# Guest Lecture

## Bodo Linz
## 02/18/20

# Bacterial Genomics:
# From sequencing reads
# to multiple genome alignment

# Guest Lecture
## Bodo Linz
## 02/18/20

## Today's lecture

- Download complete genomes from NCBI
- Split complete genome into overlapping "reads"
- Download a Short Read Archive (SRA) from NCBI
- Join paired reads from the archive
- Align joined reads/split "reads" against a reference genome
- Call SNPs, generate consensus sequence
- Generate multiple genome alignment
- Make pairwise genome comparisons using `blastn` and `MSPcrunch`, visualize in Artemis Comparison Tool
- Extract and assemble a gene sequence from an SRA

# Download a complete genome from NCBI

go to https://www.ncbi.nlm.nih.gov/genome/

type the species: *Bordetella holmesii*
Select: Genome Assembly and Annotation report

# Download a complete genome from NCBI



type the isolate: "ATCC 51541"

click on chromosome in replicons

tick "Show sequence", click "Update view"

# Download a complete genome from NCBI

```
>NZ_CP007495.1 Bordetella holmesii 44057, complete genome
CATGCCACCGAACTTCGCCTTCCAGTTGTAGTACGTTGCCTCGGAGATTCCGTGCTTGCGGCACAACTCT
GCGGGCTTGGCACCTGCATCGGCTTCCTTGAGCACGCCGATGATTTGCTCTTCCGTAAATCGTTTCTTCA
TTGCATTCCTTTGGGAACGGACTCTACATCGATTTCGTACTAATCACGGGGAGCAGGTCATACCGAGCGC
GTAACCCAGGAAACAAGGGCCCGGCATGCGTCTTCGCCCGGCATGACTCCCGCTCGATCACGCGTCAGCA
GAAAACACACCGCTTGCGCAACGCTCTAGAACGGCTAGCTGCGCAATGCTGCCCACAGTGGCGCCAGGCA
GCGTTCCGGCATCGTCGAGTTGCACGAGCAGCGGCGCGTAGGCTTGCGCCCCAACCTTCTGTGCAGAAAG
CCTCGCCTGCGTGGCGACGCTAGGGTCAG
CGACAGGGACTGGCAAAGAGGCGTTGGCA
GGAGATCGCCTCACGGGTGTAAATGTTTT
TGAGGAATAAAAGCCTGACGATGACCTAC
TCACTGTCCTGTTCGGGATGGGAAGGAGT
GGTTGACGGCGTTGTCAACGGCTTGAATTT
CGGCACTGGCGCGAACTGGCCAGCAATCG
```

# if closed genome in multiple line fasta format make single line fasta sequence

```
cat F029.fasta | awk
'BEGIN{RS=">";FS="\n"}NR>1{seq=seq"";for(i=2;i<=NF;i++)seq=seq""$i;
print ">"$1"\n"seq}' > F029g.fasta
```

# RS – Record Separator: end of record marker, default new line
        → new entry starts here
# FS – Field Separator: \n separates the fields
# for(i=2;i<=NF;i++)seq=seq""$i
# for all rows starting from 2, until the last row, seq is seq plus the next seq
# print ">"$1"\n"seq}' – print >, fasta header, \n, built sequence


# to split genomes into overlapping 400 bp reads run script

```
split_genome_to_reads.sh
```

```bash
#!/bin/bash
# split_genome_to_reads
# author Bodo Linz
# split a genome into 400 bp long overlapping reads, 20 bp steps

file="F061g.fasta"
NAMEGENOME=${file%%".fasta"}

echo ""
echo "load input file $NAMEGENOME.fasta"
echo ""
echo "split genome into 400 bp fragments"
echo "----------------------------------------------"
cat ${NAMEGENOME}.fasta | awk 'NR > 1' | fold -w400 > fake4a

# select all rows except the first     awk 'NR > 1'
# split into chunks of 400 nucleotides         fold -w400

echo ""
echo "R: add read number a."
echo "----------------------------------------------"
# Run R in '--slave' mode to incorporate in bash script
R --slave -f /home/bodo/bin/add_sequence_read_number.R
```

```
#!usr/bin/R
# add_sequence_read_number.R
# delete al objects
rm(list = ls())                    #    fake4a
                                   #    ATGTCTGATTGACCGTAGCATTGTAG
                                   #    TGAGTGCGTACCCGTACGTGACCATT
# load packages
 library(base)

# set the working dicrectory
setwd("~/bodo.2/bordetella/Bholmesii/align")

# load data in table format
data <- read.table("fake4a", header = FALSE, sep= "\t")

# row count in file data into file pos
pos <- seq.int(nrow(data))

data2 <- cbind(pos,data)   # combine files pos and data

write.table(data2, file = "fake4b", sep = "\t",
row.names = FALSE, col.names = FALSE)
```

```
# fake4b
  $1   $2
  1    "ATGTCTGATTGACCGTAGCATTGTAG"
  2    "TGAGTGCGTACCVGTACGTGACCATT"


cat fake4b | awk -v FS="\t" -v OFS="" '{print ">read"$1"a",
"\n", $2}' | tr -d '"' > ${NAMEGENOME}_reads.fa

# tr -d '"'
# change " to nothing
```

```
# Let's look at F061g_reads.fa
>read1a
ATGTCTGATTGACCGTAGCATTGTAG
>read2a
TGAGTGCGTACCCGTACGTGACCATT
```

# what we got    # what we want



forward reads
reverse reads

```
# F061g_reads.fa      so far
>read1a
ATGTCTGATTGACCGTAGCATTGTAG
>read2a
TGAGTGCGTACCCGTACGTGACCATT
```

```
# add 20 A's at beginning of sequence, then split into "reads"
cat ${NAMEGENOME}.fasta | awk 'NR > 1' | awk -v OFS="" '{print
"AAAAAAAAAAAAAAAAAAAA",$1}' | fold -w400 > fake4a

R --slave -f /home/bodo/bin/add_sequence_read_number.R

cat fake4b | awk -v FS="\t" -v OFS="" '{print ">read"$1"b",
"\n", $2}' | tr -d '"' >> ${NAMEGENOME}_reads.fa
```

# ">>" - append to this file

```
We got
>read1a
ATGTCTGATTGACCGTAGCATTGTAG
>read1b
AAAAAAAAAAATGTCTGATTGACCGT
```

now iterate with more A's

```
# add 40 A's at beginning of sequence, then split into "reads"
        read1c read2c read3c ……
# add 60 A's at beginning of sequence, then split into "reads"
        read1d read2d read3d ……
# add 80 A's at beginning of sequence, then split into "reads"
        read1e read2e read3e ……
# add 100 A's at beginning of sequence, then split into "reads"
        read1f read2f read3f ……
……
# add 380 A's at beginning of sequence, then split into "reads"
        read1t read2t read3t ……
```

