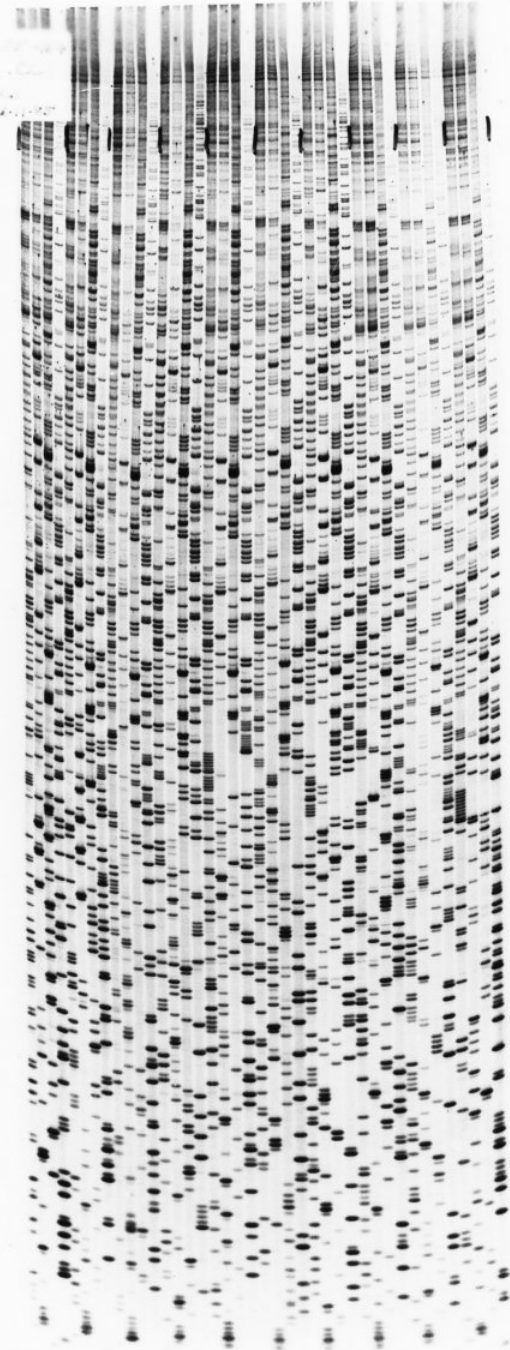# MinION sequencing

# What's next?

EvoPad Summer School 2019

Victoria Shabardina
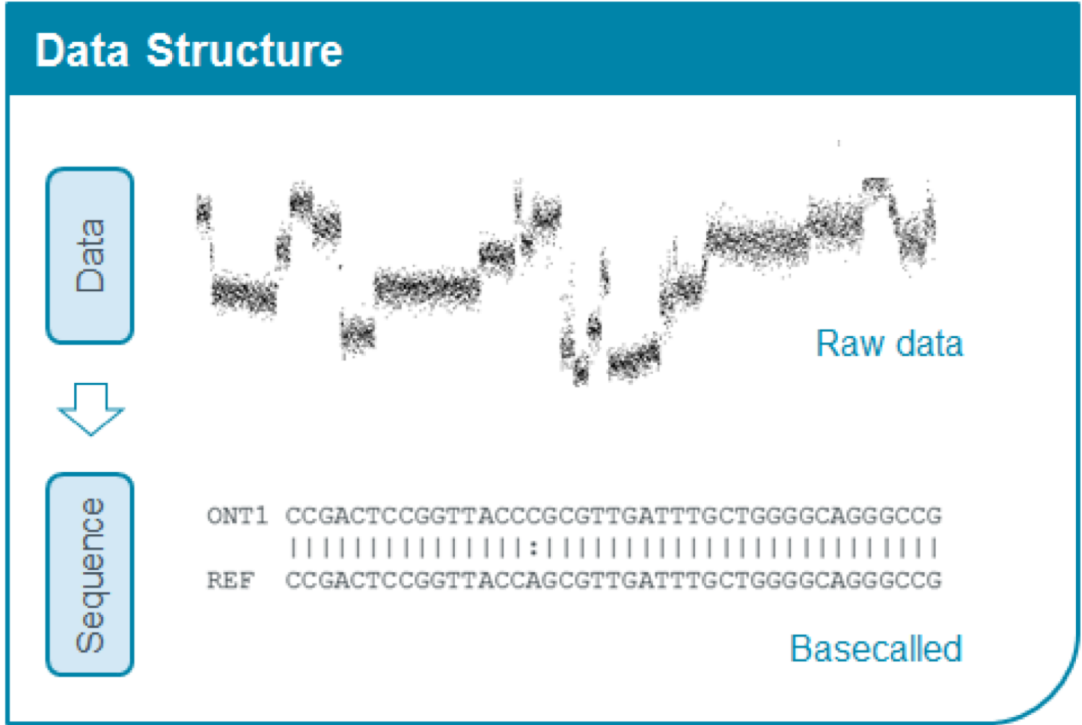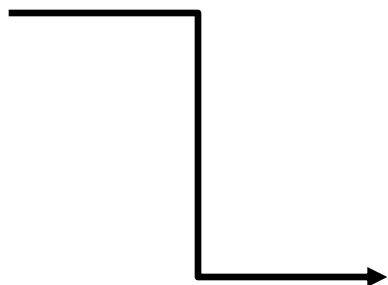
1977

