# GigaScience

## NanoPipe - a web server for nanopore MinION sequencing data analysis
### --Manuscript Draft--

| | |
|---|---|
| Manuscript Number: | GIGA-D-18-00394R1 |
| Full Title: | NanoPipe - a web server for nanopore MinION sequencing data analysis |
| Article Type: | Technical Note |

| | |
|---|---|
| Abstract: | Background<br><br>The fast-moving progress of the third generation long read sequencing technologies will soon bring the biological and medical sciences to a new era of research. Altogether the technique and experimental procedures are becoming more straightforward and available to biologists from diverse fields, even without any profound experience in DNA sequencing. Thus, the introduction of the MinIONTM device by Oxford Nanopore Technologies promises to "bring sequencing technology to the masses" and also allows quick and operative analysis in field studies. However, the convenience of this sequencing technology dramatically contrasts with the available analysis tools, which may significantly reduce enthusiasm of a "regular" user. To really bring the sequencing technology to every biologist, we need a set of user-friendly tools that can perform a powerful analysis in an automatic manner.<br><br>Findings<br><br>NanoPipe was developed in consideration of the specifics of the MinIONTM sequencing technologies, providing accordingly adjusted alignment parameters. The range of the target species/sequences for the alignment is not limited, and the descriptive usage page of NanoPipe helps a user to succeed with NanoPipe analysis. The results contain alignment statistics, consensus sequence, polymorphisms data, and visualization of the alignment. Several test cases are used to demonstrate efficiency of the tool.<br><br>Conclusions<br><br>Freely available NanoPipe software allows effortless and reliable analysis of MinIONTM sequencing data for experienced bioinformaticians, as well for wet-lab biologists with minimum bioinformatics knowledge. Moreover, for the latter group, we describe the basic algorithm of actions necessary for MinIONTM sequencing analysis from the first to last step.<br><br>Issue Section<br><br>TECHNICAL NOTE |

| | |
|---|---|
| Corresponding Author: | Wojciech Makalowski<br>Univeristy of Münster<br>Münster, GERMANY |
| Corresponding Author Secondary Information: | |
| Corresponding Author's Institution: | Univeristy of Münster |
| Corresponding Author's Secondary Institution: | |
| First Author: | Victoria Shabardina |
| First Author Secondary Information: | |
| Order of Authors: | Victoria Shabardina |

| | Tabea Kischka |
| --- | --- |
| | Felix Manske |
| | Norbert Grundmann |
| | Martin C. Frith |
| | Yutaka Suzuki |
| | Wojciech Makalowski |

| **Order of Authors Secondary Information:** | |
| --- | --- |
| **Response to Reviewers:** | Reply to reviewers' commentaries for Submission GIGA-D-18-00394 |

Dear Editor,
Thank you very much for processing our manuscript and giving the opportunity to work on its flaws and send the improved version. We also would like to thank the reviewers for the valuable criticism and attentive reading of the presented work. Their comments helped to improve our manuscript significantly. We hope that our replies and corresponding revisions of the manuscript are satisfactory and you will find the manuscript ready for publication at the GigaScience journal.

Sincerely,
Wojciech Makałowski

Reviewer 1.
Comment 1. A rather serious disagreement between the text and the web software is the set of reference genomes; the github archive appears to have no reference genomes. The manuscript claims 12 database and the 7 offered on the website don't fully overlap that set.

Reply: The manuscript and the NanoPipe menu, both, offer 9 pre-processed targets (page 8 of the manuscript): for human (hg38), RefSeq accession GCF_000001405.38; Escherichia coli, RefSeq accession GCF_000005845.2; Caenorhabditis elegans, RefSeq accession GCF_000002985.6; Droshophila melanogaster, RefSeq accession GCF_000001215.4; Mus musculus, RefSeq accession GCF_000001635.26; Arabidopsis thaliana, RefSeq accession GCF_000001735.4; Plasmodium falciparum strain 3D7, downloaded from plasmodb.org, version=2013-03-01; a representative genome for Camponotus floridanus, RefSeq accession GCF_003227725.1; and Dengue virus genome variants for serotyping (NC_001477.1, NC_001474.2, NC_001475.2 and NC_002640.1 for variant 1, variant 2, variant 3 and variant 4, respectively).
It is a good idea to keep all the information about targets in github repository: the list of targets and their source were added to the directory https://github.com/IOB-Muenster/nanopipe2/targets

Comment 2. I tried uploading a small set of E.coli reads and an E.coli reference (one of the manuscript-promised, not-on-server cases!) but I seem to be permanently parked behind two other jobs.
And Comment 3. The website has a "View Testcase" button -- that yields an error message screen.

Reply: We have tried the test case and the job run from the different computers, locations and web browsers and did not encounter any problems. Is it possible that the errors that the reviewer was experiencing had happened due to the local internet connection? What internet browser the reviewer used?

Comment 4. Similarly, the github package does not appear to contain a useful test case or any code to check the installation except running the package -- so if anything goes wrong it could be difficult for a novice to determine which of the long list of dependencies (11 Perl modules, 9 Python modules, 2 other tools)

Reply: The installation package deposited at github is aimed at users with some bioinformatics and coding knowledge and who would want to explore and modify the tool themselves. The novice users are offered to use the online version. Nevertheless,

we thank the reviewer for this comment and agree that, indeed, it is helpful to know what has gone wrong, if anything has. We introduce the following improvements in the installation package:
1) The installation script has been written with the complementary explanations, please, see https://github.com/IOB-Muenster/nanopipe2/blob/master/install.sh and https://github.com/IOB-Muenster/nanopipe2/blob/master/install.txt
2) The check.sh script at the same directory (https://github.com/IOB-Muenster/nanopipe2) has been written for a user to be able to monitor if all the required packages are present on the user's computer.
3) A testcase has been deposited to check if the installation was successful. The install.txt file on the github repository includes description how to run the testcase. If the test script gives any errors, the user should double check the installation procedure and/or contact us via "Contact" option on the NanoPipe web page).

Comment 5. The authors might also consider distributing as a docker container and/or conda package with all dependencies covered.

Reply: Docker and conda are the good solutions for an installation packages, nevertheless we think that our install.sh script on github is covering well all the necessities. Besides, the accent of the tool is on its online application.

Comment 6. Streamlining the process of creating a new target database would be desirable -- the "install.txt" file gives a 5 step protocol -- two of which should be combined ("Create the target database" & "target.fasta". The lastdb step really should be wrapped in something that checks the new database information for consistency

Reply: The createtarget.sh script has been added to the https://github.com/IOB-Muenster/nanopipe2 github directory. It simplifies the target generation process. The explanation of usage can be found in the install.txt document.

Reviewer 2.
Part 1. Notes about the ONT technology
Thank you very much for the suggestions related to the introduction part; the facts from the Clive Brown's talk are, indeed, very exciting. Although we don't think that all the technical details should be included in our manuscript, since it is aimed at the broad audience, we have modified the text accordingly to give the general impression about the technology progressing. Please, see lines 59-61 and 75-77

Part2.
Comment 1. Introduction doesn't mention Metrichor, despite using TM in where MinION is used.

Reply: thank you for noting this, we filled that gap (line 82).

Comment 2. "for whole genome sequencing by MinION TM a researcher can expect read lengths up to several thousand nucleotides"
 - longest read observed so far is 2.3 *million* nucleotides.

Reply: This information has been corrected (line 208).