We got overlapping forward reads, let's get the reverse reads

```
echo ""
echo "reverse genome"
echo "---------------------------------------------"
cat ${NAMEGENOME}.fasta | awk "NR > 1" | awk '{print $1}' >
temp.fas

cat temp.fas | tr "[ATGCatgcNn]" "[TACGtacgNn]" | rev | awk
'{print ">F061g-rev.fasta","\n",$1}' > ${NAMEGENOME}-rev.fasta
```

```
# Let's walk through:
echo ""
echo "reverse genome"
echo "--------------------------------------------"
cat ${NAMEGENOME}.fasta | awk "NR > 1" | awk '{print $1}' >
temp.fas
# awk "NR > 1" - select all rows except row 1
# awk '{print $1}' - print what we got

cat temp.fas | tr "[ATGCatgcNn]" "[TACGtacgNn]" | rev | awk
'{print ">F029g-rev.fasta","\n",$1}' > ${NAMEGENOME}-rev.fasta
# tr "[ATGCatgcNn]" "[TACGtacgNn]" - change A to T, T to A,
  C to G, G to C, a to t, t to a, etc.
# rev - reverse resulting sequence
# awk '{print ">F029g-rev.fasta","\n",$1}'
# write header ">F029g-rev.fasta", then new line ("\n")
  , then the new reverse sequence
```

```
# Repeat with reverse genome and add reads to previous file
echo "split reverse genome into 400 bp fragments"
echo "----------------------------------------------"
cat ${NAMEGENOME}-rev.fasta | awk 'NR > 1' | fold -w400 > fake4a


cat fake4b | awk -v FS="\t" -v OFS="" '{print ">read"$1"reva",
"\n", $2}' | tr -d '"' >> ${NAMEGENOME}_reads.fa
```

different
suffix

```
# add 20 A's at beginning of rev sequence, split into "reads"
cat ${NAMEGENOME}-rev.fasta | awk 'NR > 1' | awk -v OFS=""
'{print "AAAAAAAAAAAAAAAAAAAA",$1}' | fold -w400 > fake4a

R --slave -f /home/bodo/bin/add_sequence_read_number.R

cat fake4b | awk -v FS="\t" -v OFS="" '{print ">read"$1"revb",
"\n", $2}' | tr -d '"' >> ${NAMEGENOME}_reads.fa
```

iterate with more A's
We got the reads file from a complete genome.

# Download a short read archive (SRA) from NCBI

The only option: use the `sratoolkit` from NCBI

- to download sratoolkit, type:
**wget** [ftp://ftp-trace.ncbi.nlm.nih.gov/sra/sdk/current/sratoolkit.current-centos_linux64.tar.gz](ftp://ftp-trace.ncbi.nlm.nih.gov/sra/sdk/current/sratoolkit.current-centos_linux64.tar.gz)
# or wherever the program is currently located at the ncbi website

# still there!

# Download a short read archive (SRA) from NCBI

The only option: use the `sratoolkit` from NCBI

- to download sratoolkit, type:
```
wget ftp://ftp-trace.ncbi.nlm.nih.gov/sra/sdk/current/sratoolkit.current-centos_linux64.tar.gz
```
# or wherever the program is currently located at the ncbi website

- to unpack the toolkit, type:
```
tar -xzf sratoolkit.current-centos_linux64.tar.gz
```

- location of fastq-dump and other commands:

```
~/[user_name]/sra-toolkit/bin/fastq-dump
```

# Download a short read archive (SRA) from NCBI
# Where do you find the archive?



click on the BioSample, e.g. SAMN02189846

# Download a short read archive (SRA) from NCBI – from where?

# Download a short read archive (SRA) from NCBI – from where?

# Download a short read archive (SRA) from NCBI – from where?



This is the SRA for this genome: SRR935461

# Download a short read archive (SRA) from NCBI

`~/[user_name]/sra-toolkit/bin/fastq-dump`

- go to the /bin directory

- Since the documentation is pretty minimal, here is the command line to type:
`./fastq-dump --outdir ~/bodo.2/Bholmesii/fastq --skip-technical  --readids --dumpbase --split-files --clip SRR_ID`

# ./fastq-dump – start the command fastq-dump in the current directory "./"
# --outdir – specify the output directory, here `~/bodo.2/Bholmesii/fastq`
# --skip-technical – dump only biological reads, skip info such as:
`Application Read Forward -> Technical Read Forward <- Application Read Reverse - Technical Read Reverse.`

# --readids – append the real read-ID after spot ID 'accession.spot.readid'
# --dumpbase – formats sequence using base space (default other than SOLiD)
# --split-files – Save forward and reverse reads into separate files. Files will receive suffix corresponding to read number.
# --clip SRR_ID – change the SRR_ID to whatever the ID is, e.g.  SRR942665

e.g.

```
./fastq-dump --outdir ~/bodo.2/Bholmesii/fastq
--skip-technical  --readids --dumpbase --split-files --clip
SRR942665
```

Downloaded paired reads: SRR942665_1.fastq and SRR942665_2.fastq

Let's have a look at the FASTQ format, it's in 4 lines:

@SEQ_ID
NUCLEOTIDE_SEQUENCE
+ (sometimes with seqID again)
QUALITY_SCORES_FOR_ALL_NUCLEOTIDES

e.g.

@SRR942665.3.1 SOLEXA4:47:D1RLFACXX:6:1101:2945:2102 length=101
TTCTGTGGAAAGGTGAGGTCATCGACGTCGGCGTGCGCCTCGGCGCGCAGGCCCACTTTGTCCAGGC
AGTCCCAGGCCAGGGCGCGCGCATCGGCCAGGCC
+
CCCFFDFFHHFHHIGGIIAEEHHJHGIJJJJIG@AGGIHGIGEADDDDDBDDBDBBBBDDDDCDCCCBBBC
DDDDC@BDDBBDDBBBBBBB@B<@DBDABBD

quality value characters in left-to-right increasing order of quality (ASCII):

#$%&'()*+,-./0123456789:;<=>?@
ABCDEFGHIJKLMNOPQRSTUVWXYZ[\]^_`abcdefghijklmnopqrstuvwxyz{|}~

Join the paired reads:
SRR942665_1.fastq and SRR942665_2.fastq using `FLASH`

Magoc and Salzberg (2011). FLASH: fast length adjustment of short reads to improve genome assemblies. *Bioinformatics* **27**: 2957-2963.

- very accurate and fast tool to merge overlapping paired-end reads
- Merged read pairs result in unpaired longer reads
- Longer reads are preferred in genome assembly and analysis processes

```
flash <mates1.fastq> <mates2.fastq> [-m minOverlap] [-M
maxOverlap] [-x mismatchRatio]
```

```
flash SRR942665_1.fastq SRR942665_2.fastq -m 10 -M 200 -x 0.1
```
You get 5 files:
```
out.extendedFrags.fastq
out.notCombined_1.fastq
out.notCombined_2.fastq
out.hist
out.histogram
```

Joined paired reads are in: `out.extendedFrags.fastq`

rename:
**`mv out.extendedFrags.fastq Bhz5132_SRR942665_joined.fastq`**

                                           ↑         ↑

                                 `strain archiveID`

if wanted
rename: **`mv out.notCombined_1.fastq Bhz5132_SRR942665_nc1.fastq`**
rename: **`mv out.notCombined_2.fastq Bhz5132_SRR942665_nc2.fastq`**

# Congratulations,
# you got the joined reads from a Short Read Archive!

We got the joined reads from a Short Read Archive in fastq format.

Problem: We got the reads file from a complete genome in fasta format
`F061_reads.fa.`

```
# change genome reads.fa to genome reads.fastq
# run perl script fasta_to_fastq.pl
# we have multiple genomes and run the script in a loop
files=$(ls *_reads.fa)
for file in $files; do name={file%%".fa"}; perl
~/bin/fasta_to_fastq.pl ${name}.fa > ${name}.fastq | echo "done
with $name"; done

# files=$(ls *_reads.fa) - create a list that contains all files
ending at _reads.fa
# for file in $files; do name={file%%".fa"}; - for all files in
this list, use the file name without ".fa"
# perl ~/bin/fasta_to_fastq.pl ${name}.fa > ${name}.fastq
for all files run perl script fasta_to_fastq.pl with input file
${name}.fa, save as output file ${name}.fastq
```

```perl
#Copyright (c) 2010 LUQMAN HAKIM BIN ABDUL HADI
#!/usr/bin/perl
use strict;
my $file = $ARGV[0];
open FILE, $file;
my ($header, $sequence, $sequence_length, $sequence_quality);
while(<FILE>) {
        chomp $_;
        if ($_ =~ /^>(.+)/) {
                if($header ne "") {
                        print "\@".$header."\n";
                        print $sequence."\n";
                        print "+"."\n";
                        print $sequence_quality."\n";
                }
                $header = $1;
                    $sequence = "";
                    $sequence_length = "";
                    $sequence_quality = "";
        }
          else {
                    $sequence .= $_;
                    $sequence_length = length($_);
                    for(my $i=0; $i<$sequence_length; $i++) {$sequence_quality .= "I"}
          }
}
close FILE;
print "\@".$header."\n";
print $sequence."\n";
print "+"."\n";
print $sequence_quality."\n";
```

```
We will
# align multiple genomes and run all steps in a loop
# align *.fastq files from SRA or from split genomes
  to the reference genome using bowtie2
# change .sam file to .bam file using samtools view
# sort bam file using samtools sort
# call variants using bcftools mpileup / bcftools call
# remove low quality variants by setting a threshold
# remove indels, check for potential problems
# zip the variant file using bgzip
# index the variant files using bcftools index
# generate consensus sequence using bcftools consensus
# change header and make fasta sequence one line
# make multiple alignment using cat
```