2m

Elegance of the technical progress ;)

2017

## Data Structure

Data

Raw data

Sequence

ONT1 CCGACTCCGGTTACCCGCGTTGATTTGCTGGGGCAGGGCCG
     |||||||||||||||||:|||||||||||||||||||||||
REF  CCGACTCCGGTTACCAGCGTTGATTTGCTGGGGCAGGGCCG

Basecalled

# Reccurent Neural Network (RNN) – works like your brain! It can learn on the previous data and improve its performance on new data



Nanopore basecallers are trained on many sequenced data, so you can run it on your data even if you are sequencing first time

MinION sequencing is controlled by MinKNOW software → MinKNOW runs basecalling in parralel → If you stop sequencing MinKNOW stops basecalling !!!

Alternative: Guppy software (command line)

# FAST5 file format

FAST5 is a type of the Hierarchical Data Format (HDF) and designed for the storage of big datasets

It is binary – not readable by human

HDFview – the tool to see HDF files

# FASTQ file format

Each entry consists of 4 lines:

1 @header_name_of_sequence/read
2 sequence
3 +
4 quality (coded by ASCII symbols: https://en.wikipedia.org/wiki/FASTQ_format#Encoding)

Poretools is the tool developed for MinION to convert .fasta5 files to .fastq or .fasta files

Guppy basecaller allows generating directly .fast5 or .fastq files ☺

HDFviewer window:

(https://support.hdfgroup.org)

FASTQ file

Raw signal
(current from a MinION channel)

# FASTA file format

contains 2 lines:
1 >header
2 sequence

>ff6c98dd-bce9-4a2b-bf13-081841413c94_Basecall_2D_minion_20170511_Ae_aegypti
GCGCTGGTTCAGTTACATATTGCTAGGGTTAAGCAGTGGTGACCACAGATTTTTATGATTTATGGATT
CTTTTCTTCTGGCTACATTACTGGAACAGAGCCTGCTTCTCAACAGTGTTCTTATGAACGCTTCAGCTTA
GTATAAAGGC
>4c8a2487-0e13-41a6-ac7b-6cd2fbf5eb95_Basecall_2D_minion_20170511_Ae_aegypti
TACGCGGTGACAAAAACGTGCGTACCGGCAACCGCATGTTGAAACAGGAAAACGTACAAAGGACCC
TCGCAAAATGCGCGACAAAATCTGCAACGTACAACATGCGATAAACGTGCGTGAGGAGATC
>…..
………..