Comment 3. MAP005/MAP006 kits and MAP003 flow cell in line 245 suggest a very old kit (~2-3 years old).This is unusual for a paper about to be published, but is consistent with one other GigaScience
paper that I've seen (?Sara Goodwin). I'd be interested to know how long this paper languished
in pre-review doldrums until GigaScience accepted it for sending out to reviewers.
And: Line 270 for H1975 suggests SQK-LSK108/FLO-MIN107 R9; a bit more recent (but still old).

Reply: Indeed, the flow cells used for some experiments were not of the latest version. It depended on the collaboration we had had with the wet labs and the work progress. Nevertheless, the usage of not so recent flow cells in regard to the software development brings more robustness to the analysis, since the sequencing precision in this case is worse than with the newest equipment. We were able to show that the data

processing with FLO-MAP003 and FLO-MIN107 flow cells provide with the reliable results. That ensures that the modern flow cells and sequencing kits will be not worse.

Comment 4.  Line 292 -- what's the "standard procedures" for poly-A RNA extraction? Was this using poly-A  bead selection? Why not strand switch sequencing with ONT adapters (available/recommended in ONT protocols from August 2017)?

Reply: We thank the reviewer for this comment: indeed, we happened not to be very precise in the description. The details about the RNA extraction has been added: please, refer to lines 297-298. We chose to use poly-A bead extraction method, because this procedure was optimized in our collaborators' lab to succeed with the RNA extraction from ant material. As the work with insect tissues is prone to difficulties, we did not want to change the established work flow in the lab that have expertise in working with ant species.

Comment 5. Given that this is a paper about *software*, I can't see any obvious reason why
new samples were sequenced. It'd be nice to see this algorithm applied to recent large public
human datasets (e.g. nanopore-wgs-consortium: ultra-long-read runs, or full-length RNA/cDNA
runs). Reads can be subset, as necessary, to cover particular genomic regions. This will help
encourage people to use existing public datasets for their own software.

Reply: Following the suggestions of the reviewer and the editor we have added one more test case on the recently published direct RNA long read sequencing data of Vaccinia virus and its host. Please, see lines 308-319. The jobs IDs are 154401029652282 (mapped to the green monkey (host) genome) and 154400756783780 (mapped to the virus genome).

Comment 6. Text mentions that NanoPipe was used in Bangkok in 2017, but the github commits only go back to September 18 this year. Is there any reason why NanoPipe wasn't version controlled in 2017?

Reply: We were running NanoPipe through the series of tests and changes in 2017, that's why the github page was not created back then yet. The tool was at beta-version.

Comment 7. Note: IUPAC annotation for any of four nucleotides is 'N', not 'X'
 [https://github.com/IOB-Muenster/nanopipe2/blob/f2026e16b8942ec1cb60157b032a9c4bcbfebef7/modules/nanopipe2/calculate/analyze.pm#L341]
 [https://www.bioinformatics.org/sms/iupac.html]

Reply: This, indeed, is different from IUPAC symbols. The issue is clarified in the manuscript, line 228, and in the usage explanations within the NanoPipe.

| Additional Information: | |
| --- | --- |
| Question | Response |
| Are you submitting this manuscript to a special series or article collection? | No |
| Experimental design and statistics<br><br><br>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist. Information essential to interpreting the data presented should be made available | Yes |

| | |
|---|---|
| in the figure legends.<br><br>Have you included all the information requested in your manuscript? | |
| **Resources**<br><br>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.<br><br>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist? | Yes |
| **Availability of data and materials**<br><br>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the "Availability of Data and Materials" section of your manuscript.<br><br>Have you have met the above requirement as detailed in our Minimum Standards Reporting Checklist? | Yes |

# NanoPipe - a web server for nanopore MinION sequencing data analysis

Victoria Shabardina[1], Tabea Kischka[1], Felix Manske[1], Norbert Grundmann[1], Martin C. Frith[2,3,4], Yutaka Suzuki[5], and Wojciech Makałowski[1]

[1] Institue of Bioinformatics, University of Muenster, Germany

[2] Artificial Intelligence Research Center, AIST, Japan

[3] Department of Computational Biology and Medical Sciences, The University of Tokyo, Japan

[4] AIST-Waseda University Computational Bio Big Data Open Innovation Laboratory, Japan

[5] Laboratory of Systems Genomics, Department of Computational Biology and Medical Sciences, The University of Tokyo

Corresponding author: Wojciech Makalowski, wojmak@uni-muenster.de

ORCIDs: Tabea Kischka: 0000-0002-7117-5575; Martin C. Frith: 0000-0003-0998-2859; Wojciech Makalowski:0000-0003-2303-9541

**Abstract**

*Background:* The fast-moving progress of the third generation long read sequencing technologies will soon bring the biological and medical sciences to a new era of research. Altogether the technique and experimental procedures are becoming more straightforward and available to biologists from diverse fields, even without any profound experience in DNA sequencing. Thus, the introduction of the MinION™ device by Oxford Nanopore Technologies promises to "bring sequencing technology to the masses" and also allows quick and operative analysis in field studies. However, the convenience of this sequencing technology dramatically contrasts with the available analysis tools, which may significantly reduce enthusiasm of a "regular" user. To really bring the sequencing technology to every biologist, we need a set of user-friendly tools that can perform a powerful analysis in an automatic manner.

*Findings:* NanoPipe was developed in consideration of the specifics of the MinION™ sequencing technologies, providing accordingly adjusted alignment parameters. The range of the target species/sequences for the alignment is not limited, and the descriptive usage page of NanoPipe helps a user to succeed with NanoPipe analysis. The results contain alignment statistics, consensus sequence, polymorphisms data, and visualization of the alignment. Several test cases are used to demonstrate efficiency of the tool.

*Conclusions:* Freely available NanoPipe software allows effortless and reliable analysis of MinION™ sequencing data for experienced bioinformaticians, as well for wet-lab biologists with minimum bioinformatics knowledge. Moreover, for the latter group, we describe the basic algorithm of actions necessary for MinION™ sequencing analysis from the first to last step.

**Issue Section:** TECHNICAL NOTE

*Keywords:* sequencing technologies, long reads sequencing, bioinformatics software, MinION, Oxford Nanopore

**Background**

Recent years have witnessed a DNA sequencing boom due to the constantly improving technologies and, consequently, the accessibility of sequencing to a large spectrum of customers including scientists and medical practitioners. Researchers in many fields, from metagenomics to plant physiology to medicine have been implementing sequencing experiments into their research. Oxford Nanopore Technologies (ONT) essentially accelerated this process by introducing the MinION™ sequencer, a portable device with minimum requirements for technical skills and bioinformatics knowledge. Thus, DNA sequencing experiments became feasible even in field studies, in small laboratories and soon will be available for medical applications in clinics [1].

The NCBI PubMed database includes 261 scientific articles containing "Oxford Nanopore" phrase published between 2009 and 2018 (by September 1, 2018), the majority of which were published in the last three years (see Fig. S1 in Supplementary). This is pointing, both, at the increased popularity of the ONT sequencers and at the considerable improvement of the technology and sequencing quality in the last three years. For example, the R10 version of the MinION™ flow cell was recently announced and it promises to improve sequencing quality, including for homopolymer stretches. High throughput and long reads allow diverse applications of MinION™: virology [2], [3], plant pathology and agriculture [4], [5], tuberculosis studies [6], metagenomics and diet [7], veterinary research [8]; and as a portable platform: field biodiversity studies [9], detection of Ebola virus in patients on the spot [10], [11], sequencing in space [12]. Not to forget fundamental applications for long read sequencing studies, such as de novo

3

70 genome assembly, improvement of existing genome assemblies and discovery of
71 structural variants and long repeats [13]–[16]. ONT sequencing is favorable for
72 microbiology research as small-sized bacterial genomes can be covered in just one
73 MinION™ read [17], [18], thus providing high resolution in genome architecture.