**BOWTIE2**

\# download from
http://sourceforge.net/projects/bowtie-bio/files/bowtie2/2.2.3/bowtie2-2.2.3-linux-x86_64.zip

\# unzip
```
unzip bowtie2-2.2.3-linux-x86_64.zip
```

\# copy unzipped executables into $PATH (e.g. `~/bin`)
```
cd bowtie2-2.2.3
cp bowtie* ~/bin/
```

\# generate bowtie2 index files of the reference sequence(s)
\# bowtie2-build -f <reference> <reference-index>   (-f is fasta format, default fastq (-q)
Build reference genome database: `bowtie2-build -f ref.fas ref`

Install `samtools`, `bcftools` and `htslib`
In root:
Download current releases from www.htslib.org/download:
samtools-1.9          bcftools-1.9          htslib-1.9  into /home/Downloads
extract each

```
# download and install ncurses for samtools
yum install ncurses-devel

# Building and installing samtools
cd samtools-1.9
./configure --prefix=/home/bodo/bin        # optional to define directory
make
make install

# Building and installing bcftools
cd bcftools-1.9
./configure --prefix=/home/bodo/bin        # optional to define directory
make
make install

# Building and installing htslib
cd htslib-1.9
./configure --prefix=/home/bodo/bin        # optional to define directory
make
make install
```

# align *.fastq files from SRA or from split genomes to the reference genome using `bowtie2`

Is the reference genome database built? `bowtie2-build -f ref.fas ref`

# syntax

bowtie2 -x <db> -1 <mate1> -2 <mate2> -U <unpaired> -S <sam output>

```
files=$(ls *_reads.fastq)
```

```
for file in $files; do name=${file%%"_reads.fastq"}; bowtie2 -p
6 -k 2 -x ref -U ${name}_reads.fastq -S ${name}.sam | echo
"done with $name"; done
```

# `bowtie2 -p 6`         if your computer has multiple processors use -p option

# `-k 2`                with -k 2, bowtie2 searches for at most 2 distinct alignments

# change .sam file to .bam file

```
files=$(ls *.sam)
```

```
for file in $files; do name={file%%".sam"}; samtools view -S -b
${name}.sam > ${name}.bam | echo "done with $name"; done
```

# sort bam file

```
files=$(ls *.bam)

for file in $files; do name={file%%".bam"}; samtools sort
${name}.bam -o ${name}.sorted.bam | echo "done with $name";
done
```

# call variants from a sorted bam file (<u>important:</u> use the same reference file as in bowtie2)

```
files=$(ls *.sorted.bam)
for file in $files; do name=${file%%".sorted.bam"}; bcftools
mpileup -f ref.fa ${name}.sorted.bam | bcftools call -mv -o
${name}.call.vcf | echo "done with $name"; done
```

# remove indels, remove low quality variants by setting a threshold
```
files=$(ls *.call.vcf)
for file in $files; do name=${file%%".call.vcf"}; cat
${name}.call.vcf | grep -v "INDEL" | bcftools view -i
'%QUAL>=80' > ${name}.calls.vcf | echo "running $name"; done
```
```
# grep -v "INDEL" -  unselect INDELs (optional, if you want
SNPs only, otherwise do not unselect)
```
```
# bcftools view -i '%QUAL>=80' -  set quality threshold of 80
```

# remove indels – F061.call.vcf

| #CHROM | POS | ID | REF | ALT | QUAL | FILTER | INFO |
|---|---|---|---|---|---|---|---|
| CP007494 | 12370 | . | TCCCCC | TCCCCCC | 124 | . | INDEL;IDV=21;IMF=0.954545;DP=22;VDB=0.00.... |
| CP007494 | 12384 | . | GCC | GCCC | 161 | . | INDEL;IDV=23;IMF=1;DP=23;VDB=0.00131074;S... |
| CP007494 | 13477 | . | GCCC | GCCCC | 171 | . | INDEL;IDV=39;IMF=0.975;DP=40;VDB=0.734156;.... |
| CP007494 | 18817 | . | C | A | 177 | . | DP=40;VDB=0.706575;SGB=-0.693145;MQSB=1;.... |
| CP007494 | 19713 | . | TGGGG | TGGGGG | 177 | . | INDEL;IDV=39;IMF=0.975;DP=40;VDB=0.74084;S |
| CP007494 | 19862 | . | GCCCCC | GCCCCCC | 173 | . | INDEL;IDV=39;IMF=0.975;DP=40;VDB=0.730783;.... |
| CP007494 | 20286 | . | AGGGGG | AGGGGGG | 177 | . | INDEL;IDV=40;IMF=1;DP=40;VDB=0.727385;SGB=... |
| CP007494 | 23192 | . | T | C | 177 | . | DP=40;VDB=0.688692;SGB=-0.693145;MQSB=1;M.. |
| CP007494 | 23198 | . | A | G | 177 | . | DP=40;VDB=0.699478;SGB=-0.693145;MQSB=1;M.. |
| CP007494 | 23806 | . | GCCC | GCCCC | 159 | . | INDEL;IDV=40;IMF=1;DP=40;VDB=0.753943;SGB... |
| CP007494 | 23826 | . | CGGGGGG | CGGGGGGG | 119 | . | INDEL;IDV=39;IMF=0.975;DP=40;VDB=0.753946;S... |
| CP007494 | 26776 | . | GCCCC | GCCCCC | 177 | . | INDEL;IDV=39;IMF=0.975;DP=40;VDB=0.699478;S... |
| CP007494 | 28257 | . | CGGGGG | CGGGGGG | 173 | . | INDEL;IDV=39;IMF=0.975;DP=40;VDB=0.730783;S... |
| CP007494 | 28910 | . | GCCCC | GCCCCC | 171 | . | INDEL;IDV=39;IMF=0.975;DP=40;VDB=0.74084;S... |
| CP007494 | 36469 | . | A | G | 175 | . | DP=40;VDB=0.727385;SGB=-0.693145;MQSB=1;M... |

```
files=$(ls *.call.vcf)
for file in $files; do name=${file%%".call.vcf"}; cat
${name}.call.vcf | grep -v "INDEL" > test.vcf
```

