers_actinobacteria.txt >
    /data/scratch/btx604/oasisia/blobtools/nt_v5/fasta/oasisia_actin
    obacteria.fasta
14
15  samtools faidx /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia
    /genomePolishing/illumina/oasisia_pilon_step2.fasta -r /data/scr
    atch/btx604/oasisia/blobtools/nt_v5/headers_bacteria.txt > /data
    /scratch/btx604/oasisia/blobtools/nt_v5/fasta/oasisia_bacteria.f
```

asta

```
samtools faidx /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia/genomePolishing/illumina/oasisia_pilon_step2.fasta -r /data/scratch/btx604/oasisia/blobtools/nt_v5/headers_host.txt > /data/scratch/btx604/oasisia/blobtools/nt_v5/fasta/oasisia_host.fasta
```

- it works very quickly

Check if the files look alright:

```
1 wc -l /data/scratch/btx604/oasisia/blobtools/nt_v5/headers_proteobacteria.txt #5
2 grep -c "^>" /data/scratch/btx604/oasisia/blobtools/nt_v5/fasta/oasisia_proteobacteria.fasta #5
3
4 wc -l /data/scratch/btx604/oasisia/blobtools/nt_v5/headers_bacteria.txt
5 grep -c "^>" /data/scratch/btx604/oasisia/blobtools/nt_v5/fasta/oasisia_bacteria.fasta
6
7 wc -l /data/scratch/btx604/oasisia/blobtools/nt_v5/headers_actinobacteria.txt
8 grep -c "^>" /data/scratch/btx604/oasisia/blobtools/nt_v5/fasta/oasisia_actinobacteria.fasta
9
10 wc -l /data/scratch/btx604/oasisia/blobtools/nt_v5/headers_host.txt
11 grep -c "^>" /data/scratch/btx604/oasisia/blobtools/nt_v5/fasta/oasisia_host.fasta
```

- they should give the same value

Step 4: Haploidization

This is the final step of the assembly. The fasta obtained during the previous step

(oasisia_host.fasta) has high level of duplication level, in fact that version of the genome can be considered the diploid version. However, we are interested in the haploid version that will be easier to annotate and will give better and more precise results. **Purge_Dups** is a software designed for this purpose. Minimap2, an aligner tool, Busco and Quast will also be used during this phase.

Step 4.1: Installation

Purge_Dups installation:

```
1 03g
2 cd src/
3 git clone https://github.com/dfguan/purge_dups.git
4 cd purge_dups/src && make
5 git clone https://github.com/dfguan/runner.git
6 module load python
7 cd runner && python setup.py install --user
```

Minimap2 installation:

```
1 module load anaconda2
2 conda create -n Minimap2
3 source activate Minimap2
4 conda install -c bioconda minimap2
```

Step 4.2: Purge_Dups and quality checks

All the programs that we need (Busco and Quast as well) are included in the following unique and universal script. It just needs a directory “purge_dups” containing the output of blobtools (e.g. oasisia_host.fasta) and it should be submitted like this:

```
qsub purge_dups_v2.sh $1 $2
```

- \$1 is the species name (e.g. oasisia)
- \$2 is the absolute path to the PacBio raw reads obtained transforming subreads.bam into fq.gz with BAM2fastx tools 19/08/19 (e.g. /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia/00-pacbio/oasisia_pb_raw.fastq.gz)

Code to run the pipeline:

```
1  #!/bin/bash
2  #$ -wd /data/scratch/btx604/
3  #$ -j y
4  #$ -o /data/scratch/btx604/
5  #$ -pe smp 12
6  #$ -l h_vmem=30G
7  #$ -l h_rt=240:0:0
8  #$ -l highmem
9
10 #variables
11 ref_genome=/data/scratch/btx604/$1/purge_dups/"$1"_host.fasta
12 pb_fasta=$2
13 ref_genome_split=/data/scratch/btx604/$1/purge_dups/"$1"_host.sp
lit.fasta
14 self_aln_genome=/data/scratch/btx604/$1/purge_dups/"$1"_host.spl
it.self.paf.gz
15
16
17 cd /data/scratch/btx604/$1/purge_dups/
18
19 echo 'working on '$1
20
21 echo 'MINIMAP2 -----
-----'if [ -e ./purge_"$1".paf.gz ]
22 then
23     echo "purge_$1.paf.gz found"
24 else
25     module load anaconda2
26     source activate Minimap2
```