74 ONT supplies its users with the necessary software to perform base calling, i.e. converting
75 of MinION™'s electrical signals into a sequence of nucleotides: the on-run MinKNOW and
76 offline Albacore. Both applications utilize the complex, recurrent neural network (RNN)
77 algorithm, which is recently very popular in computer science. It allows the software to
78 learn from existing data and improve its performance. It is worth noting that the base
79 calling process is the central to improving the accuracy of ONT sequencing technology
80 and its algorithm is being constantly improved and updated. The output is a collection of
81 FAST5 and/or FASTQ files containing the base-called sequences. These are the files that
82 are used for any sequence analysis in bioinformatics, thus, the base-caller can be called a
83 "gate" from MinION™ into data interpretation. Nevertheless, the range of ONT provided
84 analysis tools is limited, and concerns only specific applications, excluding general
85 processing, which is left to the user. For example, the EPI2ME software suite, based on
86 Metrichor platform, includes applications for barcode analysis, metagenomics and
87 antimicrobial resistance analyses and some technical tests [19]. Several research groups
88 have been recently focusing on the development of MinION™ specific bioinformatics tools
89 [20]–[22], although most of these require considerable bioinformatics knowledge. These
90 conditions impede benefits that MinION™ based DNA sequencing could bring to medical
91 practitioners and researchers with less IT experience. To fill this gap, we developed
92 NanoPipe, a web-driven automatic pipeline that can quickly and effortlessly process data
93 produced by MinION™, as well provide necessary files for further bioinformatics analysis
94 if required.

**Methods**

95

96 NanoPipe can be conceptually divided into four stages: 1) data uploading in FASTA or

97 FASTQ formats, 2) alignment of MinION™ reads against the target sequences, 3)

98 alignment analysis, and 4) results display (see Fig. 1). It was developed with no-IT-

99 experienced users in mind, hence it provides a web-driven interface with the usage page

100 describing the main features of the tool. The start-of-analysis page is simple and intuitive.

101 A user must, first, choose a target genome from a NanoPipe's list (see further) or upload

102 their own target sequence. The next step is the essential part of the pipeline: mapping of

103 sequencing reads to the target using version 946 of the LAST sequence aligner (LAST,

104 RRID:SCR_006119). LAST accounts for the MinION™ specific sequencing errors, thus it

105 generates highly reliable results. It can determine the rates of insertion, deletion and each

106 kind of substitution in a type of data (e.g. MinION™ reads of AT-rich *Plasmodium* DNA)

107 [23]. It then uses these rates to determine the most probable alignments. LAST also finds

108 the most probable division of each read into one or more parts together with the most

109 probable alignment of each part with the last-split function [24], [25], i.e. if LAST finds a

110 better scoring alternative alignment of the read where it is being split into parts and

111 mapped to different regions of the target, such alignment is submitted to the results. This

112 is a principled way to handle complex DNA rearrangements, gene fusions in RNA,

113 chimeric host/viral sequences, etc. The tasks of detecting polymorphisms and

114 distinguishing viral serotypes can be performed more precisely for MinION™ data when

115 based on a LAST alignment, because the tool estimates the probability that each base is

116 correctly aligned. The probability is low if there is an ambiguity, i.e. the base could align to

117 more than one place. The default NanoPipe parameters for the LAST alignment are

118 efficient for most cases, but can be easily adjusted by the user. We use last-train to find

119 the optimal alignment parameters for MinION™ sequencing [23]. For advanced users, it is

120 recommended, although not necessary, to acquaint with the settings of LAST and last-

121 train [26].

122 After the alignment is completed, NanoPipe evaluates the nucleotide variation for each

123 position, and based on this analysis generates consensus sequence and a list of possible

124 single nucleotide polymorphisms (SNP). The minimum nucleotide count per position

125 should be at least ten (i.e., at least ten reads), otherwise a gap will be assigned at the

126 position in the consensus sequence. The consensus sequence is calculated based on the

127 majority rule, i.e. the nucleotide with the higher count at a particular position is assigned

128 to the consensus. If the counts for any two nucleotides differ from each other by not more

129 than by 20%, both nucleotides are included in the consensus with the use of the IUPAC

130 nomenclature [27]. Statistical evaluation of nucleotide variation is presented in a separate

131 table, and suggests a polymorphism candidate if an alternate nucleotide has coverage of

132 at least 20% of all reads.

133 Figure 1. A schematic representation of the NanoPipe workflow.

134

135 To distinguish between artifacts and true polymorphisms, we have set three additional

136 filters. First, SNP candidate should have read coverage of at least 30% from the maximum

137 coverage in the respective contig, otherwise it is not listed in the polymorphisms table

138 and not counted as a SNP. Second, the probability of SNP at the position is calculated. For

139 that, the relative nucleotide frequency is multiplied by a custom weight factor:

140 transversions are weighted by the factor of 1 and transitions – by the factor of 2. This is

141 based on the assumption that transitions are two times more likely than transversions.

142 This weighted probability for each SNP candidate is rescaled for convenience, so that the

143 maximum value is 1. The probability is displayed for each nucleotide in the

144 polymorphisms table, a SNP candidate is most likely to be the true SNP if its joint

6

145 probability is 1. Third, the analysis is refined by an assessment of the alignment quality

146 around a potential SNP by estimating a p-error. The p-error is calculated based on the

147 formula used in the LAST methods to calculate the probability of the alignment for a

148 single nucleotide [28]. In NanoPipe the p-error is estimated over a region of maximum 10

149 nucleotides before and after the SNP position (excluding the SNP itself) and for all the

150 read alignments at the region (Fig. 2). It is based on the LAST reliability score assigned to

151 each base pair of the alignment. In addition, polymorphism candidates for human and

152 *Plasmodium* are linked to the public SNP databases.

$$p - error = \frac{\sum_1^i 10^{-\left(\frac{ASCII_{value}-33}{10}\right)}}{i}$$

153

154 Figure 2. Formula used in NanoPipe for the p-error calculations, where $i$ is the total

155 number of nucleotides around the SNP for all mapped reads, i.e. $i$ = N(nucleotides around

156 SNP)*N(reads mapped to the evaluated region). ASCII values are extracted from the LAST

157 alignment and indicate the reliability of each base's alignment, see more information

158 about LAST quality symbols at [28].

159 The results are supplemented with a number of useful pages, such as alignment of

160 consensus sequences against target sequences, read length distribution and individual

161 reads' alignment length distribution, and nucleotide plots showing distribution of

162 nucleotides from all reads at each position of the consensus sequence. The latter is the

163 interactive visual representation of the results and enables the user to monitor any

164 nucleotide variations by eye. NanoPipe also provides the necessary files for browsing the

165 alignment in the IGV genome browser [29], where each individual read alignment is

166 displayed.

167 **Findings**

7

168 The experience of our working group in conducting MinION™ sequencing and analysis

169 workshops [30] for medical doctors and wet-lab researchers revealed the general gaps in

170 bioinformatics knowledge among these groups. While the experimental part of the

171 workshop was easily performed by the participants, even if mastered for the first time,

172 the processing of the sequenced data proved to be the most enduring and difficult part.