| #CHROM | POS | ID | REF | ALT | QUAL | FILTER | INFO |
|---|---|---|---|---|---|---|---|
| CP007494 | 18817 | . | C | A | 177 | . | DP=40;VDB=0.706575;SGB=-0.693145;MQSB=1; |
| CP007494 | 23192 | . | T | C | 177 | . | DP=40;VDB=0.688692;SGB=-0.693145;MQSB=1; |
| CP007494 | 23198 | . | A | G | 177 | . | DP=40;VDB=0.699478;SGB=-0.693145;MQSB=1; |
| CP007494 | 36469 | . | A | G | 175 | . | DP=40;VDB=0.727385;SGB=-0.693145;MQSB=1; |
| CP007494 | 49966 | . | G | A | 176 | . | DP=40;VDB=0.589467;SGB=-0.693144;MQSB=1; |
| CP007494 | 56749 | . | C | T | 176 | . | DP=40;VDB=0.611779;SGB=-0.693144;MQSB=1; |
| CP007494 | 101035 | . | C | G | 39.3362 | . | DP=18;VDB=0.0014347;SGB=-0.636426;RPB=0. |
| CP007494 | 101036 | . | G | T | 39.3362 | . | DP=18;VDB=0.00113077;SGB=-0.636426;RPB=0 |
| CP007494 | 101042 | . | C | G | 4.03223 | . | DP=18;VDB=0.00113077;SGB=-0.636426;RPB=0 |
| CP007494 | 101045 | . | C | A | 4.03223 | . | DP=18;VDB=0.0014347;SGB=-0.636426;RPB=0. |
| CP007494 | 101046 | . | A | T | 4.03223 | . | DP=18;VDB=0.00113077;SGB=-0.636426;RPB=0 |
| CP007494 | 101048 | . | T | G | 4.03223 | . | DP=18;VDB=0.00130514;SGB=-0.636426;RPB=0 |
| CP007494 | 101095 | . | G | A | 141 | . | DP=16;VDB=0.000162435;SGB=-0.688148;MQSB |
| CP007494 | 101182 | . | A | G | 165 | . | DP=24;VDB=0.00155009;SGB=-0.692717;MQSB= |
| CP007494 | 135659 | . | G | A | 176 | . | DP=40;VDB=0.658446;SGB=-0.693144;MQSB=1; |
| CP007494 | 158886 | . | A | G | 178 | . | DP=40;VDB=0.647495;SGB=-0.693144;MQSB=1; |
| CP007494 | 181835 | . | C | T | 177 | . | DP=40;VDB=0.720531;SGB=-0.693145;MQSB=1; |
| CP007494 | 185923 | . | T | C | 177 | . | DP=40;VDB=0.706575;SGB=-0.693145;MQSB=1; |
| CP007494 | 187739 | . | G | A | 91 | . | DP=8;VDB=0.000585975;SGB=-0.651104;MQSB= |
| CP007494 | 191456 | . | G | A | 178 | . | DP=40;VDB=0.595224;SGB=-0.693144;MQSB=1; |

# set quality treshold

| #CHROM | POS | ID | REF | ALT | QUAL | FILTER | INFO |
|---|---|---|---|---|---|---|---|
| CP007494 | 23198 | . | A | G | 177 | . | DP=40;VDB=0.699478;SGB=-0.693145;MQSB=1; |
| CP007494 | 36469 | . | A | G | 175 | . | DP=40;VDB=0.727385;SGB=-0.693145;MQSB=1; |
| CP007494 | 56749 | . | C | T | 176 | . | DP=40;VDB=0.611779;SGB=-0.693144;MQSB=1; |
| CP007494 | 101035 | . | C | G | 39.3362 | . | DP=18;VDB=0.0014347;SGB=-0.636426;RPB=0. |
| CP007494 | 101036 | . | G | T | 39.3362 | . | DP=18;VDB=0.00113077;SGB=-0.636426;RPB=0 |
| CP007494 | 101042 | . | C | G | 4.03223 | . | DP=18;VDB=0.00113077;SGB=-0.636426;RPB=0 |
| CP007494 | 101045 | . | C | A | 4.03223 | . | DP=18;VDB=0.0014347;SGB=-0.636426;RPB=0. |
| CP007494 | 101046 | . | A | T | 4.03223 | . | DP=18;VDB=0.00113077;SGB=-0.636426;RPB=0 |
| CP007494 | 101048 | . | T | G | 4.03223 | . | DP=18;VDB=0.00130514;SGB=-0.636426;RPB=0 |
| CP007494 | 101095 | . | G | A | 141 | . | DP=16;VDB=0.000162435;SGB=-0.688148;MQSB |
| CP007494 | 101182 | . | A | G | 165 | . | DP=24;VDB=0.00155009;SGB=-0.692717;MQSB= |
| CP007494 | 135659 | . | G | A | 176 | . | DP=40;VDB=0.658446;SGB=-0.693144;MQSB=1; |
| CP007494 | 158886 | . | A | G | 178 | . | DP=40;VDB=0.647495;SGB=-0.693144;MQSB=1; |
| CP007494 | 181835 | . | C | T | 177 | . | DP=40;VDB=0.720531;SGB=-0.693145;MQSB=1; |
| CP007494 | 187739 | . | G | A | 91 | . | DP=8;VDB=0.000585975;SGB=-0.651104;MQSB= |
| CP007494 | 191456 | . | G | A | 178 | . | DP=40;VDB=0.595224;SGB=-0.693144;MQSB=1; |

```
cat test.vcf | bcftools view -i '%QUAL>=80' > F061.calls.vcf
```

| #CHROM | POS | ID | REF | ALT | QUAL | FILTER | INFO |
|---|---|---|---|---|---|---|---|
| CP007494 | 23198 | . | A | G | 177 | . | DP=40;VDB=0.699478;SGB=-0.693145;MQSB=1; |
| CP007494 | 36469 | . | A | G | 175 | . | DP=40;VDB=0.727385;SGB=-0.693145;MQSB=1; |
| CP007494 | 56749 | . | C | T | 176 | . | DP=40;VDB=0.611779;SGB=-0.693144;MQSB=1; |
| CP007494 | 101095 | . | G | A | 141 | . | DP=16;VDB=0.000162435;SGB=-0.688148;MQSB |
| CP007494 | 101182 | . | A | G | 165 | . | DP=24;VDB=0.00155009;SGB=-0.692717;MQSB= |
| CP007494 | 135659 | . | G | A | 176 | . | DP=40;VDB=0.658446;SGB=-0.693144;MQSB=1; |
| CP007494 | 158886 | . | A | G | 178 | . | DP=40;VDB=0.647495;SGB=-0.693144;MQSB=1; |
| CP007494 | 181835 | . | C | T | 177 | . | DP=40;VDB=0.720531;SGB=-0.693145;MQSB=1; |
| CP007494 | 187739 | . | G | A | 91 | . | DP=8;VDB=0.000585975;SGB=-0.651104;MQSB= |
| CP007494 | 191456 | . | G | A | 178 | . | DP=40;VDB=0.595224;SGB=-0.693144;MQSB=1; |

# zip the manipulated file using `bgzip`

```
files=$(ls *.calls.vcf)
for file in $files; do name=${file%%".calls.vcf"}; bgzip
${name}.calls.vcf > ${name}.calls.vcf.gz; done
```

# index the variant files

```
files=$(ls *.calls.vcf.gz)
for file in $files; do name=${file%%".calls.vcf.gz"}; bcftools
index ${name}.calls.vcf.gz; done
```

# generate consensus sequence from variants (use the <u>same ref file</u> as in bowtie2)

```
files=$(ls *.calls.vcf.gz)
for file in $files; do name=${file%%".calls.vcf.gz"}; cat
ref.fa | bcftools consensus -o ${name}.cns.fa
${name}.calls.vcf.gz > ${name}.out; done
```

# change header to file name plus title string and make fasta sequence one line

```
files=$(ls *.cns.fa)
for file in $files; do name=${file%%".cns.fa"}; printf
">"$name" alignment against ref genome \n" > ${name}-cns.fasta

# changes the header to the actual strain name alignment …
# >CP007494 Bordetella holmesii ATCC 54514, IS masked
# to >$name alignment against ref genome
# writes only this line

cat ${name}.cns.fa \
| awk 'BEGIN{RS=">";FS="\n"}NR>1{seq="";for (i=2;i<=NF;i++)
seq=seq""$i; print seq}' >> ${name}-cns.fasta | echo "done with
$name"; done
# RS – Record Separator: end of record marker, default new line
# FS – Field Separator: \n separates the fields
# for all fields except the first, until the last field, seq is
seq plus the next seq
```