```
27
28   minimap2 -x map-pb $ref_genome $pb_fasta | gzip -c - > purge
   _$1.paf.gz
29
30   source deactivate
31   module unload anaconda2
32 fi
33
34 echo 'PBCSTAT -----
   -----'
35 if [ -e ./PB.base.cov ]
36 then
37   echo "PB.base.cov found"
38 else
39   module load python
40
41   /data/SBCS-MartinDuranLab/03-Giacomo/src/purge_dups/bin/pbcsta
   t *.paf.gz #(produces PB.base.cov and PB.stat files)
42
43   module unload python
44 fi
45
46 echo 'CALCUTS -----
   -----' if [ -e ./calcults.log ]
47 then
48   echo "calcults.log found"
49 else
50   module load python
51
52   /data/SBCS-MartinDuranLab/03-Giacomo/src/purge_dups/bin/calcut
   s PB.stat > cutoffs 2> calcults.log
```

```
53
54     module unload python
55 fi
56
57 echo 'SPLIT_FA -----
-----'
58 if [ -e "$ref_genome_split" ]
59 then
60     echo "$ref_genome_split found"
61 else
62     module load python
63
64     /data/SBCS-MartinDuranLab/03-Giacomo/src/purge_dups/bin/split_
fa $ref_genome > $ref_genome_split
65
66     module unload python
67 fi
68
69 echo 'MINIMAP2 -----
-----' if [ -e "$self_aln_genome" ]
70 then
71     echo "$self_aln_genome found"
72 else
73     module load anaconda2
74     source activate Minimap2
75
76     minimap2 -x asm5 -DP $ref_genome_split $ref_genome_split | gzi
p -c - > $self_aln_genome
77
78     source deactivate
79     module unload anaconda2
```

```
80  fi
81
82  echo 'PURGE_DUPS -----
-----'
83  if [ -e ./purge_dups.log ]
84  then
85      echo "purge_dups.log found"
86  else
87      module load python
88
89      /data/SBCS-MartinDuranLab/03-Giacomo/src/purge_dups/bin/purge_
dups -2 -T cutoffs -c PB.base.cov $self_aln_genome > dups.bed 2>
purge_dups.log
90
91      module unload python
92  fi
93
94  echo 'GET_SEQS -----
-----'
95  if [ -e ./hap.fa ]
96  then
97      echo "hap.fa found"
98  else
99      module load python
100
101      /data/SBCS-MartinDuranLab/03-Giacomo/src/purge_dups/bin/get_se
qs dups.bed $ref_genome > purged.fa 2> hap.fa
102
103      module unload python
104  fi
105
```

```
106 echo 'BUSCO -----
-----'
107 if [ -e ./busco/ ]
108 then
109     echo "hap.fa found"
110 else
111     module load busco/3.0
112     module load augustus
113     export AUGUSTUS_CONFIG_PATH=/data/SBCS-MartinDuranLab/02-Chema
/src/Augustus/config
114
115     mkdir busco
116     cd ./busco/
117
118     BUSCO.py -i ../purged.fa -m genome -o busco -f -c 4 -l /data/h
ome/btx604/datasets/metazoa_odb9
119
120     cd ..
121
122     module unload busco/3.0
123     module unload augustus
124 fi
125
126 echo 'QUAST -----
-----'
127 if [ -e ./quast/ ]
128 then
129     echo "quast found"
130 else
131     module load anaconda2
132     source activate quast
```

```
133
134 mkdir quast
135 quast \
136   ./purged.fa \
137   -o ./quast \
138   --eukaryote \
139
140 source deactivate
141 module unload anaconda2
142 fi
```

- purged.fa is the haploid version of the genome