> sign is convenient marker to browse through the FASTA files

*Mind_gaps_in_command_line and file_names!*

Basecalling with Guppy, an ONT produced tool-kit

**Guppy** can do 4 different jobs:          Commands:

**1** 1D basecalling                              *guppy_basecaller*

**2** 1D2 basecalling                             *guppy_basecaller_1d2*

**3** Debarcoding (demultiplexing)                *guppy_barcoder*

**4** Alignment                                   *guppy_aligner*

Guppy can be used on Windows, Mac OS, and Linux

# Basecalling with Guppy, an ONT produced tool-kit

**Guppy** can do 4 different jobs:                    Commands:

**1** 1D basecalling                                  *guppy_basecaller*

**2** 1D2 basecalling                                 *guppy_basecaller_1d2*

**3** Debarcoding (demultiplexing)                    *guppy_barcoder*

**4** Alignment                                       *guppy_aligner*

## Workflow:

*guppy_basecaller / MinKNOW → guppy_barcoder*

*guppy_basecaller / MinKNOW → guppy_aligner*

Basecalling with Guppy

How to use your computer efficiently?

Concider: RAM (random-access memory) and number of CPUs (central processing unit).

Guppy_basecaller (1D) uses 1GB per 1 CPU + 4 GB

4 CPUs: 1x4 + 4 = 8 GB of RAM

Guppy_basecaller_1d2 uses 2GB per 1 CPU + 4GB

# Basecalling with Guppy

*guppy_basecaller --help*

One line:
*guppy_basecaller –i input/reads.fast5  -s output/reads.fastq --flowcell FLO-MIN107*
*--kit SQK-LSK108 --qscore_filtering –q 0 --num_callers 1 --cpu_threads_per_caller 1 -r*

# Basecalling with Guppy

One line:
*guppy_basecaller –i input/reads.fast5  -s output/reads.fastq --flowcell FLO-MIN107*
*--kit SQK-LSK108 --qscore_filtering  –q 0  --num_callers 1  --cpu_threads_per_caller 1  -r*

-i  (where is your input files)
-s  (where you want to save the output)
*--flowcell*
*--kit*

*guppy_basecaller --print_workflows*

--qscore_filtering  (sorts reads into 'pass' and 'fail' folders, --min_qscore is 7 by default)
-q 0  (writes all reads per run in one FASTQ file, default is 4000 reads per file)
-r – recursive  (will go through all files in the folder)
*--num_callers* and *--cpu_threads_per_caller* tell how much of your computer power to use

# Basecalling with Guppy

Other options:

GPU run possible

--fast5_out  (output FAST5 and FASTQ files, default – only FASTQ)
--compress_fastq  (generates gzip output file)

RNAseq:
--reverse_sequence (RNA strain goes through the pore backwards)
--u_substitution (T → U)

--resume – useful if basecalling was interrupted

--calib_detect - callibration strand detection

Adapter trimming is by default

# FAST5 file structure of a basecalled read

**Basecalled data format in Guppy**

**The read .fast5 file structure looks as follows:**

```
/{attributes: file_version}
| -UniqueGlobalKey
|   -tracking_id {attributes: standard tracking-id fields}
|   -channel_id {attributes: channel_number, digitisation, offset, range, sampling_rate}
|   -context_tags {attributes: set when the experiment is configured}
| -Raw
|   -Reads
|     -Read_42 {attributes: start_time, duration, read_number, start_mux, read_id}
|       -Signal {samples}
| -Analyses/
|   -Segmentation_000 {attributes: name, version, time_stamp}
|     -Summary/
|       -segmentation {attributes: has_template, has_complement, duration_template, first_sample_template, num_events_tem
|   -Basecall_1D_000 {attributes: name, version, time_stamp}
|     -BaseCalled_template
|       -Events {annotated event data}
|       -Fastq {embedded fastq file}
|     -BaseCalled_complement
|       -Events {annotated event data}
|       -Fastq {embedded fastq file}
|     -Summary
|       -basecall_1d_template {attributes: called_events, event_stride, mean_qscore, sequence_length, strand_score, stay_prob,
```

# NanoPipe – interactive tool for MinION sequencing analysis

*Developed in the Institute of Bioinformatics, University of Münster, Germany*

http://bioinformatics.uni-muenster.de/tools/nanopipe

**Sequencing**



Electric signal in FAST5 files

**Basecalling**

FASTQ files

**Bioinformatics analysis**

**NanoPipe workflow: what does "*pipe*" mean?**

Bioinformatics analysis

MinION reads

Target

LAST alignment

nucleotides count per position

consensus

alignment statistics

MAF and BAM files

consensus-target alignment

nucleotide plots

polymorphisms

# Key steps of NanoPipe: LAST

LAST sequence aligner maps MinIon-produced reads to a target (selected region, exon, gene, genome)

Martin Frith
University of Tokyo,
Division of Biosciences



Align reads to genome

Query

Target    Genome

Other aligners: BLAST (psi-BLAST, delta-BLAST), HMMER, MUSCLE, MAFFT...