# join all consensus sequences into a multiple genome alignment

format: $1(Title)    $2(Seq)

```
cat ${name}.cns.fa | perl ~/bin/mergelines2.pl | awk -v FS=" "
-v OFS="\t"} '{print $1,$6}' | tr -d ">" > Bholmesii.phy
```

```
# cat ${name}.cns.fa – opens all consensus sequences
# perl ~/bin/mergelines2.pl – merges every 2 lines into 1
```

```
# from     >$name alignment against ref genome
           Sequence

# to    >$name alignment against ref genome       Sequence
```

```
# awk -v FS=" " -v OFS="\t"} '{print $1,$6}'
# fields are separated by spaces,
# print $1 (>$name) and $6 (Sequence)
```

```
# tr -d ">" – delete ">"
```

```
  $1        $2
  Strain    Sequence
```

→ Generated a multiple genome alignment

format: $1(Title)    $2(Seq)

# How to perform a <u>pairwise</u> genome comparison and display in ACT?

1. Whole Genome Blast – genome comparison
2. MSPcrunch – change blast format to Artemis input



Parkhill et al 2003 Nat Genet 35: 32-40

# How to perform a <u>pairwise</u> genome comparison and display in ACT?

1. Whole Genome Blast – genome comparison
2. MSPcrunch – change blast format to Artemis input

### `Blastall`

go to: `ftp://ftp.ncbi.nlm.nih.gov/toolbox/ncbi_tools/old`

select toolbox folder, e.g. `20120620`

click on `ncbi.tar.gz` to download

go to "Downloads" on your computer

to unpack type: `tar -xvzf ncbi.tar.gz`

to `make` type: `./ncbi/make/makedis.csh`

change directory: `cd ncbi/bin`

copy everything to: `/home/[user]/bin` (change to your bin directory)


### `MSPcrunch`

Get MSPcrunch from:

`http://sonnhammer.sbc.su.se/download/software/MSPcrunch+Blixem/`

install (or get the compiled program from me)

# How to perform a <u>pairwise</u> genome comparison and display in ACT?

1. Whole Genome Blast – genome comparison
2. MSPcrunch – change blast format to Artemis input

```
# need fasta files of both genomes
# generate data base, use "formatbd"
formatdb -i genome1.fasta -p F -o T

# -i: input Fasta file
# -p: T input type protein, F nucleotide sequence
# -o: T output database NCBI styled, F none

# output:
# genome1.nhr
# genome1.nin
# genome1.nsd
# genome1.nsi
# genome1.nsq
```

# How to perform a genome comparison and display in ACT?

1. Whole Genome Blast – genome comparison
2. MSPcrunch – change blast format to Artemis input

```
# need fasta files of both genomes
# run blastn
# syntax: blastall -p [program] -d [database] -i
[subject genome] -b [max hits] -v [max hits] -o
[output file]

blastall -p blastn -d genome1.fasta -i genome2.fasta
-o genome1-genome2.out -v 1000000 -b 1000000

Query: 4599606 tggtgaggtcgggcgaatcgtcca
                ||||||||||||||| ||||||||||||
Sbjct: 4074107 tggtgaggtcggacgaatcgtcca


Query: 4599666 caggagcttgttgcattgcgatgc
                ||||    ||||||||||||||||||||||
Sbjct: 4074047 cagggacttgttgcattgcgatgc
```

**How to perform a genome comparison and display in ACT?**

1. Whole Genome Blast – genome comparison
2. MSPcrunch – change blast format to Artemis input

```
# take blast output and change format to table
MSPcrunch -d genome1-genome1.out > genome1-
genome2.cmp

what you get:
score % sim from     to       gen1   from  to    genome2
10689 99.58 181497 183650 AXSJ      1   2154 Bb_RB50
 8233 99.82 183699 185350 AXSJ   2143   3794 Bb_RB50


 # so, we got:
 #      genome1.fasta (or genome1.gbk)
 #      genome1-genome2.cmp
 #      genome2.fasta (or genome2.gbk)
```

# Load your files in ACT



genome1.fasta (or genome1.gbk)

genome1-genome2.cmp

genome2.fasta (or genome2.gbk)

https://www.sanger.ac.uk/science/tools/artemis-comparison-tool-act

# ACT – Artemis Comparison Tool



**B. pertussis**

**B. bronchiseptica**

**B. parapertussis**

Parkhill et al 2003 Nat Genet 35: 32-40

## Let's shift gears:
run genome comparison against multiple genomes in a loop

genome1: BhinziiL60.fasta        vs        genome2:    BhinziiF582.fa
BhinziiH568.fa
BhinziiNCTC.fa
Bhinzii5132.fa
Bhinzii1277.fa
BhinziiCA90.fa

```bash
#!/bin/bash
# multiple_genomes_to_ACT.sh
# Author Bodo Linz
# run BLASTn and MSPcrunch for several genomes

DATABASE=BhinziiL60.fasta
BLASTALL=~/bin/blastall          # define location of program blastall
MSPCRUNCH=~/bin/MSPcrunch        # define location of program MSPcrunch
GENOME1=${DATABASE%%".fasta"}    # database name without ".fasta"

# has the database already been formatted?

if [ -f ${DATABASE}.nhr -a ${DATABASE}.nin -a ${DATABASE}.nsd -a
${DATABASE}.nsi -a ${DATABASE}.nsq ]; then \
     echo "The database is already formatted"
else
     formatdb -i ${DATABASE} -p F -o T
     echo "Done formatting the database $GENOME1.fasta"
fi
```

## Let's shift gears:
run genome comparison against multiple genomes in a loop

genome1: BhinziiL60.fasta        vs        genome2:

> BhinziiF582.fa
> BhinziiH568.fa
> BhinziiNCTC.fa
> Bhinzii5132.fa
> Bhinzii1277.fa
> BhinziiCA90.fa

```
# list the genomes to compare
files=$(ls Bhinzii*.fa)# generate list of files

# BLAST the target sequence against the reference genome
echo "Running blastn of $GENOME1 against
$files"
echo "-------------------------------------------"
echo ""


for file in $files; do GENOME2=${file%%".fa"}; $BLASTALL -p
blastn -d $DATABASE -i $GENOME2.fa -o $GENOME1-$GENOME2.out;
done


# loop: for all file(s) in list $files; do something; done


echo "Done with BLAST of $GENOME1 against
$files"
echo "-------------------------------------------"
```

# Let's shift gears:
run genome comparison against multiple genomes in a loop

genome1: BhinziiL60.fasta
genome2: BhinziiF582.fa, BhinziiH568.fa, BhinziiNCTC.fa, Bhinzii5132.fa, Bhinzii1277.fa, BhinziiCA90.fa

```
# Now: do the same for MSPcrunch
# list the BLAST output files
files=$(ls Bhinzii*.out)        # BhinziiL60-BhinziiF582.out etc.

# transform the blast output to ACT *.cmp table in MSPcrunch
echo "Running MSPcrunch of files
$files"
echo ""
echo "-------------------------------------------------"
echo ""
for file in $files; do name=${file%%".out"}; $MSPCRUNCH -d
$name.out > $name.cmp; done

echo "Done with MSPcrunch."
echo "-------------------------------------------------"
echo ""
echo "Done. Run ACT to visualize the genome comparison."
echo ""
```

works well for completed genomes

Problem: not suitable for genomes present as contigs
     SADLY: most genomes are incomplete
     EXAMPLE: *Acinetobacter baumannii* at ncbi genomes