173 Therefore, keeping in mind the weakest spots of non-bioinformatician researchers, we

174 will describe here some features of NanoPipe usage in details.

175 To start a new analysis, a query (sequencing reads) and a target (a genome or a region of

176 interest) should be provided. The query (1D or 1D$^2$) can be uploaded via the NanoPipe's

177 web interface from the user's computer in one FASTA or FASTQ file, also in archived

178 format. As the ONT-provided base callers output sequenced reads in multiple files, it is

179 important to know how to merge them into one file, in an operating system of the user

180 (Windows, Mac OS, Linux). As an alternative, NanoPipe can also handle multiple FASTA or

181 FASTQ files when they are archived (zipped), for example, with "zip" command on Linux

182 or WinZip tool on Windows. The maximum size of the query file should not exceed 3 GB.

183 The target can be chosen from a drop-down menu or uploaded by the user in FASTA

184 format. The list of NanoPipe precompiled targets includes reference genomes for human

185 (hg38), RefSeq accession GCF_000001405.38; *Escherichia coli*, RefSeq accession

186 GCF_000005845.2; *Caenorhabditis elegans*, RefSeq accession GCF_000002985.6;

187 *Droshophila melanogaster*, RefSeq accession GCF_000001215.4; *Mus musculus*, RefSeq

188 accession GCF_000001635.26; *Arabidopsis thaliana*, RefSeq accession GCF_000001735.4;

189 *Plasmodium falciparum* strain 3D7, downloaded from plasmodb.org, version=2013-03-01;

190 a representative genome for *Camponotus floridanus*, RefSeq accession GCF_003227725.1;

191 and Dengue virus genome variants for serotyping (NC_001477.1, NC_001474.2,

192 NC_001475.2 and NC_002640.1 for variant 1, variant 2, variant 3 and variant 4,

193 respectively). Otherwise, DNA sequences of whole genomes, transcriptomes or genes for

194 different organisms can be accessed via particular species databases or in the

195 corresponding databases at NCBI [31]. The NanoPipe prepared targets include the

196 precompiled, best fitting alignment parameters and substitution matrices (based on last-

197 train calculations). The substitution matrix is an important part of any alignment and

198 contains information about the mismatch cost between any pair of nucleotides. Any

199 uploaded targets will be used with the NanoPipe default parameters that are suitable for

200 most cases. To avoid "noise", i.e. mapping of too short reads, the user can set the read

201 length limit, for example to 200 nucleotides (depends on the experiment type and

202 purpose). A unique name can be assigned to a job and used later within one month to

203 retrieve the results from the NanoPipe server. After one month the data will be deleted

204 from the server.

205 Depending on the query and target size, a NanoPipe analysis can last from several

206 minutes to several hours, the server's memory used for the calculations is limited to 16

207 GB. The user will receive the notification via email (if the email address was provided)

208 when his/her job is completed. The summary of the completed analysis depicts the LAST

209 parameters that were used and the mapping statistics, i.e. how many reads were mapped

210 altogether and the reads distribution per chromosomes/scaffolds of the target. This table

211 can be sorted in increasing/decreasing order. Mapping distribution statistics show how

212 long the reads in the query were, and allow estimating whether the sequencing resulted

213 in the expected read lengths. Thus, for whole genome sequencing by MinION™ a

214 researcher can expect read lengths up to over a million nucleotides; targeted sequencing

215 results in read lengths corresponding to the target length; transcriptome sequencing or

216 RNAseq experiments should provide reads with length typical to particular species'

217 transcript lengths (around one-to-two thousand nucleotides for most organisms).

218 Alignment distribution statistics inform about the quality of the alignment, i.e. whether

219 the whole read or a part of it could be mapped to the target.

220 Alignment results are visualized in NanoPipe in several ways: ordinal line by line pairwise

221 alignment, a BAM file and a graphical representation via nucleotide plots. BAM and

222 indexed BAM (BAI) files can be easily downloaded from the results page and further used

223 for an interactive genome browsers, for example, IGV. The target FASTA files are required

224 by IGV, as well, and can be downloaded together with BAM files. IGV is free software and

225 can be accessed from [32]. Nucleotide plots represent the colored mapping scheme at

226 each position of the sequence, each nucleotide marked with its specific color. This enables

227 easy monitoring for gaps and possible nucleotide substitutions. Navigating along the plot

228 is enabled via right and left shift, as well by entering a nucleotide's coordinate in the

229 search field. For long target regions, nucleotide plots provide a zoom-out preview at the

230 bottom of the page.

231 Each chromosome/region of the target is supported with an individual consensus which

232 can be seen and downloaded in FASTA format. Positions that cannot be defined (not

233 enough information in the input data) are designated as "N", gaps are designated as "-"; if

234 a position is occupied by any of the four nucleotides it is designated as "X"; other

235 ambiguous positions are designated using IUPAC nomenclature.

236 The polymorphisms table lists SNP candidates and provides joint probabilities for each

237 candidate (maximum = 1), as well raw counts for each nucleotide at the target position. If

238 the data are available, the corresponding SNP IDs (identifiers) will be retrieved from the

239 existing databases (currently available only for human [33] and *Plasmodium* [34]). The

240 alignment quality around an SNP candidate is reported as p-error, the higher the p-error

241 (maximum 1), the lower the alignment quality. Low alignment quality might indicate a

242 region of sequencing errors around an individual SNP and, thus, signify the lower

243 reliability of the candidate detection. However, a cluster of closely located SNPs within a

244 distance of less than 10 nucleotides would have a similar effect. Therefore, the p-error is

245 an additional parameter that might be taken into consideration by the user. Detailed

246 analysis of the nucleotide plot or the alignment in the IGV viewer around questionable

247 candidates may help in making a decision, including consideration of biological relevance.

248 To view the result pages for nucleotide plots, consensus, polymorphisms and alignment

249 the user needs to choose a particular chromosome/region. This approach accelerates the

250 data display and prevents the web browser from overloading.

251

## Study cases

### *Plasmodium* polymorphism detection

254 The targeted regions were first amplified using the standard PCR protocol [35]. The

255 resulting amplicons were sequenced with the MinION™ using the ONT sequencing kit

256 SDK-MAP005 or SQK-MAP006 for the library preparation and the flow cell version FLO-

257 MAP003. The sequences were aligned against the *Plasmodium* genome (*P. falciparum*

258 reference genome based on 3D7 strain). Exhaustive discussion of the ONT utility for

259 *Plasmodium* SNP calling and the library preparation methods are presented elsewhere

260 [35]. NanoPipe succeeded in mapping 99.9% of all query reads and detecting the expected

261 mutations (see Fig. 3). The specific characteristic of the *Plasmodium* genome, multiple

262 AT/T repeats, can be easily observed in the nucleotide plot. The full analysis of a 21 MB

263 (FASTA format) query against the *P. falciparum* genome with NanoPipe took less than

264 four minutes. This example demonstrates that the analysis of just 10157 MinION™ reads

265 with a high AT content on NanoPipe results in reliable data. Detailed screen shots of the

11

266  NanoPipe results for this study case can be found in the Supplementary material (Figures

267  S2 – S10).

268

269  Figure 3. Sample case 1. A: The polymorphisms table displays the three SNP candidates:

270  two of them (at positions 403625 and 404407) are expected mutations leading to K76T

271  and A220S amino acid changes and, as a consequence, to altered resistance of the parasite

272  to chloroquine, mefloquine and quinine drugs [35]. B: Nucleotide plot. The purple arrow

273  points to the G>T substitution at the position 404407 (GCC>TCC codon change). The

274  orange arrow highlights an AT-rich region, the typical feature of the *P. falciparum*

275  genome.