File formats that are important to know when working with mapped reads:
        .maf
        .sam
        .bam

# **.maf** file format – shows pairs of aligned sequences with the coordinates

```
# LAST version 923
#
# a=15 b=3 A=15 B=4 e=101 d=58 x=100 y=44 z=100 D=1e+06 E=8.45451e+07
# R=10 u=0 s=2 S=1 M=0 T=0 m=10 l=1 n=10 k=2 w=1000 t=4.40086 j=3 Q=0
# /bioinf/projects/nanopipe2/targets/plasmodium/target
# Reference sequences=3 normal letters=5914
# lambda=0.21831 K=0.309523
#
#     A   C   G   T
# A   4 -15  -4 -22
# C -18  10 -20 -15
# G  -8 -18   9 -18
# T -23 -11 -17   4
#
# Coordinates are 0-based.  For - strand matches, coordinates
# in the reverse complement of the 2nd sequence are used.
#
# name start alnSize strand seqSize alignment
#
# m=0.01 s=101
#
a score=168 mismap=0.00337
s Pf3D7_07_v3:403089-404828:+:PfCRT_1                                     1338 42 + 1740 AAATATATAAATAAATAAATATATATATATATATATATAT
s 004749f7-eefe-432b-9793-7018b5abcc8c_Basecall_1D_template:1D_001:template 1942 42 - 2131 AAATATATAAATAAATAAATATATATATATATATATATAT
p                                                                              $(,0367899999999999999999999999888876530,($

a score=147 mismap=2.05e-05
s Pf3D7_07_v3:404757-406466:-:PfCRT_2                                     1158 104 + 1710 AAATAAAATAATAATATATATAATATATATAATATTGTTTTATTTAATTATTATTGTTA-AAATATAT----ATAAATGTCTCTTATAATTTT
s 02efe5aa-8e63-4121-840e-681be1549390_Basecall_1D_template:1D_001:template  943  95 - 1712 AAATAAAATA-TCATATATA------ATATAATAC-G-TTTATTTAATTATTATTGTTGCAAATATATTATAACAAA----CCTTATAATTTT
p                                                                              "#$$$$$$$$$$$$$%%%%%%%%%%%'(((())))))))-159=@DGJLMMMMMMMMLKKKKKKKKJJIIIIIIIIIIIIIIIIJN0000000000

a score=145 mismap=1.76e-05
s Pf3D7_07_v3:404757-406466:-:PfCRT_2                                     918 46 + 1710 TGTCGATAATCTATAAAAAG-CATAGAAAATGAAAAATTATATGGTT
s 02efe5aa-8e63-4121-840e-681be1549390_Basecall_1D_template:1D_001:template 732 45 - 1712 TGTCGATAATCTATAAAA-GTTATAGAAAACG-AAAATCATATGGTT
p                                                                              %.3<DHKMOPPPPPPPPONNNNNNNNNNNMLKIIIHHGEEDB?;4*%
```

# .bam and .sam file format – informative about alignment of sequences to a target

.bam is a binary file (humans can not read), it contains a lot of information about the alignment. Important part is the HEADER and FLAGs specifications (https://broadinstitute.github.io/picard/explain-flags.html)

.sam format is the readable version of .bam (some info is casted away):

# Summary: files' formats

Files' formats used for storing sequences:

**FAST5** – big, binary (cant read), contain a lot of metadata
**FASTQ** – readable by human, contains sequences and sequence quality
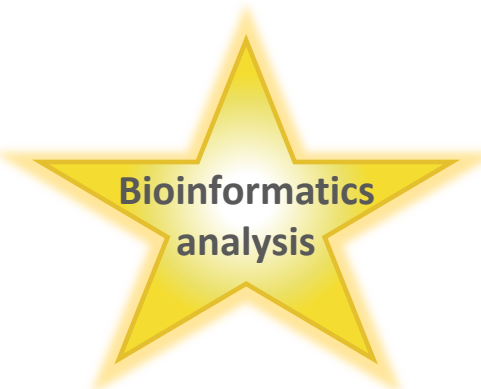**FASTA** – readable, contains sequences

Files' formats used for storing results of sequence alignment:

**maf** – contains pairs of aligned sequences with the alignment's coordinates; for example, used by LAST
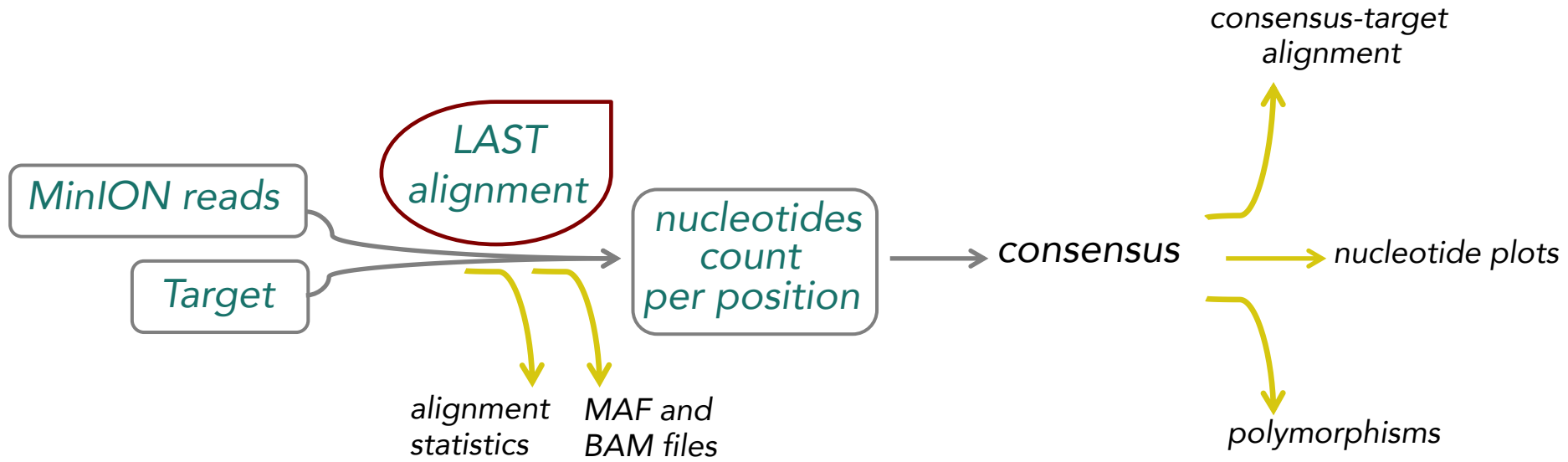**bam** – binary format, includes aligned sequences, coordinates, information about bioinformatics processing, quality, …
**sam** – human readable version of .bam, much bigger in size

**FASTQ, FASTA and bam files are widely used in all DNA/RNA bioinformatics analysis**

# NanoPipe workflow: what does "*pipe*" mean?

**Bioinformatics analysis**

*MinION reads*

*Target*

*LAST alignment*

*alignment statistics*

*MAF and BAM files*

*nucleotides count per position*

*consensus*

*consensus-target alignment*

*nucleotide plots*

*polymorphisms*

**NanoPipe helps us to…**

- See if our sequencing worked: how many reads were mapped to the target and where exactly, what part of each read mapped (Alignments length distribution)

- Detect insertions/deletions and single nucleotide variations

- Visualization of the experiment in NanoPipe and in IGV-viewer (bam and indexed bam files)

- FASTA file with the consensus sequence

Useful links:

Our lab ;p -
http://bioinformatics.uni-muenster.de


Genome viewers -
https://software.broadinstitute.org/software/igv/

https://genome.ucsc.edu/