# Let's download genomes

as contigs to run `blastall` and `MSPcrunch`

go to https://www.ncbi.nlm.nih.gov/genome/

type the species: *Acinetobacter baumannii*

Select: Genome Assembly and Annotation report
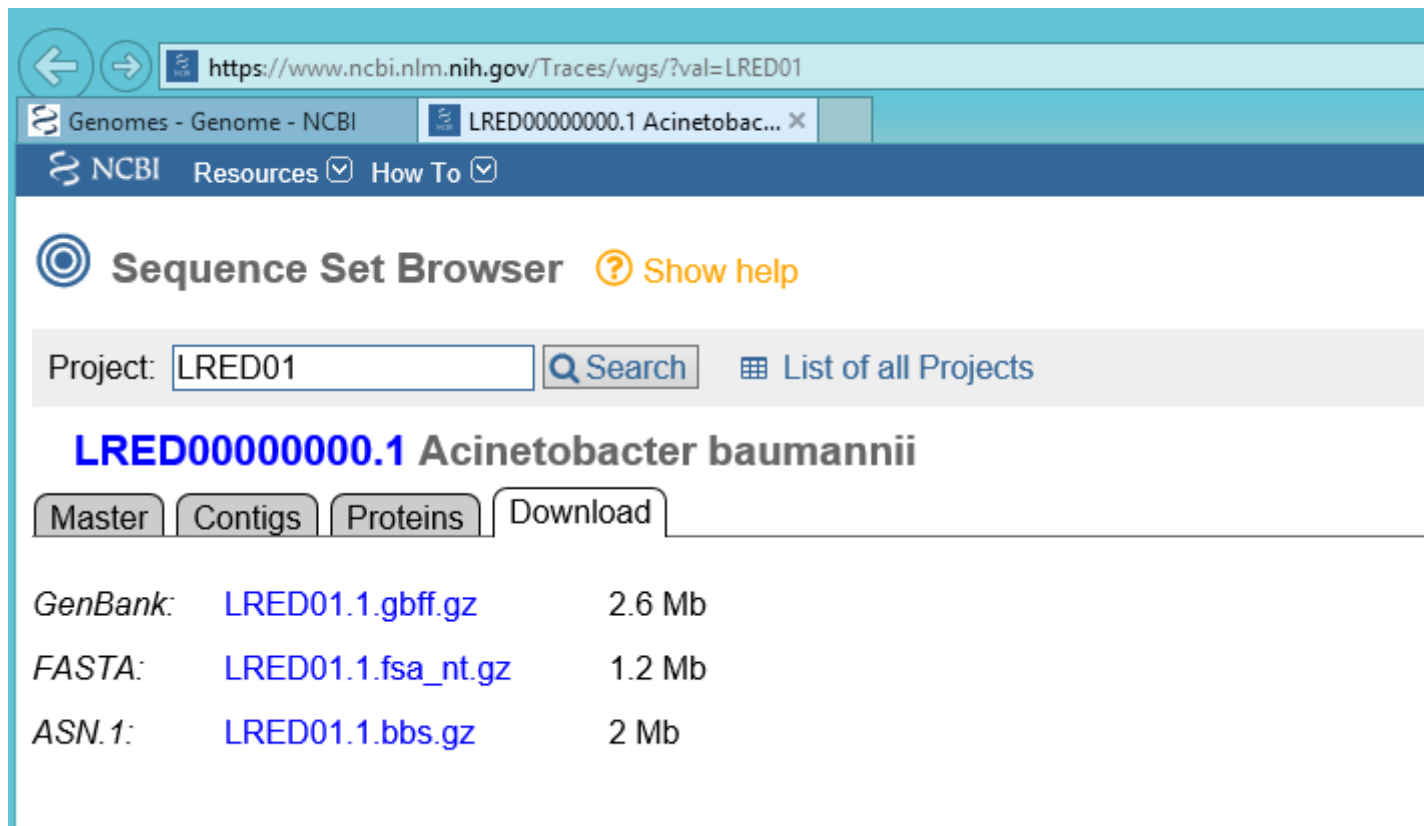
type the isolate: AB4052

click on LRED01 in WGS

# Let's download genomes

click on LRED01.1.fsa_nt.gz, download

unpack: `gzip LRED01.1.fsa_nt.gz`

rename: `mv LRED01.1.fsa_nt LRED01.1.fsa`

# We get

```
>gi|1015746545|gb|LRED01000001.1| Acinetobacter baumannii strain AB4052 LV45_contig000001, whole genome shotgun sequence
ACAAACCCGGTACGGTTCAATTAGATGGTGAATTTGCGCAAAATATTTTTGATACAGCGAAATTCTTAAA
AGGTCAGGGCAAAGTCGATCAACTTAAAGCCGATTATAAAGGCAATGTGAATTCTTCATTTTTGCAGCCT
TAAGGAGTTGTCATGAGTGTACTAGAAGCCAAACATATTCATCTGACTTTTCCTAAACAGCAAAAGCCAG
TTTTACAAGACATTAACCTAACCATTGAAGAAGGTTCTTTAACCGTGATTTTAGGTGAGTCGGGTTGTGG
CAAAACAACTTTGCTTAATATCTTGGCAGGGTTTCAAAAGCCGAGTTCAGGTGATGTGCTTGTAAATCAT
GAAGTCGTAACTGGACCAGATGTAACTCGTGCTGTTGTATTTCAAGATCACGCCTTACTTCCTTGGTTGA
ATGTTGCAGATAATGTTGGCTTCGCTTTGCAGTTAAAAGGTTTAAAGCGCGCGGATATCGAAGCACAAGT
GAACGCAATTTTAAAAATTGTGGGTTTAAGTCACGTTGAAAAAGCGAATATCTGGGAACTTTCCGGTGGT
ATGAAAAACGTGTTGGTATTGCCAGAGCTTTGATCAGTCACGCGCCGTTTATTTATTAGATGAACCTT
TTGCCGCATTAGATGCTTTTACGCGTGAAAACATGCAGCAGTTAGTGCTCGATTTATGGATTCAACAAAA
TAAAAGCTTCTTTTTGATTACTCATGACATTGAAGAAGCATTATTGCTCAGCAATCAGTTAGTTCTGATG
ACGGCGCATCCAGGCAAAATTGTAGAAACTCTACACCTCGATTTTGCCCAACGGTACCGTCAGGGTGAGT
CTATTCGCTCAATTAAATCGGATTCTCAATTTATTCAGCTCAGAGAACAGCTATTTGAAAGTTTAAGGGC
ACAAAAACAAAGCGGTAAGGAGGCGTTACCTACATGAACACTAAAGATAACGTCTATGAATATGACAAAA
CAGAGCTTAAACCTGAGTTAAATGTGCAAACAGAAAATGCTTCATTTCTATCATCATTTTTTGAGAAGCA
TCGTACTTTGGTGGTCAGCATAATCAGTGTGGGAAGTGTAGTTGCACTCTGGTTCCTCATTACTGCTTTG
CATGTTGTACCTGAACTGTTTTTTACCGAGTCCACAGGCAGTCTGGCAAAAATTTATATCGGTCAGCCAAG
AAGGCTTTATGAAAGCAACTTTGTGGCAACATTTGGCAGCCAGCATTTCTCGTGTATTTTTAGCTTTGAT
TGCTGCCGTGGTGATTGGTGTTCCGCTGGGTTTGTGGATGGGGCTGAACAAATGGGTTCGTGCTGTTCTA
GATCCTTTGGTTGAATTATTACGTCCAATTCCACCGTTAGCTTATTTGCCATTACTTGTTATTTGGTTCG
GTATTGGTGAAACCACAAAAGTACTTTTGATTTTCTTCTCGATTTTGGCGCCAGTCATTATTAGTAGTGC
GCATGGTGTGTTAAGCCATCAGCTTAATCGTGAACGTGCGGCATTGTCATTAGGGGCAAGCCAGTCACAA
GTCTTTTGGCATGTCATTTTACCAACGGCTTTGCCTCATATTATTACCGGTATTCGTATTGGTCTTGGGG
TGGGCTGGTCAACATTAGTTGCAGCTGAGTTGGTTGCAGCGGACCGTGGTATTGGTTTTATGGTGCAATC
AGCAGCACAGTTCTTAATTACCGATACGGTGATTCTGGGCATTATTGTGATTGCGATTGTCGCAGTTAGT
TTTGAGCTGTTTTTACGTTGGTTACAAAAACAGTTTTCTCCTTGGTATGGTCAGCAGTTGTAGTAAAGAA
GATGAATACAGTAGTAGCAAACTTAAATATAGAAGTGATCAAGCCTACCATTGGCGCAATTATTCACAAT
ATTGATTTGAATGCGTTAAATGAACAGACAACGCAACAAATCCAGCAGGCTTTGCTTGATCATCAGGTCA
TTTTTTTTCGAAAGCAACAATTAGCACCACAAGCACAAGCAGACTTGGCACGTAGTTTTGGTACATTGCA
TGTGCACCCGATTTATCCTTCAATTGAAGATGTACCTGAGGTGATGGTGCTCGACAGTTGGAAACAAGAT
TTGCGTGACAATGAACTTTGGCACACAGATGTGACTTTTAGTAAAACTCCACCTTTAGGTTGTGTGTTGC
AAGCTATTAAAATTCCACCTGTAGGTGGTGACACGTTGTGGTCGAGCAACACAGCAGCTTTTAAAGGACT
TCCGCTTGAGTTACAGCGAAAACTACGTGGCTTAACTGCAACCCACGATATTCGTAAGTCTTTTCCGCTT
GAGCGTTTTGCCCATAACGAAGAAGAACGTGAAAAGCTTTTGCAAACCTTTAAGCGTAACCCACCAGTGG
TTCATCCAGTGGTGCGTACTCATCCGGTTACAGGCGAGCCTTTGTTGTTTGTAAGTGAGGGCTTTACCAC
TCGCATTAATGAGTTACCCGAACAAGAAAGTGAGCAATTACTTAATTTCTTGTTTGAACATGCGACCCAA
GAGCAATTTCATTTACGCTGGAAATGGCAAGACGGTGACGTCGCGATTTGGGATAACCGTTGCACACAAC
ATAAAGCATTATTTGATTACGGAGATGCTCATCGAATTATGCACCGTGCAACTATTAACGGTGATGTGCC
ATTTTATAAAGAAGAACAACAGCCAGAGTTAGCAGAGGCTTAATTTCTTTAATTATTCTTTGTTTCAATT
CCAACGCAGCGTTTTGAGTTGGAATTGAAACAGTAACTGTTTAGCTCATTCCAAATCCTGACAATATGCC
TGTGTAATTTTTTACAGGAGGTGAGGCCCAATCACCAACTTTGCTGGTTTTTAAATTTAACTGAACTAAC
ATTTCAGCTTGTTTAACTGCTGCTGCAACGCCGTCTATCACAATTACGCCCAACTCGTTTTGTAGCTTTA
TGCATAAATCGCTCATACCTGCACAGCCCAAAACAATTGCATCGCTTTTGTCTTCCGCTAGGGCTTTTTT
GCACTCATCTCGTATGGTTCGATAAGCATCTGAGTCAGGAAGCTCCAACTCTTCAACTGCAATGTCACAA
GCTCGAACATTTTTGCAAAATGGCGTAGCCCCGTAGCGATGAGCCAGATGCCAGCTCATATTCACTGTGC
```