276

**277  Targeted sequencing of EGFR transcript from human lung adenocarcinoma cell line**

**278  H1975**

279  The region of the human EGFR cDNA corresponding to exons 17-22 was amplified using

280  the primers CTAAGATCCCGTCCATCGCC (forward) and ACATATGGGTGGCTGAGGGA

281  (reverse). The library preparation was performed following the manufacturer's

282  recommendation using the SQK-LSK108 kit from ONT and then sequencing was done with

283  the MinION™ using flow cell version FLO-MIN107 R9. The raw 1D reads (900 MB, FASTA

284  format) were uploaded to NanoPipe using the human reference genome as a target; the

285  analysis took 135 minutes. 79.9% of all reads were mapped to the EGFR gene on

286  chromosome 7. Most of the remaining reads (13.5%) were mapped to chromosome 11, to

287  regions corresponding to the cortactin gene (NM_138565) that is overexpressed in

288  different cancers [36], and the gene encoding subunit 2 of the splicing factor 3b protein

289  complex (NM_006842). This gene might be differentially expressed in tumor tissues [37],

290  nevertheless, there is no definite research about possible roles of this gene in cancers. It

12

291 may be the case that the gene is overexpressed in cell line H1975 and, therefore, its

292 transcript was sequenced as an abundant contaminant, similar to the cortactin transcript.

293 The four mutations expected within this region in this cell line [38] were detected, see Fig.

294 4. This example demonstrates the suitability of NanoPipe for cancer sequencing analysis.

295

296 Figure 4. A. Polymorphisms results. The two expected nucleotide substitutions are silent:

297 G>A (CAG>CAA = Gln) at position 55181370 and T>C (ACT>ACC = Thr) at position

298 55198724; two other substitutions at position 55181378 (C>T leading to the amino acid

299 change T745M) and at position 55191822 (T>G leading to the amino acid change L813R)

300 are responsible for the sensitivity to anticancer drugs, in particular to gefitinib and

301 erlotinib [38]. B. Nucleotide plot. The zoom-out scheme at the bottom depicts seven

302 alignment picks, they represent seven sequenced exons of the transcript. The nucleotide

303 plot is showing the first of these picks, pointed to by the orange arrows.

304

**RNAseq analysis of the ant species *Camponotus maculatus***

305

306 Monarch Total RNA Miniprep kit (NEB) was used for the poly-A RNA extraction according

307 to the manufacturer's recommendations. The library was prepared using the SQK-PCS108

308 kit from ONT and then sequenced with MinION™ using FLO-MIN107 R9 flow cell version

309 following the manufacturer's recommendations. The raw 1D reads in FASTA format were

310 uploaded to NanoPipe, the reference genome of *Camponotus floridanus* was used as the

311 target genome, as neither genome nor transcriptome of *C. maculatus* are available yet. The

312 analysis of the 1.5 GB query ran for 3 hours. Out of 1814750 raw reads 1773747 (97.7 %)

313 were mapped to the target, spanning 431 scaffolds out of 657, including 150 scaffolds

314 with coverage of more than 1000 reads per scaffold. This result is consistent with the

315  number of reads sequenced, and the fact that the target was the genome of a different

316  species, thus, NanoPipe can be used for studying newly sequenced species.

317

**Direct RNAseq of poxvirus isolated from the host cells of green monkey**

319  We examined the recent MinION™ direct RNA sequencing data of the *Vaccinia virus WR*

320  mRNA isolated from the kidney fibroblast cells of *Chlorocebus sabeus* (for details see

321  [39]). Using the virus genome and the monkey genome as targets in the two separate

322  runs, we could separate the reads coming from the two organisms. From the tested 29846

323  reads, 1314 were mapped to the virus genome (GenBank accession LT966077.1) and

324  14714 reads were mapped to the *C. sabeus* genome (GenBank accession

325  GCA_000409795.2), which is consistent with the published results. The option of the

326  direct target upload was also tested on the small and big size input files (the virus genome

327  size is 198 KB and the green monkey genome is of 2.82 GB). The both runs were finished

328  successfully in 12 seconds and in 51 minutes, respectively, including the construction of

329  the target databases. NanoPipe proved to be useful in analyzing the mixed long reads of

330  virus and its host.

331

**Discussion**

333  We have developed NanoPipe, a web-driven application that enables easy analysis of

334  MinION™ sequencing data and is suitable for, both, experienced bioinformaticians and

335  biologists with limited IT knowledge. NanoPipe provides users with the consensus

336  sequence of the studied DNA and a list of putative single nucleotide polymorphisms. In

337  this work, we used four different experimental datasets: targeted sequencing of the AT-

338  rich genome of *P. falciparum*, targeted sequencing of EGFR human cDNA, direct RNA

339  sequencing of *V. virus* and its host, green monkey *C. sabeus*; and RNAseq of the ant  *P.*

14

*maculatus* that has so far no fully sequenced genome nor transcriptome. These four test cases represent different tasks that biologists from different fields of study could be interested in: sequencing and SNP detection in a repeat-rich genomes with low sequencing yield (study case with *P. falciparum*); mutation detection in human cancer samples (study case with lung adenocarcinoma cell line H1975); RNAseq of a newly sequenced organism (study case with *C. maculatus*) and of a mixed species sample (virus-host study). NanoPipe proved to be a reliable tool for all these tasks: It detected all expected mutations for *Plasmodium* and the EGFR transcript region and succeeded in mapping more than 90% RNAseq reads from *C. maculatus* to the reference genome. In the case with the virus-host direct RNA sequencing NanoPipe successfully segregated the mixed reads between the virus genome and its host genome. We have tested NanoPipe during two workshops at the International Summer School for MinION™ sequencing in Bangkok (2017) and Manado (2018) with participants having no bioinformatics experience from the medical and biological research fields. NanoPipe proved to be efficient and understandable for these users. Nevertheless, our team also noted during these practical studies that the participants lacked some basic knowledge, necessary for conducting successful sequencing experiments. Therefore, we described here in detail the general scheme of working with long read sequencing data. The flexibility of NanoPipe allows researchers to study any sequence of interest (including cancer samples in human patients), as the list of targets is not limited. With this software, we hope to make MinION™ sequencing even more accessible for medical researchers and biologists without sophisticated IT resources and expertise. The next step for the NanoPipe project will include more sophisticated statistics for SNP detection and evaluation. Also, we are planning to cover the microbiology field and implement an option for metagenomics analysis.

**Availability of data and source code**

Project name: NanoPipe

Project home page:  http://bioinformatics.uni-muenster.de/tools/nanopipe, github page:

https://github.com/IOB-Muenster/nanopipe2. The NanoPipe package for local

installation is available at the NanoPipe github page. The explanation on installation and

check-install procedures can be found in the github directory.

Operating system: Unix

Programming language: Javascript, Python, Perl

License: Apache License 2.0

Other: there are no limitations for web browser type. The local version can be installed on

Unix operating system

RRID: SCR_016852

The raw sequence data for the test runs and test jobs can be accessed from the homepage

http://bioinformatics.uni-muenster.de/share/NanoPipe_test_data/

Test data description:

1. Data set *P. falciparum* sequencing

   1.1 MinION ™ targeted sequencing for *Plasmodium falciparum* sequencing.

   1.2 PfCRT and K-propeller genes were sequenced using long-read technologies within

the experiments of testing MinION ™ sequencer for SNP detection in *Plasmodium*.