>fasta header contig 1
sequence
>fasta header contig2
sequence
>fasta header contig 3
sequence
etc.

# Let's assume we ran `blastall` and `MSPcrunch`: complete genome against genome in contigs

This is what we get:



All hits against the <u>first contig</u>

# Solution: modify the genome format

Solution 1:   keep only the first fasta header
                     remove all following fasta headers

```
>AHAJ01000001.1 Acinetobacter baumannii AB5711 ctg7180000006434, whole genome shotgun sequence
TGCCGCGCACTTAAAAAAGTTCGTAGATGAAATGGGTTTAACTAACATCCAAATCATGATCCCATTCGTA
CGTACAGTGTCTGAAGCAAAACGCGTCATTGAGTTATTTAGCTCAAAATTGGCTTGAAGCGTGGTGAGAA
TGGCTTAAAAGTCATCATGATGTGTGAATTACCAACTAATGCATTTGTTAGCTGAACAATTCCTTGAACT
ACTTCGATGGCTTCTACTATCGGTTCCAAACGGACTTAACTCAGGTTAACACTTGGTCTTTGACCGTGAC
TCTGGTATTGTTTCTCACTTGTTCGATGAGCGTGATGCTGCTGTAAAAGCTCTCCTTTCAATGGCAATTC
ATGCTTGTCGTAAAGCTGGTAAATATGTCGGTATCTGTGGTCAAGGACCATCAGACCACCCAGACCTTGC
AAAATGGTTAATGGAGCAAGGCATTGAATCAGTATCTCTTAACCCTGACTCGGTTTTAGACACATGGTTC
TTCCTTGCTGAA
AGTTCTGCAAGTGCTTTTTGATTTGCGTCTTCGGGATAAAGTCGAGGTGTATCCGGAAAAGTTTCGTCTA
GGTAGCGAGCGATACGGGTACTGTCTTGTATACGCTGCCCTTTATGGTCAATAACAGGTACTTTGCCCAC
TTTACTGAGCAAAGGAACTTTCGCTCCAAGAATGCCGTTGTAATTAATCGTTTCGTATGGGATTCCCTTA
AATTTCAAAGCTCTTGCAACTTTTTGGCAAAATGGAGAAATTTCCCATTGATGCAAAATAATATCCGACA
TTTATTCACCTTTATTTTTAATTGCCTGTTTTGCTCTCAGTTCCTTTTTGGAACTAATTATTAAATATAC
AGAATGTCTTTTTAAGTCAAACTATTTTTGATGACGACCAAGTTTCAAAATATAAAAAAAAGACGC
```

```
printf ">AHAJ01000001.1\n" > AHAJ01.fa                    >AHAJ01000001.1
# print everything between " "
# and save as file AHAJ01.fa
cat AHAJ01.1.fsa | grep -v ">"  >> AHAJ01.fa
# >> add to file AHAJ01.fa and save

# What does the grep command do?
```

# Solution: modify the genome format

**OR (a little more sophisticated)**

```
printf ">AHAJ01000001.1\n" > AHAJ01.fa
cat AHAJ01.1.fsa \
    | awk '{
        if(substr($1,1,1) == ">"){
                printf "";
        }else{
                printf "%s",$1;
                printf "\n";
        }
    }' >> AHAJ01.fa

# substr: substring
# if $1 at position 1 for 1 character = ">", print nothing
# else print
# printf "%s" – take the first of the following arguments ($1) and
print it as a string (s), "%d" – as a number (decimal)
# then print "\n"
# >> add to file AHAJ01.fa
```

# Note the different headers

```
>AHAJ01000001.1 Acinetobacter baumannii AB5711 ctg7180000006434, whole genome shotgun sequence
TGCCGCGCACTTAAAAAAGTTCGTAGATGAAATGGGTTTAACTAACATCCAAATCATGATCCCATTCGTA
```

```
>gi|1015746545|gb|LRED01000001.1| Acinetobacter baumannii strain AB4052 LV45_contig000001, whole genome shotgun sequence
ACAAACCCGGTACGGTTCAATTAGATGGTGAATTTGCGCAAAATATTTTTGATACAGCGAAATTCTTAAA
```

```
# modify genome input file to format ">LRED01000001.1"
cat $GENOME2 \
| awk '{
        if(substr($1,1,3) == ">gi"){
                printf ">";
                printf substr($1,19,14);
                printf "\n";
        }else{
                printf "%s",$1;
                printf "\n"
        }
}' \
> AB4052_genome.fasta
```

Ready for blast and MSPcrunch ….