   1.3 The sequencing was performed for the MinION™ sequencing workshop, 2017 in

Bangkok, Thailand at the Mahidol University.

   1.4 Raw reads in a single FASTA file, archived using zgip command.

   1.5 *P. falciparum* culture, strain 3D7.

   1.6 Size: 21 MB, archived: 4.3 MB.

390

391 2. Data set for EGFR sequencing

392    2.1 MinION targeted sequencing of the CDS for EGFR human gene, spanning exons 17-

393 22.

394    2.2 These data were generated to test MinION capacity for detecting SNP in human

395 cancer DNA

396    2.3 The sequencing was performed for the MinION™ sequencing workshop, 2017 in

397 Bangkok, Thailand at the Mahidol University.

398    2.4 Raw reads in a single FASTA file, archived using zgip command.

399    2.5 Human lung adenocarcinoma cell line H1975 (RRID:CVCL_1511)

400    2.6 Size: 900 MB, archived: 259 MB.

401

402 3. Data set for *Camponotus maculatus* sequencing

403    3.1 MinION RNAseq for ant *C. maculatus*

404    3.2 Poly-A RNA sequencing was performed for the given ant species at first time

405    3.3 The sequencing was performed in the Institute of Bioinformatics, Muenster,

406 Germany; the library was prepared in the group of Juergen Gadau, University of Muenster.

407    3.4 Raw reads in a single FASTA file, archived using zgip command.

408    3.5 Larvae and adult individuals from the lab culture

409    3.6 Size: 1.5 GB, archived: 309 MB.

410

411 Supporting data and snapshots of the GitHub are also archived in the *GigaScience* GigaDB

412 repository[40].

413

414 **Supplementary materials.**

17

415  Figure S1: Number of scientific publications that contain the phrase "Oxford Nanopore" in

416  their abstracts.

417  Figures S2-S10: The results pages for the *P. falciparum* test case.

418

**Declarations**

**List of abbreviations**

421  CDS: coding sequence; GB: GigaByte; ID: Identifier; IT: Information Technology; MB: Mega

422  Byte; ONT: Oxford Nanopore Technologies; SNP: Single Nucleotide Polymorphism

423

**Ethics approval and consent to participate**

425  Not applicable

426

**Consent for publication**

428  Not applicable

429

**Competing interests**

431  The authors declare that they have no competing interests

432

**Funding**

434  The work has been funded by Institute of Bioinformatics, University of Muenster and

435  University Clinic Muenster (UKM)

436

**Authors' contribution**

438  VS validated the software, participated in the concept design and composed the

439  manuscript draft; TK designed the software and participated in the code writing; FM

wrote the code for polymorphism analysis and contributed to the manuscript writing; NG wrote the major part of the pipeline code and designed the software; MF participated in the improvement of the software, including LAST optimization, and contributed to the manuscript writing; YS contributed to the analysis and provided the access to the MinION™ sequencing summer school; WM designed the concept and supervised the work. All the authors read and approved the manuscript.

## References

[1]     ONT, "Flongle.". https://nanoporetech.com/products/flongle. [Accessed: 05-Sep-2018].

[2]     T. A. Sasani, K. R. Cone, A. R. Quinlan, and N. C. Elde, "Long read sequencing reveals poxvirus evolution through rapid homogenization of gene arrays," *bioRxiv*, p. 245373, Jan. 2018.

[3]     T. Karamitros, I. Harrison, R. Piorkowska, A. Katzourakis, G. Magiorkinis, and J. L. Mbisa, "De Novo Assembly of Human Herpes Virus Type 1 (HHV-1) Genome, Mining of Non-Canonical Structures and Detection of Novel Drug-Resistance Mutations Using Short- and Long-Read Next Generation Sequencing Technologies," *PLoS One*, vol. 11, no. 6, p. e0157600, Jun. 2016.

465 [4] A. Bronzato Badial, D. Sherman, A. Stone, A. Gopakumar, V. Wilson, W. Schneider,

466 and J. King, "Nanopore Sequencing as a Surveillance Tool for Plant Pathogens in

467 Plant and Insect Tissues," *Plant Dis.*, vol. 102, no. 8, pp. 1648–1652, Aug. 2018.

468 [5] M. B. Fleming, E. L. Patterson, P. A. Reeves, C. M. Richards, T. A. Gaines, and C.

469 Walters, "Exploring the fate of mRNA in aging seeds: protection, destruction, or

470 slow decay?," *J. Exp. Bot.*, vol. 69, no. 18, pp. 4309–4321, Aug. 2018.

471 [6] A. Bainomugisa, T. Duarte, E. Lavu, S. Pandey, C. Coulter, B. J. Marais, and L. M.

472 Coin, "A complete high-quality MinION nanopore assembly of an extensively drug-

473 resistant Mycobacterium tuberculosis Beijing lineage strain identifies novel

474 variation in repetitive PE/PPE gene regions," *Microb. Genomics*, Jun. 2018.

475 [7] W. Pearman, A. N. H. Smith, G. Breckell, J. Dale, N. E. Freed, and O. K. Silander, "New

476 tools for diet analyses: nanopore sequencing of metagenomic DNA from stomach

477 contents to quantify diet in an invasive population of rats," *bioRxiv*, p. 363622, Jul.

478 2018.

479 [8] S. Theuns, B. Vanmechelen, Q. Bernaert, W. Deboutte, M. Vandenhole, L. Beller, J.

480 Matthijnssens, P. Maes, and H. J. Nauwynck, "Nanopore sequencing as a

481 revolutionary diagnostic tool for porcine viral enteric disease complexes identifies

482 porcine kobuvirus as an important enteric virus," *Sci. Rep.*, vol. 8, no. 1, p. 9830,

483 Dec. 2018.

484 [9] H. Krehenwinkel, A. Pomerantz, J. B. Henderson, S. R. Kennedy, J. Y. Lim, V. Swamy,

485 J. D. Shoobridge, N. H. Patel, R. G. Gillespie, and S. Prost, "Nanopore sequencing of

486 long ribosomal DNA amplicons enables portable and simple biodiversity

487 assessments with high phylogenetic resolution across broad taxonomic scale,"

488 *GigaScience*. 2019. In press.

489 [10] J. Quick, N. J. Loman, S. Duraffour, J. T. Simpson, E. Severi, L. Cowley, J. A. Bore, R.