# Let's walk through

```
>gi|1015746545|gb|LRED01000001.1| Acinetobacter
```

```
cat $GENOME2 \
| awk '{
        if(substr($1,1,3) == ">gi"){
# if at pos $1 the substring starting from character 1 for 3 characters
# equals (exactly) ">gi"
                printf">";
                printf substr($1,19,14);
                printf"\n";
# then print ">"
# then print the substring of 14 characters starting from character 19
# which is "LRED01000001.1"
# then print "\n" (carriage return)
        }else{
                printf"%s",$1;
                printf"\n"
# if criterion is not met, print all lines, then print "\n"
        }
        }' \
> AB4052_genome.fasta
We Get: >LRED01000001.1
        >AHAJ01000001.1 Acinetobacter baumannii AB5711 ctg7180000…
        → We took care of the different headers
```

# Thank you.

# Question?

**Let's back up some:**
How to get a specific gene sequence
from a Short Read Archive

Download Short Read Archive (SRA) from NCBI

assemble the
genome from reads

extract reads for
a specific gene

extract the gene
from the genome

assemble the gene
sequence from the reads
YASRA:
**Y**et **A**nother **S**hort **R**ead **A**ssembler

<u>Extract the reads for a certain membrane transporter gene</u>
(locus_tag BB1335 in *B. bronchiseptica* RB50)

<u>to check for a frameshift mutation in a *B. hinzii* genome</u>

Expected length without frameshift:       1416 bp
Expected length with -1 frameshift:       1415 bp

- We use **`lastZ`** and **`YASRA`**

Harris, R.S. (2007) Improved pairwise alignment of genomic DNA. Ph.D. Thesis, The Pennsylvania State University. (http://www.bx.psu.edu/~rsharris/lastz/)
download and install

**- Download SRA**
- Run FLASH to join the reads

```
flash SRR942665_1.fastq SRR942665_2.fastq -m 10 -M 100 -x 0.1
```

**rename:** `mv out.extendedFrags.fastq Bhz5132_SRR942665_joined.fastq`
**rename:** `mv out.notCombined_1.fastq Bhz5132_SRR942665_nc1.fastq`
**rename:** `mv out.notCombined_2.fastq Bhz5132_SRR942665_nc2.fastq`

## Let's dig in:

```
cat SRR942665_joined.fastq | lastz BB1335.fa[nameparse=darkspace]
/dev/stdin[nameparse=-full] --yasra90 --coverage=75
--ambiguous=iupac --format=general:name1,zstart1,end1,
name2,strand2,zstart2,end2,nucs2,quals2
| grep -v "^#"
| awk -v FS="\t" '{print $0,$4}'
| uniq -u -f 8
| awk -v FS="\t" -v OFS="\t" '{print $1,$2,$3,$4,$5,$6,$7,$8,$9}'
| sort -k 2,2n -k 3,3n
| ~/bodo.1/bin/YASRA-2.33/src/assembler -r -o -c -h /dev/stdin
> Bhinzii5132_BB1335_consensus.fa
```

➢WOW!


➢DON'T PANIC !!!


➢Let's walk through …

```
cat SRR942665_joined.fastq  # open file
| lastz BB1335.fa[nameparse=darkspace] /dev/stdin[nameparse=-
full]    # call the program lastz, which aligns the reads against
sequence BB1335.fa, our target gene
--yasra90 --coverage=75  # min identity 90%, min length 75%
--ambiguous=iupac  # IUPAC Nucleotides allowed
--format=general:name1,zstart1,end1,
name2,strand2,zstart2,end2,nucs2,quals2 # format
# name1,zstart1,end1 - our target sequence BB1335.fa
# name2,nucs2,quals2 - sequencing reads to align
| grep -v "^#"  # don't select reads that start with bad quality
| awk -v FS="\t" '{print $0,$4}' # print all $ plus $4 again
| uniq -u -f 8  # take only lines where field 8 ($8 = nucs2) is
  a unique sequence = if duplicated sequence take only once
| awk -v FS="\t" -v OFS="\t" '{print $1,$2,$3,$4,$5,$6,$7,$8,$9}'
# print all fields again
| sort -k 2,2n -k 3,3n
# sort by increasing position in target, first start then end
| ~/bodo.1/bin/YASRA-2.33/src/assembler -r -o -c -h /dev/stdin
# run the assembler
> Bhinzii5132_BB1335_consensus.fa
# save
```

# Created consensus sequence: Bhinzii5132_BB1335_consensus.fa

```
>Contig1_BB1335_0_1415
ATGCTATCGACCATATTTTCGTTTTCCTCGCTGTACTTCGCCACGCTGTTGATGTTGATC
GGCACGGGCCTGTTCAACACCTATATGGGCCTGACCCTGACGGCGAAATCCGTCAACGAA
GTCTGGATCGGCTCCATGATCGCAGGGTATTACCTCGGCCTGGTCTGCGGGGCGCGGCTG
GGCCACAAACTCATCATCCGGGTGGGCCATATCCGGGCCTTCGTGGCCTGCGCGGCCGTG
GCCACCAGCATGATCCTGCTGCAGGCCCAGATCGACTACCTGCCCATCTGGCTGCTGCTG
CGCCTGGTCTCGGGCATCATGATGGTGACCGAATTCATGGTCATCGAAAGCTGGCTCAAC
GAACAAACCGAAAACCGCCAGCGCGGCCGCGTATTCTCGGTGTACATGGTGGTCTCCGGC
CTGGGCACGGTGCTGGGACAGCTGGCGCTCACGCTCTACGGCGCGCTGGACGACGGGCCG
CTCATCCTGGTGGCCATGTGCCTGGTCCTGTGCCTGGTGCCCATCGCCGTGACGGCGCGC
TCGCACCCGCCCACGCCGCGTCCGGCGCCGCTGGACTTCTTCTTTTTCGTCAAGCGCGTG
CCGCTGGCCATGACGGTCCTGTTCGTGGCCGGCAACCTGAGTGGCGCCTTCTACGGGCTG
GCCCCGGTCTATGCCGCCAAGCATGGCCTGCAGACTTCCCAGGTGGCCTTGTTCGTCGCC
GTGTCCGTCACCGCCGGCCTGCTGTCGCAATGGCCCATCGGCTGGCTGTCCGACCGCGTC
AATCGCGCCGGCCTGATCCGTTTAACGCCGCCGTGCTGGTGCTGCTGCCCACGCTGATGT
GGGGCTGGCTGGACCTGCCTTTCTGGCTGCTGCTCTGCCTCTCGGCGCTGCTGGGCGTGC
TGCAGTTCACCCTCTATCCGCTGGGCGCGGCCCTGGCCAATGACCATGTGGAGGCCGAGC
GCCGGGTGAGCCTGAGCGCCGTGCTGCTGATGGTCTACGGGGTGGGCGCCTGCCTGGGCC
CGCTGGTCGCCGGCATCCTCATGTCGCTCGGCGGGCACGCCATGTACTACGTCTTCGTGC
CGGCCTGCGCCCTTATCCTGGTCTGGCGCGTGCGGCCCAGCGCCGTCACTGGCGTGCACC
AGGTCGAGGAGGCGCCGGTGCAATTCGTGCCCATGCCCGACACGCTGCAGTCCTCGCCCG
CCATGGTGGCCTTGGATCCCCGTGTGGATCCCGAGGTGGACCCGGCCATGGAGATGGTCA
CGCCCGAGGCCGGCGTGGTGCAGCCGCCGCCGCCGGCCGCCGAACCCGCTGCCGGCACGG
CGGCCTTCGACAACGTCGTGGCCGAGCCGGGCGAGCCGGCCACCGTCCTGTCCGCAGACG
GCGCGCCGAGTCCGCGCACAGGGACGGACGCCTGA
```

# How many nucleotides?
Easiest solution:

```
cat Bhinzii5132_BB1335_consensus.fa \
| grep -v ">" | tr -d "\n" | wc -L

# output: 1415
# grep -v ">" - select lines that do not contain ">"
# → only sequence without fasta header
# tr -d "\n" - translate carriage return "\n" to nothing
# → concatenates all sequence lines
# wc - word count
# wc -L returns number of characters in longest line
# Result: 1415
# That means, we are dealing with the frameshift gene variant
```