490 Koundouno, G. Dudas, A. Mikhail, N. Ouédraogo, B. Afrough, A. Bah, J. H. J. Baum, B.

491 Becker-Ziaja, J. P. Boettcher, M. Cabeza-Cabrerizo, Á. Camino-Sánchez, L. L. Carter,

492 J. Doerrbecker, T. Enkirch, I. G.- Dorival, N. Hetzelt, J. Hinzmann, T. Holm, L. E.

493 Kafetzopoulou, M. Koropogui, A. Kosgey, E. Kuisma, C. H. Logue, A. Mazzarelli, S.

494 Meisel, M. Mertens, J. Michel, D. Ngabo, K. Nitzsche, E. Pallasch, L. V. Patrono, J.

495 Portmann, J. G. Repits, N. Y. Rickett, A. Sachse, K. Singethan, I. Vitoriano, R. L.

496 Yemanaberhan, E. G. Zekeng, T. Racine, A. Bello, A. A. Sall, O. Faye, O. Faye, N.

497 Magassouba, C. V. Williams, V. Amburgey, L. Winona, E. Davis, J. Gerlach, F.

498 Washington, V. Monteil, M. Jourdain, M. Bererd, A. Camara, H. Somlare, A. Camara,

499 M. Gerard, G. Bado, B. Baillet, D. Delaune, K. Y. Nebie, A. Diarra, Y. Savane, R. B.

500 Pallawo, G. J. Gutierrez, N. Milhano, I. Roger, C. J. Williams, F. Yattara, K.

501 Lewandowski, J. Taylor, P. Rachwal, D. J. Turner, G. Pollakis, J. A. Hiscox, D. A.

502 Matthews, M. K. O. Shea, A. M. Johnston, D. Wilson, E. Hutley, E. Smit, A. Di Caro, R.

503 Wölfel, K. Stoecker, E. Fleischmann, M. Gabriel, S. A. Weller, L. Koivogui, B. Diallo,

504 S. Keïta, A. Rambaut, P. Formenty, S. Günther, and M. W. Carroll, "Real-time,

505 portable genome sequencing for Ebola surveillance," *Nature*, vol. 530, no. 7589,

506 pp. 228–232, Feb. 2016.

507 [11] P. Mbala Kingebeni, C.-J. Villabona-Arenas, N. Vidal, J. Likofata, J. Nsio-Mbeta, S.

508 Makiala-Mandanda, D. Mukadi, P. Mukadi, C. Kumakamba, B. Djokolo, A. Ayouba, E.

509 Delaporte, M. Peeters, J.-J. Muyembe Tamfum, and S. Ahuka Mundeke, "Rapid

510 confirmation of the Zaire Ebola Virus in the outbreak of the Equateur province in

511 the Democratic Republic of Congo: implications for public health interventions,"

512 *Clin. Infect. Dis.*, Jun. 2018.

513 [12] S. L. Castro-Wallace, C. Y. Chiu, K. K. John, S. E. Stahl, K. H. Rubins, A. B. R. McIntyre,

514 J. P. Dworkin, M. L. Lupisella, D. J. Smith, D. J. Botkin, T. A. Stephenson, S. Juul, D. J.

515     Turner, F. Izquierdo, S. Federman, D. Stryke, S. Somasekar, N. Alexander, G. Yu, C.

516     E. Mason, and A. S. Burton, "Nanopore DNA Sequencing and Genome Assembly on

517     the International Space Station," *Sci. Rep.*, vol. 7, no. 1, p. 18022, Dec. 2017.

518   [13]  E. A. Solares, M. Chakraborty, D. E. Miller, S. Kalsow, K. Hall, A. G. Perera, J. J.

519     Emerson, and R. S. Hawley, "Rapid Low-Cost Assembly of the Drosophila

520     melanogaster Reference Genome Using Low-Coverage, Long-Read Sequencing.,"

521     *G3 (Bethesda).*, p. g3.200162.2018, Jul. 2018.

522   [14]  S. Deschamps, Y. Zhang, V. Llaca, L. Ye, G. May, and H. Lin, "A chromosome-scale

523     assembly of the sorghum genome using nanopore sequencing and optical

524     mapping," *Nat Commun.* 2018 Nov 19;9(1):4844. doi: 10.1038/s41467-018-

525     07271-1.

526   [15]  B. Pradhan, T. Cajuso, R. Katainen, P. Sulo, T. Tanskanen, O. Kilpivaara, E. Pitkänen,

527     L. A. Aaltonen, L. Kauppi, and K. Palin, "Detection of subclonal L1 transductions in

528     colorectal cancer by long-distance inverse-PCR and Nanopore sequencing," *Sci.*

529     *Rep.*, vol. 7, no. 1, p. 14521, Dec. 2017.

530   [16]  L. F. K. Kuderna, E. Lizano, E. Julia, J. Gomez-Garrido, A. Serres-Armero, M.

531     Kuhlwilm, R. A. Alandes, M. Alvarez-Estape, T. Alioto, M. Gut, I. Gut, M. H. Schierup,

532     O. Fornas, and T. Marques-Bonet, "Selective single molecule sequencing and

533     assembly of a human Y chromosome of African origin," *bioRxiv*, p. 342667, Jun.

534     2018.

535   [17]  K. Chacon-Vargas, A. A. Chirino, M. M. Davis, S. A. Debler, W. R. Haimer, J. J. Wilbur,

536     X. Mo, B. W. Worthing, E. G. Wainblat, S. Zhao, and J. G. Gibbons, "Genome

537     Sequence of *Zymomonas mobilis* subsp. *mobilis* NRRL B-1960," *Genome Announc.*,

538     vol. 5, no. 30, Jul. 2017.

539   [18]  A. C. Reis, K. Kroll, M. Gomila, B. A. Kolvenbach, P. F. X. Corvini, and O. C. Nunes,

22

540       "Complete Genome Sequence of *Achromobacter denitrificans* PR1," *Genome*

541       *Announc.*, vol. 5, no. 31, Aug. 2017.

542   [19]   Metrichor, "METRICHOR LTD." https://metrichor.com/. [Accessed: 05-Sep-2018].

543   [20]   N. J. Loman, J. Quick, and J. T. Simpson, "A complete bacterial genome assembled

544       de novo using only nanopore sequencing data," *Nat. Methods*, vol. 12, no. 8, pp.

545       733–735, Aug. 2015.

546   [21]   A. Srivathsan, B. Baloğlu, W. Wang, W. X. Tan, D. Bertrand, A. H. Q. Ng, E. J. H. Boey,

547       J. J. Y. Koh, N. Nagarajan, and R. Meier, "A MinION-based pipeline for fast and cost-

548       effective DNA barcoding," *Mol Ecol Resour.* 2018 Apr 19. doi: 10.1111/1755-

549       0998.12890.

550   [22]   R. M. Leggett, D. Heavens, M. Caccamo, M. D. Clark, and R. P. Davey, "NanoOK:

551       multi-reference alignment analysis of nanopore sequencing data, quality and

552       error profiles," *Bioinformatics*, vol. 32, no. 1, p. btv540, Sep. 2015.

553   [23]   M. Hamada, Y. Ono, K. Asai, M. C. Frith, and J. Hancock, "Training alignment

554       parameters for arbitrary sequencers with LAST-TRAIN," *Bioinformatics*, vol. 33,

555       no. 6, p. btw742, Dec. 2016.

556   [24]   M. C. Frith and R. Kawaguchi, "Split-alignment of genomes finds orthologies more

557       accurately," *Genome Biol.*, vol. 16, no. 1, p. 106, Dec. 2015.

558   [25]   M. C. Frith and S. Khan, "A survey of localized sequence rearrangements in human

559       DNA," *Nucleic Acids Res.*, vol. 46, no. 4, pp. 1661–1673, Feb. 2018.

560   [26]   Frith M., "LAST: genome-scale sequence comparison." http://last.cbrc.jp/.

561       [Accessed: 05-Sep-2018].

562   [27]   IUPAC, "IUPAC Codes." https://www.bioinformatics.org/sms/iupac.html.

563       [Accessed: 05-Sep-2018].

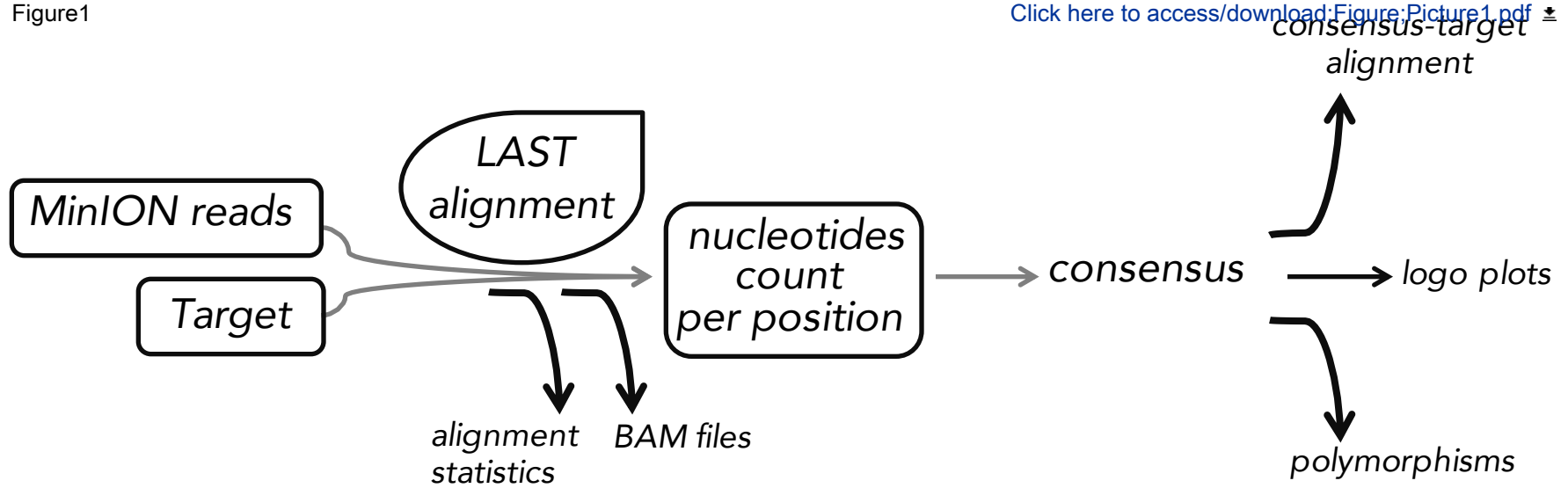564   [28]   Frith M., "last-split." http://last.cbrc.jp/doc/last-split.html. [Accessed: 05-Sep-

565    2018].

566    [29]    H. Thorvaldsdottir, J. T. Robinson, and J. P. Mesirov, "Integrative Genomics Viewer

567            (IGV): high-performance genomics data visualization and exploration," *Brief.*

568            *Bioinform.*, vol. 14, no. 2, pp. 178–192, Mar. 2013.

569    [30]    Y. Suzuki, "JSPS, Training Seminar for MinION." http://www.cb.k.u-

570            tokyo.ac.jp/suzukilab/minion/summary.html. [Accessed: 09-Oct-2018].

571    [31]    NCBI, "National Center for Biotechnology Information."

572            https://www.ncbi.nlm.nih.gov/. [Accessed: 05-Sep-2018].

573    [32]    IGV-page, "Home | Integrative Genomics Viewer."

574            http://software.broadinstitute.org/software/igv/. [Accessed: 05-Sep-2018].

575    [33]    NCBI.dbSNP, "dbSNP." https://www.ncbi.nlm.nih.gov/projects/SNP/. [Accessed:

576            10-Oct-2018].

577    [34]    PlasmoDB, "PlasmoDB : The Plasmodium Genomics Resource."

578            http://plasmodb.org/plasmo/. [Accessed: 10-Oct-2018].

579    [35]    L. R. Runtuwene, J. S. B. Tuda, A. E. Mongan, W. Makalowski, M. C. Frith, M.

580            Imwong, S. Srisutham, L. A. Nguyen Thi, N. N. Tuan, Y. Eshita, R. Maeda, J.

581            Yamagishi, and Y. Suzuki, "Nanopore sequencing of drug-resistance-associated

582            genes in malaria parasites, Plasmodium falciparum," *Sci. Rep.*, vol. 8, no. 1, p. 8286,

583            Dec. 2018.

584    [36]    C. Gattazzo, V. Martini, F. Frezzato, V. Trimarco, E. Tibaldi, M. Castelli, M. Facco, F.

585            Zonta, A. M. Brunati, R. Zambello, G. Semenzato, and L. Trentin, "Cortactin, another

586            player in the Lyn signaling pathway, is over-expressed and alternatively spliced in

587            leukemic cells from patients with B-cell chronic lymphocytic leukemia.,"

588            *Haematologica*, vol. 99, no. 6, pp. 1069–77, Jun. 2014.

589    [37]    T. A. Lantto, I. Laakso, H. J. D. Dorman, T. Mauriala, R. Hiltunen, S. Kõks, and A.

590      Raasmaja, "Cellular Stress and p53-Associated Apoptosis by Juniperus communis

591      L. Berry Extract Treatment in the Human SH-SY5Y Neuroblastoma Cells.," *Int. J.*

592      *Mol. Sci.*, vol. 17, no. 7, Jul. 2016.

593  [38]  A. Suzuki, M. Suzuki, J. Mizushima-Sugano, M. C. Frith, W. Makalowski, T. Kohno, S.

594      Sugano, K. Tsuchihara, and Y. Suzuki, "Sequencing and phasing cancer mutations

595      in lung cancers using a long-read portable sequencer.," *DNA Res.*, vol. 24, no. 6, pp.

596      585–596, Dec. 2017.

597  [39]  D. Tombácz, I. Prazsák, A. Szűcs, B. Dénes, M. Snyder, and Z. Boldogkői, "Dynamic

598      Transcriptome Profiling Dataset of Vaccinia Virus Obtained from Long-read

599      Sequencing Techniques," *Gigascience*, Nov. 2018. doi: 10.1093/gigascience/giy139

600  [40]  Shabardina V; Kischka T; Manske F; Grundmann N; Frith MC; Suzuki Y;

601      Makalowski W (2019): Supporting data for "NanoPipe - a web server for nanopore

602      MinION sequencing data analysis" GigaScience Database.

603      http://dx.doi.org/10.5524/100551

604

Figure1

Figure3

Figure 3

Figure4

A

| Position | A | C | G | T | Target | Matches in dbSNP | P-error (local alignment quality) | raw A | raw C | raw G | raw T |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 55173958 | - | - | 1.0 | - | a | | 0.0027 | 104623 | 1304 | 38751 | 753 |
| 55181370 | 1.0 | - | - | - | g | rs1050171: G/A+G/C | 0.0008 | 114375 | 5449 | 75568 | 3831 |
| 55181378 | - | - | - | 1.0 | c | rs121434569: C/T | 0.0007 | 1423 | 65474 | 2980 | 128968 |
| 55191822 | - | - | 1.0 | - | t | rs121434568: T/A+T/G | 0.0013 | 19653 | 10419 | 75524 | 59396 |
| 55192839 | - | - | 1.0 | - | a | rs376176117: A/T | 0.0695 | 70026 | 211 | 18566 | 37 |
| 55198724 | - | 1.0 | - | - | t | rs1140475: T/A+T/C | 0.0563 | 1723 | 72518 | 896 | 49645 |

B

Click here to access/download
**Supplementary Material**
NanoPipe_Supplementary.pdf